

Running Head: VISUAL SPEECH, DETECTION, AND ATTENTION

1

Developmental Shifts in Detection and Attention

for Auditory, Visual, and Audiovisual SpeechSusan Jerger, Ph.D.^{a,b}Markus F. Damian, Ph.D.^cCassandra Karl, M.S. (pending),^{a,b}Hervé Abdi, Ph.D.^a^aSchool of Behavioral and Brain Sciences, GR4.1, University of Texas at Dallas, 800 W. Campbell Rd,

Richardson, TX 75080,

S. Jerger: sjerger@utdallas.edu, H. Abdi: herve@utdallas.edu, C. Karl: cnkarl07@gmail.com^bCallier Center for Communication Disorders, 811 Synergy Park Blvd., Richardson, TX 75080^cUniversity of Bristol, School of Experimental Psychology, 12a Priory Road, Room 1D20, Bristol BS8 1TU,United Kingdom, m.damian@bristol.ac.ukCorresponding Author: Susan Jerger, School of Behavioral and Brain Sciences, GR4.1, University of Texas at Dallas, 800 W. Campbell Rd, Richardson, TX 75080, sjerger@utdallas.edu, Phone: 512-216-296

Conflict of Interest Statement: There are no relevant conflicts of interest.

This work was supported by the National Institute on Deafness and Other Communication Disorders, grant DC-000421 to the University of Texas at Dallas

Abstract

Purpose. Successful speech processing depends on our ability to detect and integrate multisensory cues yet there is minimal research on multisensory speech detection and integration by children. To address this need, we studied the development of speech detection for auditory (A), visual (V), and audiovisual (AV) input.

Method. Participants were 115 typically-developing children clustered into age groups between 4-14 years. Speech detection (quantified by response times, RTs) was determined for one stimulus, /buh/, presented in A, V, and AV modes (articulating vs. static facial conditions). Performance was analyzed not only in terms of traditional mean RTs but also in terms of the faster vs. slower RTs (defined by 1st vs. 3rd quartiles of RT distributions). These time regions were conceptualized respectively as reflecting optimal detection with efficient focused attention vs. less optimal detection with inefficient focused attention due to attentional lapses.

Results. Mean RTs indicated better detection 1) of multisensory AV speech than A speech only in 4-5-yr-olds, and 2) of A and AV input than V input in all age groups. The faster RTs revealed that AV input did not improve detection in any group. The slower RTs indicated that 1) the processing of silent V input was significantly faster for the articulating than static face, and 2) AV speech or facial input significantly minimized attentional lapses in all groups except 6-7-yr-olds (a peaked U-shaped curve). Apparently, the AV benefit observed for mean performance in 4-5-yr-olds arose from effects of attention.

Conclusions. The faster RTs indicated that AV input did not enhance detection in any group, but the slower RTs indicated that AV speech and dynamic V speech (mouthing) significantly minimized attentional lapses and thus did influence performance. Overall, A and AV inputs were detected consistently faster than V input; this result endorsed stimulus-bound auditory processing by these children.

VISUAL SPEECH, DETECTION, & ATTENTION

3

1
2
3 When children engage in face-to-face conversations, they typically detect, discriminate, and identify
4 audiovisual speech sounds. Detection is the awareness that an audiovisual speech event occurred,
5 discrimination is the awareness that two audiovisual speech sounds differ from each other, and
6 identification is the labelling of the speech sounds. These different levels of speech perception tap
7 different levels of linguistic processing, which are, at least to some extent, hierarchical, and children
8 must detect and discriminate speech sounds before they can identify and label them (e.g., Aslin & Smith,
9 1988; Jerger, Martin, & Damian, 2002; McClelland & Elman, 1986; Stevenson, Sheffield, Butera, Gifford,
10 & Wallace, 2017). Gogate, Walker-Andrews, and Bahrck's (2001) model of early word acquisition—as it
11 relates to audiovisual speech—is an example of this hierarchical perceptual analysis. The model
12 proposes that when infants detect the redundancies between speech sounds and their corresponding lip
13 movements/mouth shapes, they can more readily discriminate similar-sounding phonological patterns,
14 such as “pin” and “tin,” and thus can recognize/label each pattern and associate it with its concept.
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29

30 In short, lower-level multisensory processes underpin higher-level multisensory speech perception
31 and word recognition skills, and altered lower-level processes can have cascading effects onto these
32 higher levels of processing. This relation is illustrated by the speech, language, and educational
33 difficulties observed in children with early onset hearing impairments and by the delayed expressive
34 language skills observed in children with early onset visual impairments (e.g., Briscoe, Bishop, &
35 Norbury, 2001; Eimas & Kavanagh, 1986; Jerger, Damian, Tye-Murray, Dougherty, Mehta, & Spence,
36 2006; McConachie & Moore 1994).
37
38
39
40
41
42
43
44

45 Despite the unquestionable contribution of detection and discrimination abilities to multisensory
46 speech perception and word recognition, these lower levels of multisensory speech processing,
47 particularly detection, are less well-studied in children than the higher-level speech recognition skills.
48 The extant discrimination literature indicates that visual speech (i.e., the articulatory gestures of talkers)
49 benefits phoneme discrimination in individuals ranging in age from infancy (e.g., Teinonen, Aslin, Alku &
50
51
52
53
54
55
56
57
58
59
60

VISUAL SPEECH, DETECTION, & ATTENTION

4

1
2
3 Csibra, 2008) to adulthood (e.g., Files, Tjan, Jiang & Bernstein, 2015). In children, visual speech improves
4
5 feature contrast discrimination (e.g., vi vs. zi, a place feature contrast, Hnath-Chisolm, Laipply, &
6
7 Boothroyd, 1998), vowel phoneme monitoring (Fort, Spinelli, Savariaux, & Kandel, 2012), and phoneme
8
9 discrimination for visually distinct contrasts (e.g., ba vs. ga, LaLonde & Holt, 2015; but see Boothroyd,
10
11 Eisenberg, & Martinez, 2010, for an exception).

12
13
14 With regard to age, improvements in the benefits from visual speech have been observed for
15
16 syllable/ nonword discrimination up to 7-yrs by Hnath-Chisolm et al. (1998) but up to 10-yrs by Fort et al.
17
18 (2012). In distinction to these results, however, Jerger, Damian, McAlpine, and Abdi (2017) recently
19
20 demonstrated that visual speech altered discrimination in all age groups from 4- to 14-yrs. These
21
22 researchers administered a same-different syllable-discrimination task, with the contrast of the critical
23
24 syllable pair requiring children to discriminate a syllable with an intact /b/ onset (e.g., /b/i) from the
25
26 same syllable but with a non-intact (spliced out) /-b/ onset (/-/b/i). Results showed that the presence or
27
28 absence of visual speech was critical for perception: the addition of visual speech to auditory speech
29
30 caused children to vote “same” when they listened to the intact: non-intact syllable pairs (e.g., /b/i : /-
31
32 b/i), a configuration implying that visual speech caused the non-intact onsets to be perceived as intact.
33
34 The degree of this “Visual Speech Fill-In Effect” for the non-intact onsets predicted the children's
35
36 receptive vocabulary skills.
37
38
39
40

41 In concert with the speech discrimination literature, the extant multisensory speech detection
42
43 literature indicates that adults detect audiovisual speech better than auditory speech (Bernstein, Auer, &
44
45 Takayanagi, 2004; Tjan, Chao, & Bernstein, 2013; Grant, 2001; Grant & Seitz, 2000; Kim & Davis, 2003 &
46
47 2004; LaLonde & Holt, 2016; Tye-Murray, Spehar, Myerson, Sommers, & Hale, 2011), and that infants
48
49 detect equivalent phonetic information in auditory and visual speech and changes in any mode
50
51 (auditory, visual, or audiovisual speech) for at least some conditions (e.g., Kuhl & Meltzoff, 1982;
52
53 Lewkowicz, 2000). In children, there is only one study, which reported that 6 – 8-yr-olds showed an
54
55
56
57
58
59
60

VISUAL SPEECH, DETECTION, & ATTENTION

5

1
2
3 adult-like detection advantage for audiovisual relative to auditory speech (LaLonde & Holt, 2016).
4

5 Although there is a dearth of information about multisensory speech detection by children, there is a
6
7 tenable child literature on the detection of non-speech multisensory inputs, such as a noise and a light.
8
9 This literature used *simple response times* to assess how quickly children can detect a pre-identified
10
11 sensory target and execute a preprogrammed motor response: Faster detection for the multisensory
12
13 compared to the uni-sensory inputs indicates multisensory facilitation. This literature reports that
14
15 children roughly 7-yrs and older detect simultaneous auditory and visual nonspeech inputs faster than
16
17 either uni-sensory input (Barutchu, Crewthe, & Crewther, 2009; Barutchu et al., 2011; Barutchu et al.,
18
19 2010; Brandwein et al., 2011; Gilley, Sharma, Mitchell, & Dorman, 2010). However, the degree of
20
21 facilitation is smaller and more variable in children than in adults up to about 14 – 15 years of age.
22
23
24

25 In short, proficient speech detection is critical for children to have access to the audiovisual cues that
26
27 underpin speech and language development, yet multisensory speech detection remains understudied
28
29 in children. To help address this gap in the literature, we studied the development of speech detection
30
31 as quantified by *simple response times* for uni-sensory speech (auditory or visual) vs. multisensory
32
33 speech (audiovisual) in children from 4 – 14-yrs-of-age. The stimulus in our study consisted of the
34
35 utterance “buh” presented in auditory only, visual only, and audiovisual modes. A primary research
36
37 question was whether children show enhanced detection of audiovisual speech relative to the uni-
38
39 sensory inputs.
40
41
42

43 Such enhanced detection is supported by evoked potential evidence in adults revealing that inputs
44
45 from the auditory and visual modalities interact at both early and late stages of sensory processing (e.g.,
46
47 Baart, Stekelenburg, & Vroomen, 2014; Molholm et al., 2002; van Wassenhove, Grant, & Poeppel, 2005).
48
49 This pattern of evoked potential findings has been interpreted to indicate that multisensory speech
50
51 perception is a *multi-staged* process with general spatial and temporal audiovisual speech
52
53 correspondences interacting early in processing and phonetic audiovisual speech features interacting
54
55
56
57
58
59
60

VISUAL SPEECH, DETECTION, & ATTENTION

6

1
2
3 later in processing (Baart et al., 2014; see also Schwartz et al., 2004). We should acknowledge that these
4
5 proposed stages of multisensory speech perception clearly occur before the behavioral response times
6
7 of individuals, which makes it difficult (as pointed out by Schroger & Widmann, 1998) to specify the
8
9 stage(s) of processing at which the auditory and visual inputs are interacting. Our experimental design—
10
11 the children responded to only one pre-identified speech syllable “buh” presented in the auditory,
12
13 visual, or audiovisual modes—clearly minimized the need for phonetic processing to identify the input.
14
15 That said, as speech input unfolds, it automatically activates corresponding phonological representations
16
17 according to the match between the evolving input and the representations in memory (e.g., Marslen-
18
19 Wilson & Zwitserlood, 1989; McClelland & Elman, 1986). Thus, the auditory and visual speech inputs of
20
21 this research may interact at any or all stages of analysis (see also Davis & Kim, 2004; Reisberg, McLean,
22
23 & Goldfield, 1987).
24
25
26

27
28 Another aspect of our experimental design was that the visual input consisted of either the dynamic
29
30 visual speech that produced the auditory “buh,” or the talker's static face. We included a static face not
31
32 only as a control condition but also because different types of previous studies have observed some
33
34 interesting differences between dynamic vs. static faces. First, accuracy on a task monitoring for an
35
36 auditory syllable in a carrier phrase is significantly better when adults view the talker's dynamic
37
38 articulating face vs. a static face (Davis & Kim, 2004). Second, although a dynamic articulating face and a
39
40 visual symbol both enhance the detection of auditory speech in adults, the dynamic articulating face
41
42 produces a relatively greater degree of multisensory facilitation (Bernstein et al., 2004; but see Tjan et
43
44 al., 2013). Third, dynamic faces—relative to static faces—enhance the recognition of emotional
45
46 expressions by adults and of unfamiliar faces by infants (Alves, 2013; Otsuka et al., 2009) possibly
47
48 because (as proposed by O'Toole, Roark, & Abdi, 2002) motion may enhance the perceptual processing
49
50 of faces and thus produce richer mental representations. And, fourth, a dynamic articulating face
51
52 generates more extensive cortical activation than a static face on fMRI scans (Calvert & Campbell, 2003;
53
54
55
56
57
58
59
60

VISUAL SPEECH, DETECTION, & ATTENTION

7

Campbell et al., 2001). Overall, the preponderance of this evidence predicts that performance in children may benefit more from the dynamic articulating face than from the static face.

Finally, we should note that dynamic faces are also more ecologically valid because they correspond to everyday social interactions, and this, in turn, may make them more attention provoking. In fact, some investigators propose that visual speech may act as a type of alerting mechanism that boosts attention, which helps children detect and process information faster (Campbell, 2006; Wickens, 1974). Thus, we also expect some differential effects of attention on the dynamic vs. static faces.

Attention is a key consideration because simple response time tasks as used herein are easy and monotonous—characteristics that are gold standards for assessing sustained attention (e.g., Betts, McKay, Maruff, & Anderson, 2006; Langner & Eickhoff, 2013; Manly et al., 2001). Sustained attention may be defined as “the ability to self-sustain mindful, conscious processing of stimuli whose repetitive, non-arousing qualities would otherwise lead to habituation and distraction.” (Robertson, Manly, Andrade, Baddeley, & Yiend, 1997, p. 747). Typically, younger children find it more difficult to sustain attention, and so they may find a simple response task particularly taxing because of their immature frontal cortex, which may limit the use of more automatic strategies (Thillay et al., 2015).

Children continue to improve their capacity to sustain attention up to the preteen/teenage years, with much of the developmental change occurring before 10 – 11 years (e.g., Betts et al., 2006; Dye & Bavelier, 2010; Manly et al., 2001; Thillay et al., 2015). Because of their immature sustained attention, younger children are more likely to experience difficulties in maintaining task goals, and this will increase the number of momentary lapses of attention and produce a larger number of slowed responses. Thus, the number of slowed responses is considered an index of these momentary attentional lapses (Key, Gustafson, Rentmeester, Hornsby, & Bess, 2017; Lewis, Reeve, Kelly, & Johnson, 2017; Venker et al., 2007; Weissman, Roberts, Visscher, & Woldorff, 2006). We predict that these occasional lapses producing slowed responses will create slower mean performance (based on all trials) in the younger

1
2
3 children than in the preteen-teenagers. To the extent that dynamic faces are more richly encoded and
4
5 more attention-provoking than static faces, we predict that performance for the dynamic face will show
6
7 fewer slowed responses. Below we describe how we assessed our data on the development of speech
8
9 detection (as defined by response times) for uni-sensory vs. multisensory inputs with two
10
11 complementary analyses.
12

13
14 Traditionally, the analysis of simple response times relies on a measure of central tendency—
15
16 typically the mean (see Laurienti, Burdette, Maldjian, & Wallace, 2006; Miller, 1988). Thus, in the first
17
18 analysis, we analyzed mean response times in the children divided into chronological age groups. In the
19
20 second analysis, however, we augmented this traditional approach by an analysis of the *faster vs. slower*
21
22 response times. The second analysis was motivated by the observation that mean performance does not
23
24 yield a pure measure of detection because, as noted above, the children's ability to detect sensory input
25
26 depends on their ability to sustain focused attention (e.g., Barutchu et al., 2009; Betts et al., 2006; Thillay
27
28 et al., 2015; see Footnote 1). Researchers studying age-related changes in elderly individuals have also
29
30 wrestled with the limitations of mean performance (e.g., Rabbitt & Goward, 1994; Rabbitt, Osman,
31
32 Moore, & Stollery, 2001; Tse, Balota, Yap, Duchek, & McCabe, 2010). Results in this arena that studied
33
34 faster vs. slower response times suggested: that elderly participants' fastest times are minimally affected
35
36 by increasing chronological age, and that differences in mean performance with age may
37
38 disproportionately reflect differences in the number of slowed times (see, e.g., Rabbitt et al., 2001). In
39
40 our second analysis, we interpreted results based on the rationale that optimal detection and efficient
41
42 sustained focused attention is located in the faster times, and less optimal detection with inefficient
43
44 sustained focused attention due to attentional lapses is located in the slower times (see Tse et al., 2010,
45
46 and Zhou & Krott, 2016, for similar reasoning). Both analyses are introduced by Data Analytic Sections
47
48
49
50
51
52 and Research Questions.
53

54 **Method**

55
56
57
58
59
60

Participants

Participants were 115 native English-speaking children ranging in age from 4;2 to 14;6 yrs (51% boys). The racial distribution was 84% White, 9% Asian, and 7% Black, with 9% reporting Hispanic ethnicity. Hearing sensitivity, visual acuity, auditory word recognition (Ross & Lerman, 1971), vocabulary skills (Dunn & Dunn, 2007), and visual perception (Beery & Beery, 2004) were within normal limits (age-based when appropriate) in all participants. Normal hearing sensitivity was defined as bilaterally symmetrical thresholds of ≤ 20 dB Hearing Level (HL) at all test frequencies between 500 and 4000 Hz (ANSI, 2010). Normal binocular visual acuity (including children with corrected vision) was defined as 8 correct out of 10 targets (5 each at 20/20 and 20/25 acuity) on the Lea Symbols presented in a light box that provided self-calibrating uniform illumination for testing (e.g., Becker, Hubsch, Graf, & Kaufmann, 2002; Good-lite Company, www.goodlite.com).

Participants were divided into four groups based on age (4 – 5-year-olds: $M = 4;11$, $SD = 0.52$, $N = 32$; 6 – 7-year-olds: $M = 7;0$, $SD = 0.59$, $N = 25$; 8 – 10-year-olds: $M = 9;3$, $SD = 0.89$, $N = 31$; and 11 – 14-year-olds: $M = 12;5$, $SD = 1.17$, $N = 27$). Advances in linguistic skills have been proposed to underlie developmental changes in sensitivity to visual speech (e.g., Desjardins, Rogers, & Werker, 1997; Erdener & Burnham, 2013; Jerger, Damian, Spence, Tye-Murray, & Abdi, 2009), and our age groups represented four different linguistic stages:

4 – 5-yr-olds: immature picture-book readers and immature speakers with articulatory deficiencies for complex sounds such as /sh/;

6 – 7-yr-olds: beginning readers whose phonology systems are reorganizing from phonemes as coarticulated indistinct speech sounds to phonemes as separable distinct written sounds and maturing speakers with good articulatory proficiency although with some dysfluencies;

8 – 10-yr-olds: maturing readers with blossoming mastery of phonemes as written and spoken sounds and strong articulatory skills; and

11 – 14-yr-olds: mature readers and speakers.

1
2
3 Adults were not included because results in the 11 – 14-yr-olds and young adults did not differ
4
5 statistically. Because auditory response times vary as a function of loudness, we should note that
6
7 average hearing sensitivity (pure tone average score at 500 Hz, 1000 Hz, and 2000 Hz) was similar across
8
9 the groups, ranging from 5.41 dB HL in the 4 – 5-year-olds to 2.24 dB HL in the 11 – 14-year-olds.

11 ***Materials and Instrumentation: Stimuli and Response Times***

12
13
14 ***Recording.*** The stimulus “buh” was recorded—as part of a set of Quicktime movie files for
15
16 associated projects—by an 11-yr-old boy actor with clearly intelligible speech without pubertal
17
18 characteristics (f_0 of 203 Hz). His full facial image and upper chest were recorded, and he started and
19
20 ended each utterance with a neutral face/closed mouth. The color video signal was digitized at 30
21
22 frames/s with 24-bit resolution at a 720 × 480 pixel size. The auditory signal was digitized at a 48 kHz
23
24 sampling rate with 16-bit amplitude resolution. The video track was routed to a high-resolution
25
26 computer monitor and the auditory track was routed through a speech audiometer to a loudspeaker
27
28 atop the monitor (see Jerger, Damian, Tye-Murray, & Abdi, 2014, for further details). For this project, the
29
30 stimulus started with the frame containing the auditory onset, and the talker’s lips in this beginning
31
32 frame remained closed but were no longer in a neutral position.
33
34
35

36
37 ***Stimulus.*** The stimulus “buh” was presented in three modes: audiovisual (AV), auditory only (A), and
38
39 visual only (V). For the AV presentation, children saw and heard the talker; for the A presentation, the
40
41 computer screen was blank; and for the V presentation, the loudspeaker was muted. Testing in these
42
43 three modes was carried out in two separate conditions: one with a dynamic face articulating the
44
45 utterance and one with an artificially static face (i.e., the child heard the same auditory track but the
46
47 video track was edited, with Adobe Premiere Pro, to contain only the talker's still face and upper chest of
48
49 the first frame). Hence, the two facial conditions consisted of presenting these two sets of items: 1) AV
50
51 dynamic face, V dynamic face, and A (no face); or 2) AV static face, V static face, and A (no face). The A
52
53 stimuli are the same in both facial conditions, thus allowing us to estimate test-retest reliability.
54
55
56
57
58
59
60

1
2
3 We formed one list of 39 test items (13 in each mode) for each facial (dynamic and static) condition
4
5 (each list was presented forwards and backwards to yield two variations). The items of each list were
6
7 randomized with the constraint that /buh/ was presented once in each mode for each triplet of items
8
9 (e.g., two-triplet sequence = A/ AV/ V/ V/ A/ AV). This design assured that any changes in performance
10
11 due to personal factors (e.g., fatigue, practice) would be equally distributed over all modes.
12
13

14 **Response Times.** To obtain response times, the computer triggered a counter/timer (resolution less
15
16 than one ms) at the initiation of a stimulus. The stimulus continued until pressure on a response
17
18 (telegraph) key stopped the counter/timer. The response board contained two keys separated by a
19
20 distance of approximately 12 cm. A green square beside each key designated the start position for the
21
22 child's hand. The key corresponding to the response (right vs. left) was counterbalanced across
23
24 participants, and a small temporary box covered the unused key.
25
26

27 **Procedure**

28
29 Testing was carried out within a double-walled sound-treated booth. The data of this study were
30
31 gathered in one session of a multiple-day experimental protocol (e.g., Jerger et al., 2014; Jerger, Damian,
32
33 Tye-Murray, & Abdi, 2016 & 2017; Jerger, Damian, Parra, & Abdi, 2017; Jerger et al., 2017). The
34
35 presentation order of the facial conditions was counterbalanced across participants in each age group.
36
37 One facial condition (either dynamic or static) was administered, followed by about 30 minutes of other
38
39 testing, followed by the administration of the other facial condition. For the formal testing, a tester sat
40
41 at a computer workstation and initiated each trial, in an arrhythmic manner, when the child appeared
42
43 ready by pressing a touch pad (out of child's sight). A co-tester sat alongside each child to help keep the
44
45 child "on task" at least overtly at the start of each trial—defined as sitting attentively and looking at the
46
47 monitor with his or her hand on the start position. The children sat at a distance of 71 cm directly in
48
49 front of an adjustable height table containing the computer monitor and loudspeaker. The children's
50
51 view of the talker's face subtended a visual angle of 7.17° vertically (eyebrow to chin) and 10.71°
52
53
54
55
56
57
58
59
60

1
2
3 horizontally (eye level). The children heard the auditory input at an intensity of approximately 70 dB SPL.
4

5 The children were told that they would sometimes hear, sometimes see, and sometimes hear and
6 see a boy. When the boy was talking, he would always be saying “buh.” When they saw the boy, they
7 were told that they would see a movie of the boy (dynamic face) for one facial condition and a photo of
8 the boy (static face) for the other facial condition. Before each condition, the children were shown the
9 stimulus for each mode (A, V, and AV). They were told to push the key as fast as possible to the onset of
10 any of these targets with a whole hand response (the tester illustrated and the child imitated). The
11 children were told to always start with their hand on the green square—and as soon as they hit the key,
12 to be sure to put their hand back on the square and get ready for the next target. Prior to the
13 administration of each facial condition, practice trials were administered until response times had
14 stabilized across a two-triplet sequence. Flawed trials (i.e., on rare occasions, the equipment
15 malfunctioned or the child moved out of position to do something after trial started) were deleted on-
16 line and re-administered at the end of the list.
17
18
19
20
21
22
23
24
25
26
27
28
29
30

31 32 **Analysis of Mean Response Times** 33

34 ***Data Analysis*** 35

36 We compared mean performance in each mode for each facial condition. Mean values are preferred
37 because median values can provide biased estimates for response time distributions with different
38 skewness and/or different or small sample sizes (Miller, 1988; Whelan, 2008). The mean values are
39 reported in the text/graphs because they clearly show how performance differed between the age
40 groups and the modes, but, for all inferential statistical analyses, the individual values were log
41 transformed to normalize the distribution (Heathcote, Popiel, & Mewhort, 1991; Whelan, 2008). The
42 Bonferroni correction controlled the familywise alpha (Abdi, Edelman, Valentin, & Dowling, 2009).
43
44
45
46
47
48
49
50

51 To determine whether AV speech produced faster detection for each facial condition, we evaluated
52 the difference between response times in the AV mode minus the fastest uni-sensory mode as per the
53
54
55
56
57
58
59
60

1
2
3 fixed favored dimension model for multidimensional stimuli (e.g., Biederman & Checkosky, 1970;
4
5 Mordkoff & Yantis, 1993; Stevenson et al., 2014). Both the dynamic and static faces were viewed as
6
7 multidimensional AV stimuli because individuals can accurately match unfamiliar voices to both dynamic
8
9 and static unfamiliar faces well above chance; this pattern of results indicates that voices share source-
10
11 identity information with both types of faces (Krauss, Freyberg, & Morsella, 2002; Mavica & Barenholtz,
12
13 2013; Smith, Dunn, Baguley, & Stacey, 2016 a and b; but see Lachs & Pisoni, 2004). Accurate voice-face
14
15 matching would be particularly prominent in our children because they were familiar with the talker's
16
17 face and voice from the other tasks they performed in our multiple-day experimental protocol. We
18
19 predicted that the A response times would comprise the fastest uni-sensory mode because our pilot
20
21 data in children and an extensive literature in adults indicate that response times are faster for the A
22
23 than V mode (e.g., Diederich & Colonius, 2004; Harrar et a, 2014; Vickers, 2007; Woodworth &
24
25 Schlosberg, 1954). Our research questions were: 1) Do children respond faster to A than V input as
26
27 indicated in the adult literature? 2) Do children respond faster to AV than to the fastest uni-sensory
28
29 input? 3) Do children's response times differ in the facial conditions? And 4) Are children's response
30
31 times reliable?
32
33
34
35

36 **Results**

37 **Mean Response Times**

38
39
40
41 Figure 1 compares response times in the A, V, and AV modes for the static and dynamic faces in the
42
43 four age groups and in the entire group. Statistical analyses (summarized in Table 1) were performed
44
45 with a mixed-design analysis of variance (ANOVA) with one between-participant factor (Age Group: 4–5-
46
47 yrs, 6–7-yrs, 8–10-yrs, and 11–14-yrs) and two within-participant factors (Mode: V, A, and AV; Facial
48
49 Condition: static vs. dynamic). Results revealed a significant Age Group effect, which occurred because
50
51 response times (collapsed across Mode and Facial Condition) were slower in the younger than in the
52
53 older children: Mean response times were 814 ms in the 4 – 5-yr-olds but 508 ms in the 11 – 14-yr-olds.
54
55
56
57
58
59
60

1
2
3 A significant Mode effect was also observed, which occurred because response times (collapsed across
4
5 Age Group and Facial Condition) were significantly faster for the A and AV modes (592 ms and 577 ms)
6
7 than for the V mode (752 ms). A straightforward interpretation of this latter result was complicated,
8
9 however, by a significant Mode \times Facial Condition interaction, which occurred because mean response
10
11 times (collapsed across Age Group; see *All* in Figure 1) were faster for the dynamic than the static face
12
13 for V input (728 ms – 776 ms) but not for A and AV input (respectively 597 ms – 587 ms for A and 575
14
15 ms – 578 ms for AV).
16
17

18
19 Insert Figure 1 and Table 1
20

21 Below, as we turn to analyzing whether the uni-sensory inputs differed, the above results inform us
22
23 about the V vs. A modes. The significant Mode effect indicated that the A response times were faster
24
25 than the V response times. The significant Mode \times Facial Condition interaction indicated that this
26
27 difference between the V and A response times was greater for the static face (189 ms) than the
28
29 dynamic face (131 ms). There was no significant interaction involving the Age Groups; thus (as shown in
30
31 Figure 1) these significant differences characterized all Groups. Below, we addressed whether the AV
32
33 and A modes differed in any of the age groups or facial conditions.
34
35

36 **AV vs. A Modes.** To probe whether responses to AV input were faster than responses to A input (the
37
38 fastest uni-sensory input), we carried out planned orthogonal contrasts for each facial condition in each
39
40 age group (Abdi & Williams, 2010). Results indicated that the dynamic face (i.e., dynamic AV speech) was
41
42 associated with faster responses only in the 4 – 5-yr-olds, $F_{\text{contrast}}(1, 110) = 9.73$, $MSE = .001$, $p = .002$,
43
44 partial $\eta^2 = .042$. No other significant contrast was observed.
45
46

47 **Reliability.** To assess test-retest performance for the A response times, we reformatted the data to
48
49 represent the first vs. second tests (the two facial conditions were counterbalanced such that each
50
51 occurred as the first test half of the time). The response times were statistically evaluated with a mixed-
52
53 design ANOVA with one between-participant factor (Age Group: 4–5-yrs, 6–7-yrs, 8–10-yrs, and 11–14-
54
55
56
57
58
59
60

1
2
3 yrs) and one within-participant factor (Test: first vs. second). Results indicated that there was no
4
5 significant effect of test nor any test \times group interaction. A follow-up simple regression analysis (Abdi et
6
7 al., 2009) in the entire group indicated that the children's A response times for the first and second tests
8
9 were significantly correlated, $r = .840$, $F(1, 114) = 270.12$, $p < .0001$. The slope of the regression line was
10
11 0.768, which indicates there was a 0.768 unit change in the second-session responses for each one unit
12
13 change in the first-session responses. The variance (mean square) residual, or the degree of variability of
14
15 the individual data about the regression line, was 0.004. The mean auditory response times for the first
16
17 and second test sessions in the entire group were 599 ms ($SD = 190$ ms) and 585 ms ($SD = 153$ ms), and
18
19 the individual difference scores for the first test minus the second test averaged 15 ms, with a 95%
20
21 confidence interval ranging from -5 ms to 35 ms.
22
23
24

25 **Summary.** The children's mean response times became significantly faster as age increased, a result
26
27 which agrees with previous findings (e.g., Goodenough, 1935; Jerger, Martin, & Pirozzolo, 1988). The
28
29 children also responded faster to the A than the V input, a pattern consistent with the literature noted
30
31 above. This A-faster-than-V pattern of results was observed in 97% – 98% of the children for the two
32
33 facial conditions. With regard to whether the children responded faster to AV than A input, the addition
34
35 of V speech was associated with faster responses but only in the 4 – 5-yr-olds. The AV-faster-than-A
36
37 pattern of results in the dynamic facial condition was observed in 78% of the 4 – 5-yr-olds. A silent V
38
39 speech (i.e., mouthing) effect was also observed in that responses in the V mode were faster for the
40
41 dynamic facial condition than the static facial condition. This mean pattern of results was observed in
42
43 67% of the children. Evaluation of test-retest performance established highly reliable results.
44
45
46

47 **Analysis of Faster vs. Slower Response Times**

48 **Data Analysis**

49
50
51 Mean performance in the above analyses may reflect a shift of the entire response time
52
53 distribution or a shift of only the slow tail or the skewness of the distribution (e.g., see Balota & Yap,
54
55
56
57
58
59
60

VISUAL SPEECH, DETECTION, & ATTENTION

16

1
2
3 2011; Rabbitt et al., 2001). We explored possible differences in the faster vs. slower times with response
4
5 time distributions computed by Vincentile analysis, a nonparametric technique that preserves the
6
7 component distributions' shapes and does not make assumptions about the underlying distribution (see
8
9 Jiang, Rouder, & Speckman, 2004; Ratcliff, 1979). Vincentile analysis is especially recommended because
10
11 it provides stable estimates even with a small number of response times per participant/condition.
12
13

14 To obtain the Vincentile distributions, each child's response times—for each mode/facial condition—
15
16 were rank-ordered and then initially divided into sequential bins of 10% (deciles). A *cumulative*
17
18 *distribution function* (CDF) was obtained for each age group by averaging each of the bins across the
19
20 participants in that group for each facial condition/mode. Figure 1A (Appendix) portrays the CDFs for the
21
22 A, AV, and V modes in the static (1A_a) and dynamic (1A_b) facial conditions for all age groups. In adults,
23
24 CDFs such as these are explored with ex-Gaussian analyses of the response distributions, but we did not
25
26 have a sufficient number of trials to conduct this type of analysis (Heathcote et al., 1991). Thus, we
27
28 computed another set of Vincentile distributions by dividing each child's rank-ordered response times—
29
30 for each mode/facial condition—into sequential bins of 25% (quartiles). Statistically we investigated
31
32 whether our effects of interest appeared in the faster and/or slower response times by analyzing the
33
34 25th and 75th (i.e., 1st and 3rd) quartiles of the Vincentile CDFs. Again, our assumptions for interpreting
35
36 the results are that optimal detection and efficient focused attention is located in the faster times (1st
37
38 quartile), and less optimal detection with inefficient focused attention due to attentional lapses is
39
40 located in the slower times (3rd quartile). We were interested in whether the pattern of mean results
41
42 reported above was observed at both quartiles (results influenced by both detection and attention) or at
43
44 only one of the quartiles (results influenced by only detection or attention). To assess this, we carried
45
46 out contrast analyses (Abdi & Williams, 2010) on the log transformed response times at the 1st/faster
47
48 and 3rd/slower quartiles for each facial condition in each age group with a Bonferroni correction to
49
50 control the familywise alpha. Our focused research questions were: 1) Do the A vs. V inputs differ in the
51
52
53
54
55
56
57
58
59
60

1
2
3 age groups at both quartiles or only the 1st/faster or 3rd/slower quartile? 2) Do the AV vs. fastest uni-
4 sensory inputs differ in any age group at one or both quartiles? And 3) Does the facial condition affect
5 these results?
6
7
8

9 **Results**

10 **Faster vs. Slower Response Times**

11
12 **V vs. A Modes.** Figure 2 shows the mean difference scores (V response times – A response times) in
13 the age groups at each quartile for the static and dynamic facial conditions. Table A1 (Appendix)
14 presents the F_{contrast} results for the V vs. A modes. The large positive difference scores in Figure 2 along
15 with the statistical results documented that the V response times were significantly slower than the A
16 response times in all age groups at both quartiles for both facial conditions. Relative to the V input, the A
17 input was detected faster and with significantly fewer attentional lapses (see also CDFs in Appendix).
18 Faster A-than-V responses were observed in about 97% of children for both facial conditions at both
19 quartiles.
20
21
22
23
24
25
26
27
28
29
30

31
32 Insert Figure 2
33

34 As indicated by the asterisks in Figure 2 and as documented by the F_{contrast} results for the dynamic vs
35 static faces in Table A2 (Appendix), dynamic V speech—relative to a static face—decreased the mean
36 difference scores significantly at the 3rd quartile/slower responses but not at the 1st quartile/faster
37 responses, with the exception of results in the 8 – 10-yr-olds which did not differ for the facial conditions
38 at either quartile. These results indicate that dynamic V speech captured attention and reduced
39 attentional lapses more than the static face, with about 75% of children not including the 8 – 10-yr-olds
40 showing this pattern of results. Reasons for the different pattern of results in the 8 – 10-yr-olds are
41 unclear, and indeed about 60% of these children showed the typical pattern of results for the dynamic
42 vs. static facial conditions.
43
44
45
46
47
48
49
50
51
52

53
54 **AV vs. A Modes.** Figure 3 shows the mean difference scores (AV response times – A response times)
55
56
57
58
59
60

VISUAL SPEECH, DETECTION, & ATTENTION

18

1
2
3 in the age groups at each quartile for the static and dynamic facial conditions. Table A3 (Appendix)
4 presents the F_{contrast} results for the AV vs. A modes. Statistical findings in Table A3 and the differences
5 scores in Figure 3 for the 1st/faster quartile showed that multisensory AV input did not improve
6 detection in any age group. With regard to the 3rd/slower quartile, AV dynamic speech captured and
7 benefited attention in the 4 – 5-yr-olds and the 11 – 14-yr-olds, and static facial input benefited
8 attention in the 8 – 10-yr-olds. This pattern of results was observed in about 75% of children in each of
9 these age groups. Finally, as indicated by the asterisk in Figure 3 and as documented by the F_{contrast}
10 results for the dynamic vs static faces in Table A4 (Appendix), differences between the facial conditions
11 achieved statistical significance only in the 4 – 5-yr-olds at the 3rd/slower quartile, with 55% of these
12 children showing a greater difference score for the dynamic face.
13
14
15
16
17
18
19
20
21
22
23
24

25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Insert Figure 3

Because we know little about the influence of attention on AV multisensory speech perception by children, we re-assessed the results in Figure 3 at the 3rd/slower quartile with a mixed-design ANOVA with one between-participant factor (Age Group: 4–5-yr, 6–7-yr, 8–10-yr, and 11–14-yr) and two within-participant factors (Mode: A vs. AV; Facial Condition: static vs. dynamic; see Footnote 2). As always, the individual values were log transformed to normalize the distribution (Heathcote et al., 1991; Whelan, 2008), and the Bonferroni correction controlled the familywise alpha (Abdi et al., 2009). However, two considerations influenced how we carried out the current Bonferroni correction. First, a standard omnibus ANOVA is a non-specific, global test that seeks any differences within or between factors (even ones that are not of interest) and suffers from low statistical power relative to procedures that decompose the systematic variance into meaningful contrasts (Rosenthal, Rosnow, & Rubin, 2000). Second, false negatives can be a more fundamental problem than false positives in an area with little evidence because they may retard further meaningful growth of knowledge (Fiedler, Kutzner, & Krueger, 2012). Thus, as recommended when some F values in an omnibus ANOVA are more important than

1
2
3 others *a priori*, we allocated the individual α 's per family of tests unequally for the Bonferroni correction
4
5 (Abdi, 2007). We tested the critical Mode \times Facial Condition \times Age Group interaction with an $\alpha = .04$ and
6
7 shared the remaining .01 between the other F tests, which were evaluated with an $\alpha = .0017$.
8

9
10 Statistical findings are summarized in Table 2. Results revealed a significant Age Group effect, which
11
12 occurred because response times (collapsed across Mode and Facial Condition) were slower in the
13
14 younger than in the older children as noted previously. A significant Mode effect was also observed,
15
16 which occurred because response times (collapsed across Age Group and Facial Condition) were
17
18 significantly faster for the AV than the A mode (591 ms and 615 ms). A straightforward interpretation of
19
20 this latter result was complicated, however, by a significant Mode \times Facial Condition \times Age Group
21
22 interaction, which indicated that the relationship between the AV and A response times differed for the
23
24 Facial Conditions but in inconsistent ways across the Age Groups. Critically, this interaction points out
25
26 that the relationship between the AV and A response times varied across the Age Groups. To probe this
27
28 pattern of interaction, we conducted *t*-tests on the difference between the AV vs A response times in
29
30 each Age Group for each Facial Condition. Results are summarized in Table 3. Results mirrored the
31
32 previously obtained F_{contrast} findings. The significant differences between the AV and A response times
33
34 indicated that AV dynamic speech benefited attention in the 4 – 5-yr-olds and the 11 – 14-yr-olds, and
35
36 static facial input benefited attention in the 8 – 10-yr-olds. In short, facial input (either AV dynamic
37
38 speech or a static face) significantly influenced attention in all age groups, excepting the 6 – 7-yr-olds.
39
40
41
42

43 Insert Table 2 and Table 3

44 45 **Discussion**

46
47
48 Everyday tasks depend on our ability to detect and integrate information from multiple sensory
49
50 modalities. Despite the acknowledged importance of this lower level of processing for speech, however,
51
52 we know little about children's multisensory speech detection abilities. The purpose of this research was
53
54 to study the development of speech detection for A, V, and AV inputs in children from 4 – 14-yr of age.
55
56
57
58
59
60

Our experimental design featured two novel approaches. First, our V input consisted of both static and dynamic faces, which allowed us to determine whether effects on performance reflected a facial effect or an articulating-face-specific effect (influenced only by the dynamic face). Second, we assessed development not only in terms of the traditional mean response times but also in terms of the faster vs. slower response times. We should acknowledge that some of the slower response times in these children may have been reflecting motivational factors rather than attentional lapses (see Reinvang, 1998). This research, however, minimized this possibility by having a co-tester who tried to keep the children engaged in the task. We should also note that there were only 13 trials per condition (78 trials total) due to the limited testing time available with young children. Importantly, however, we selected a technique (Vincentizing) that is especially suitable for analyzing data with only a few observations per condition (i.e., it has been shown that the Vincentizing provides stable estimates even with only 10 – 20 trials per participant/condition, see Jiang et al., 2004; Ratcliff, 1979).

We discuss the results below in terms of the uni-sensory inputs (V vs. A) and the multisensory input vs. the fastest uni-sensory input (AV vs. A). A focus is to understand how the results for the 1st/faster and 3rd/slower quartiles contributed to the interpretation of mean performance in the children. These two time regions were respectively conceptualized as reflecting optimal detection with efficient focused attention vs. less optimal detection with inefficient focused attention due to attentional lapses.

V vs. A Inputs. Mean performance in the age groups indicated significantly faster A than V response times and significantly faster V responses for the silent dynamic face (i.e., mouthing) than the static face. The A-faster-than-V outcome agrees with long-term previous findings in adults (Diederich & Colonius, 2014; Harrar et al., 2014; Vickers, 2007; Woodworth & Schlosberg, 1954; see Brandwein et al., 2011, and Gilley et al., 2010, for exceptions). Analysis of the faster vs. slower response times indicated that 1) A input relative to V input not only facilitated the children's ability to detect the input but also reduced their attentional lapses whereas 2) silent dynamic V speech (mouthing) relative to a static face only

VISUAL SPEECH, DETECTION, & ATTENTION

21

1
2
3 reduced attentional lapses. This latter finding supports the proposal that a dynamic face may be more
4
5 richly encoded and thus more attention-provoking than a static face (Calvert & Campbell, 2003;
6
7 Campbell et al, 2001; O'Toole et al., 2002). Overall, the pattern of results implies that the changes in
8
9 mean performance could be reflecting effects of detection and/or attention.
10

11
12 The significantly faster speed of processing for A than V inputs strongly supports stimulus-bound
13
14 auditory processing and an automatic capture of attention by A input in these children (e.g., Sloutsky &
15
16 Napolitano, 2003; Napolitano & Sloutsky, 2004). These results are reminiscent of the auditory distraction
17
18 literature in adults (e.g., Macken, Phelps, & Jones, 2009; Watkins, Dalton, Lavie, & Rees, 2007), which
19
20 emphasizes the capacity of A input to capture attention despite adults' attempts to "not listen." Such
21
22 findings have impactful implications for speech and language development in children. As an example—
23
24 if we view the speech input more narrowly as A only and the V input more broadly as environmental
25
26 objects—pretend that a parent looks and points to an object while saying "lamp" to his or her
27
28 preschoolers. The V input in this example is permanent, but the A input is fleeting. If the children fail to
29
30 see the "lamp" at first glance, they can easily see it by taking another look. If, however, the children fail
31
32 to hear the word at first listen, they cannot easily hear it by taking another listen. Thus, the automatic
33
34 capture of attention by A input in young children may critically nurture speech and language
35
36 development because it helps children perceive words that are "written on the wind."
37
38
39
40

41
42 The unequal detection of the A and V dimensions of speech in this research may reflect, at least to
43
44 some degree, the conscious behaviors demanded by our experimental protocol. That said—to the extent
45
46 that these results generalize to AV speech perception with its more unconscious detection of the A and V
47
48 dimensions—these results may inform the interpretation of studies that manipulated the onsets of the A
49
50 and V cues and found that individuals are more likely to synthesize these cues when the V speech starts
51
52 before the A speech than vice versa. For example, in adults AV interactions occur even when the V
53
54 speech leads the A speech by 170 ms to 180 ms (Munhall, Gribble, Sacco, & Ward, 1996; van
55
56
57
58
59
60

1
2
3 Wassenhove, Grant, & Poeppel, 2007). In contrast, when A speech leads the V speech, AV interactions
4
5 occur only up to an asynchrony of 30 ms (e.g., van Wassenhove et al., 2007). This pattern of AV
6
7 interactions for asynchronous speech appears to be adult-like by 7-yr of age, although children do not
8
9 show the same degree of AV interactivity as adults (Hillock-Dunn, Grantham, & Wallace, 2016). A
10
11 greater tolerance for V-speech-leading asynchronies seems to have ecological validity because V cues
12
13 frequently start before A cues in everyday speech (e.g., Bell-Berti & Harris, 1981). That said, the current
14
15 research suggests that the greater tolerance of V-speech-leading asynchronies may also be reflecting
16
17 people's slowness in detecting V speech relative to A speech.
18
19

20
21 **AV vs. A Inputs.** Mean performance showed that response times were faster for dynamic AV input
22
23 than A input but only in the 4 – 5-yr-olds. Analysis of the faster and slower times, however, indicated
24
25 that AV dynamic speech did not influence detection (i.e., responses at the 1st/faster quartile) in any
26
27 group. These results disagree with the one previous study of speech detection by children, which
28
29 reported adult-like benefits from AV speech in 6 – 8-yr-olds on a task requiring detection of speech in
30
31 noise (Lalonde & Holt, 2016). Our results also show a different developmental course from the one
32
33 characterizing the detection advantage for nonspeech multisensory A and V inputs. The non-speech child
34
35 literature in the Introduction was provided because there are few multisensory speech detection studies
36
37 in children. We should note, however, that this non-speech A and V literature cannot be *directly* related
38
39 to the AV speech findings because speech dimensions/cues are processed in an interdependent
40
41 (conjoined) manner (Garner, 1974; Green & Kuhl, 1989; Jerger, Martin, Pearson, & Dinh, 1995; Jerger et
42
43 al., 1993; Tomiak, Mullennix, & Sawusch, 1987) whereas arbitrarily-paired inputs such as a noise and a
44
45 light are typically processed in an independent (separable) manner (e.g., Garner, 1974; Marks, 2004).
46
47 Thus, our different results are difficult to interpret due to the pronounced task differences along with
48
49 different perceptual processing structures that preclude an unambiguous comparison of speech vs. non-
50
51 speech research.
52
53
54
55
56
57
58
59
60

1
2
3 With regard to the 3rd quartile/slower response times, AV dynamic speech captured attention and
4 thus significantly minimized slowed responses relative to A speech in the 4 – 5-yr-olds and 11 – 14-yr-
5 olds. This AV effect seems reminiscent of the U-shaped curve we observed previously in which AV
6 phonologically-related speech distractors primed picture naming in 4 – 5-yr-olds and 10 – 14-yr-olds but
7 not in children of in-between ages (Jerger et al., 2009). The current results, however, additionally
8 revealed that AV static facial input significantly minimizes attentional lapses and thus, slowed responses
9 in the 8 – 10-yr-olds as well. In short, V speech or facial input relative to A speech significantly impacted
10 results in all age groups except in the 6 – 7-yr-olds (a peaked U-shaped curve).
11
12
13
14
15
16
17
18
19
20

21 Previously Jerger et al. (2009) related their U-shaped results to dynamic systems theory (e.g., Smith
22 & Thelen, 2003), which proposes that: 1) multiple factors typically underlie developmental change, and
23 2) a lack of any effect in children may be reflecting a period of transition (not a lack of effect) during
24 which immature knowledge and processing subsystems are reorganized and restructured into more
25 mature, elaborated, and robust forms. During these developmental transitions, processing systems are
26 less robust, and children cannot easily use their cognitive resources; consequently, during these
27 transitional stages, children's performance can be unstable and affected by methodological approaches
28 and task demands (Evans, 2002).
29
30
31
32
33
34
35
36
37
38

39 We propose that the developmental shifts in AV performance for the slowed times reflect different
40 stages of reorganization and transition. With regard to the 4 – 5-yr-olds and the 11 – 14-yr-olds, we
41 should note that alike performance in these groups may not be reflecting alike underlying mechanisms.
42 Whereas performance in the 11 – 14-yr-olds is mature and reflects dynamic AV speech capturing
43 attention and minimizing attentional lapses, performance in the 4 – 5-yr-olds is immature and may be
44 reflecting a dynamic AV speech effect and/or other factors. For example, 3-yr-olds and thus perhaps 4 –
45 5-yr-olds attend preferentially to dynamic over static faces (Libertus, Landa, & Haworth, 2017), and
46 younger children with less mature articulatory proficiency observe V speech more, perhaps to cement
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 their knowledge of the acoustic consequences of articulatory gestures (Desjardins et al., 1997; Dodd,
4
5 McIntosh, Erdener, & Burnham, 2008).

6
7 Performance in the 6 – 7-yr-olds did not show any influence of either type of face, but performance
8
9 in the 8 – 10-yr-olds revealed the minimization of attentional lapses by AV static facial input—an effect
10
11 which may reflect the simultaneous or correlated onsets interacting to produce a more emphatic onset-
12
13 alerting signal. As noted previously, voices share source-identity information with both the dynamic and
14
15 static faces (Krauss et al., 2002; Mavica & Barenholtz, 2013; Smith et al., 2016 a and b). We propose that
16
17 the different results in the 6 – 7-yr-olds and 8 – 10-yr-olds occurred because the relevant knowledge and
18
19 processing subsystems, particularly phonology, were reorganizing between roughly 6 – 9 years of age
20
21 into more mature resources for a wider range of activities (see Jerger et al., 2009, for discussion and
22
23 references). Phonological processes are particularly relevant because, even though this task minimized
24
25 phonological processing demands, speech input automatically activates corresponding phonological
26
27 representations as it unfolds as noted previously (e.g., Marslen-Wilson & Zwitserlood, 1989; McClelland
28
29 & Elman, 1986). Thus, the A and V inputs of this research may interact at multiple stages of analysis,
30
31 which can also be influenced by cognitive resources such as attention (e.g., Davis & Kim, 2004; Reisberg
32
33 et al., 1987). Finally, we should acknowledge that both this research and the Jerger et al. (2009) research
34
35 studied response times. The measurement of processing speed can be a more sensitive measure of task
36
37 proficiency. That said, all methods of identifying and quantifying multisensory interactions have
38
39 advantages and disadvantages (Stevenson et al., 2014).

40 41 42 43 44 45 **Conclusions**

46
47 These results emphasized the pronounced ability of both AV speech and silent dynamic V speech
48
49 (mouthing) to minimize attentional lapses and thus influence detection. Such findings demonstrate the
50
51 usefulness of V speech even in situations that do not involve impoverished A input. Another primary
52
53 result was that response times were always faster to A and AV input than to V input. Our overall results
54
55
56
57
58
59
60

VISUAL SPEECH, DETECTION, & ATTENTION

strongly endorsed stimulus-bound auditory processing by these children. Such findings are good news for children who must listen to learn.

For Peer Review

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Acknowledgements

This research was supported by the National Institute on Deafness and Other Communication Disorders, grant DC-000421 to the University of Texas at Dallas. We thank the children and parents who participated and the researchers who assisted, namely Aisha Aguilera, Carissa Dees, Nina Dinh, Nadia Dunkerton, Alycia Elkins, Brittany Hernandez, Demi Krieger, Rachel Parra McAlpine, Michelle McNeal, Jeffrey Okonye, and Kimberly Periman of University of Texas at Dallas (data collection, analysis, presentation), and Derek Hammons and Scott Hawkins of University of Texas at Dallas and Drs. Brent Spehar and Nancy Tye-Murray of Washington University School of Medicine (computer programming, stimuli recording/editing).

For Peer Review

References

- Abdi, H., & Williams, L. (2007). Bonferroni and Sidak corrections for multiple comparisons. In N. Salkind (Ed.), *Encyclopedia of measurement and statistics* (pp. 103-107). Thousand Oaks, CA: Sage.
- Abdi, H., Edelman, B., Valentin, D., & Dowling, W. (2009). *Experimental design and analysis for psychology*. New York: Oxford University Press.
- Abdi, H., & Williams, L. (2010). Contrast analysis. In N. Salkind (Ed.), *Encyclopedia of research design* (pp. 243-251). Thousand Oaks, CA: Sage.
- Alves, N. (2013). Recognition of static and dynamic facial expressions: A study review. *Estudos de Psicologia*, 18, 125-130.
- American National Standards Institute. (2010). *Specifications for audiometers*. ANSI/ASA S3.6-2010 (R2010). New York: American National Standards Institute.
- Aslin, R., & Smith, L. (1988). Perceptual development. *Annual Review of Psychology*, 39, 435-473.
- Baart, M., Stekelenburg, J., & Vroomen, J. (2014). Electrophysiological evidence for speech-specific audiovisual integration. *Neuropsychologia*, 53, 115-121.
- Balota, D., & Yap, M. (2011). Moving beyond the mean in studies of mental chronometry: The power of response time distributional analyses. *Current Directions in Psychological Science*, 20, 160-166.
- Barutchu, A., Crewther, D., & Crewther, S. (2009). The race that precedes coactivation: Development of multisensory facilitation in children. *Developmental Science*, 12, 464-473.
- Barutchu, A., Crewther, S., Fifer, J., Shivdasani, M., Innes-Brown, H., Toohey, S., . . . Paolini, A. (2011). The relationship between multisensory integration and IQ in children. *Developmental Psychology*, 47, 877-885.
- Barutchu, A., Danaher, J., Crewther, S. G., Innes-Brown, H., Shivdasani, M. N., & Paolini, A. G. (2010). Audiovisual integration in noise by children & adults. *Journal of Experimental Child Psychology*, 105, 38-50.
- Becker, R., Hubsch, S., Graf, M., & Kaufmann, H. (2002). Examination of young children with Lea symbols. *British Journal of Ophthalmology*, 86, 513-516.
- Beery, K., & Beery, N. (2004). *The Beery-Buktenica Developmental Test of Visual-Motor Integration with Supplemental Developmental Tests of Visual Perception and Motor Coordination*. (5th ed.). Minneapolis:

VISUAL SPEECH, DETECTION, & ATTENTION

28

1
2
3 NCS Pearson, Inc.

4
5 Bell-Berti, F., & Harris, K. (1981). A temporal model of speech production. *Phonetica*, 38, 9–20.

6
7 Bernstein, L., Auer Jr, E., & Takayanagi, S. (2004). Auditory speech detection in noise enhanced by lipreading.
8
9 *Speech Communication*, 44, 5-18.

10
11 Betts, J., McKay, J., Maruff, P., & Anderson, V. (2006). The development of sustained attention in children:
12
13 The effect of age and task load. *Child Neuropsychology*, 12, 205-221.

14
15 Biederman, I., & Checkosky, S. (1970). Processing redundant information. *Journal of Experimental Psychology*,
16
17 83, 486-490.

18
19 Boothroyd, A., Eisenberg, L., & Martinez, A. (2010). An on-line imitative test of speech-pattern
20
21 contrast perception (OlimSpac): Developmental effects in normally hearing children.
22
23 *Journal of Speech, Language, and Hearing Research*, 53, 531-542.

24
25 Brandwein, A., Foxe, J., Russo, N., Altschuler, T., Gomes, H., & Molholm, S. (2011). The development of
26
27 audiovisual multisensory integration across childhood and early adolescence: A high-density electrical
28
29 mapping study. *Cerebral Cortex*, 21, 1042-1055.

30
31 Briscoe, J., Bishop, D., & Norbury, C. (2001). Phonological processing, language, and literacy: A comparison of
32
33 children with mild-to-moderate sensorineural hearing loss and those with specific language impairment.
34
35 *Journal of Child Psychology and Psychiatry*, 42, 329-340.

36
37 Calvert, G., & Campbell, R. (2003). Reading speech from still and moving faces: The neural substrates of
38
39 visible speech. *Journal of Cognitive Neuroscience*, 15, 57-70.

40
41 Campbell, R. (2006). Audio-visual speech processing. In K. Brown, A. Anderson, L. Bauer, M. Berns, G. Hirst, &
42
43 J. Miller (Eds.), *The encyclopedia of language and linguistics* (pp. 562-569). Amsterdam: Elsevier.

44
45 Campbell, R., MacSweeney, M., Surguladze, S., Calvert, G., McGuire, P., Suckling, J., . . . David, A. (2001).
46
47 Cortical substrates for the perception of face actions: An fMRI study of the specificity of activation for seen
48
49 speech and for meaningless lower-face acts (gurning). *Cognitive Brain Research*, 12, 233-243.

50
51 Davis, C. & Kim, J. (2004). Audio-visual interactions with intact clearly audible speech. *The Quarterly Journal of*
52
53
54
55
56
57
58
59
60

VISUAL SPEECH, DETECTION, & ATTENTION

29

- 1
2
3 *Experimental Psychology*, 57A, 1103-1121.
4
5 Desjardins R, Rogers J, & Werker J. (1997). An exploration of why preschoolers perform differently than do
6
7 adults in audiovisual speech perception tasks. *Journal of Experimental Child Psychology*, 66:85–110.
8
9 Diederich, A., & Colonius, H. (2004). Bimodal and trimodal multisensory enhancement: Effects of stimulus
10
11 onset and intensity on reaction time. *Perception & Psychophysics*, 66, 1388-1404.
12
13 Dodd B, Mcintosh B, Erdener D, & Burnham D. (2008). Perception of the auditory-visual illusion in speech
14
15 perception by children with phonological disorders. *Clinical Linguistics & Phonetics*, 22:69–82
16
17 Dunn, L., & Dunn, D. (2007). *The Peabody Picture Vocabulary Test-IV* (Fourth ed.). Minneapolis, MN: NCS
18
19 Pearson, Inc.
20
21 Eimas, P. & Kavanagh, J. (1986). Otitis media, hearing loss, and child development: a NICHD conference
22
23 summary. *Public Health Reports*, 101, 289–293.
24
25 Erdener, D., & Burnham, D. (2013). The relationship between auditory-visual speech perception and
26
27 language-specific speech perception at the onset of reading instruction in English-speaking children
28
29 *Journal of Experimental Child Psychology*, 114, 120-138.
30
31 Evans J. (2002). Variability in comprehension strategy use in children with SLI: a dynamical systems
32
33 account. *International Journal of Language and Communication Disorders*, 37, 95–116.
34
35 Fiedler, K., Kutzner, F., & Krueger, J. (2012). The long way from α -error control to validity proper:
36
37 Problems with short-sighted false-positive debate. *Perspectives on Psychological Science*, 7, 661-669.
38
39 Files, B., Tjan, B., Jiang, J., & Bernstein, L. (2015). Visual speech discrimination and identification of
40
41 natural and synthetic consonant stimuli. *Frontiers in Psychology*, 6, 878.
42
43 <https://doi.org/10.3389/fpsyg.2015.00878>.
44
45 Fort, M., Spinelli, E., Savariaux, C., & Kandel, S. (2010). The word superiority effect in
46
47 audiovisual speech perception. *Speech Communication*, 52, 525-532.
48
49 Garner, W. (1974). *The processing of information and structure*. Potomac, MD: Lawrence Erlbaum.
50
51 Gilley, P., Sharma, A., Mitchell, T., & Dorman, M. (2010). The influence of a sensitive period for auditory-visual
52
53
54
55
56
57
58
59
60

VISUAL SPEECH, DETECTION, & ATTENTION

30

1
2
3 integration in children with cochlear implants. *Restorative Neurology and Neuroscience*, 28, 207-218.

4
5 Gogate L, Walker-Andrews A, & Bahrick L. (2001). The intersensory origins of word comprehension: an
6
7 ecological-dynamic systems view. *Developmental Science*, 4, 1–37.

8
9 Goodenough, F. (1935). The development of the reactive process from early childhood to maturity. *Journal of*
10
11 *Experimental Psychology*, 18, 431-450.

12
13 Grant, K. (2001). The effect of speechreading on masked detection thresholds for filtered speech. *Journal of*
14
15 *the Acoustical Society of America*, 109, 2272-2275.

16
17 Grant, K., & Seitz, P. (2000). The use of visible speech cues for improving auditory detection of spoken
18
19 sentences. *Journal of the Acoustical Society of America*, 108, 1197-1208.

20
21 Green, K., & Kuhl, P. (1989). The role of visual information in the processing of place and manner features in
22
23 speech perception. *Perception & Psychophysics*, 45, 34–42.

24
25 Harrar, V., Tammam, J., Perez-Bellido, A., Pitt, A., Stein, J., & Spence, C. (2014). Multisensory integration and
26
27 attention in developmental dyslexia. *Current Biology*, 24, 531-535.

28
29 Heathcote, A., Popiel, S., & Mewhort, D. (1991). Analysis of response time distributions: An example using
30
31 the Stroop task. *Psychological Bulletin*, 109, 340-347.

32
33 Hillock-Dunn, A., Grantham, D., & Wallace, M. (2016). The temporal binding window for audiovisual
34
35 speech: Children are like little adults. *Neuropsychologia*, 88, 74–82

36
37 Hnath-Chisolm, T., Laipply, E., & Boothroyd, A. (1998). Age-related changes on a children's
38
39 test of sensory-level speech perception capacity. *Journal of Speech, Language, and Hearing*
40
41 *Research*, 41, 94-106.

42
43 Jerger, S., Damian, M., McAlpine, R., & Abdi, H. (2018). Visual speech fills in both discrimination and
44
45 identification of non-intact auditory speech in children. *Journal of Child Language*, 45, 392-414.

46
47 Jerger, S., Damian, M., Parra, R., & Abdi, H. (2017). Visual speech alters the discrimination and identification
48
49 of non-intact auditory speech in children with hearing loss. *International Journal of Pediatric*
50
51 *Otorhinologyngology*, 94, 127-137.

VISUAL SPEECH, DETECTION, & ATTENTION

31

- 1
2
3 Jerger, S., Damian, M. F., Spence, M. J., Tye-Murray, N., & Abdi, H. (2009). Developmental shifts in children's
4 sensitivity to visual speech: A new multimodal picture-word task. *Journal of Experimental Child Psychology*,
5 102, 40-59.
6
7
8
9 Jerger, S., Damian, M., Tye-Murray, N., & Abdi, H. (2014). Children use visual speech to compensate for non-
10 intact auditory speech. *Journal of Experimental Child Psychology*, 126, 295-312.
11
12
13 Jerger, S., Damian, M., Tye-Murray, N., & Abdi, H. (2016). Phonological priming in children with hearing loss:
14 Effect of speech mode, fidelity, and lexical status. *Ear & Hearing*, 37, 623-633.
15
16
17
18 Jerger, S., Damian, M., Tye-Murray, N., & Abdi, H. (2017). Children perceive speech onsets by ear and eye.
19 *Journal of Child Language*, 44, 185-215.
20
21
22 Jerger, S., Damian, M., Tye-Murray, N., Dougherty, M., Mehta, J., & Spence, M. (2006). Effects of childhood
23 hearing loss on organization of semantic memory: Typicality and relatedness. *Ear & Hearing*, 27, 686-702.
24
25
26 Jerger, S., Martin, R., & Damian, M. (2002). Semantic and phonological influences on picture naming by
27 children and teenagers. *Journal of Memory and Language*, 47, 229-249.
28
29
30 Jerger, S., Martin, R., & Pirozzollo, F. (1988). A developmental study of the auditory Stroop effect. *Journal of*
31 *Brain and Language*, 35, 86-104.
32
33
34
35 Jiang, Y., Rouder, J., & Speckman, P. (2004). A note on the sampling properties of the Vincentizing (quantile
36 averaging) procedure. *Journal of Mathematical Psychology*, 48, 186-195.
37
38
39 Key, A., Gustafson, S., Rentmeester, L., Hornsby, B., & Bess, F. (2017). Speech-processing fatigue in children:
40 Auditory event-related potential and behavioral measures. *Journal of Speech, Language, and Hearing*
41 *Research*, 60, 2090-2104.
42
43
44
45 Kim, J., & Davis, C. (2003). Hearing foreign voices: Does knowing what is said affect visual-masked-speech
46 detection. *Perception*, 32, 111-120.
47
48
49
50 Kim, J., & Davis, C. (2004). Investigating the audio-visual speech detection advantage. *Speech Communication*,
51 44, 19-30.
52
53
54 Krauss, R., Freyberg, R., & Morsella, E. (2002). Inferring speakers' physical attributes from their voices. *Journal*
55 *of Experimental Social Psychology*, 38, 618-625.
56
57
58
59
60

VISUAL SPEECH, DETECTION, & ATTENTION

32

- 1
2
3 Kuhl, P., & Meltzoff, A. (1982). The bimodal perception of speech in infancy. *Science*, 218, 1138–1141.
- 4
5 Lachs, L., & Pisoni, D. (2004). Crossmodal source identification in speech perception. *Ecological Psychology*,
6
7 16, 159–187.
- 8
9 Lalonde, K., & Holt, R. (2016). Audiovisual speech perception development at varying levels of perceptual
10
11 processing. *Journal of the Acoustical Society of America*, 139, 1713–1723.
- 12
13 Lalonde, K., & Holt, R. (2015). Preschoolers benefit from visually salient speech cues. *Journal*
14
15 *of Speech, Language, and Hearing Research*, 58, 135-150
- 16
17
18 Langner, R., & Eickhoff, S. (2013). Sustaining attention to simple tasks: A meta-analytic review of the neural
19
20 mechanisms of vigilant attention. *Psychological Bulletin*, 139, 870–900.
- 21
22 Laurienti, P., Burdette, J., Maldjian, J., & Wallace, M. (2006). Enhanced multisensory integration in older
23
24 adults. *Neurobiology of Aging*, 27, 1155-1163.
- 25
26
27 Lewis, F., Reeve, R., Kelly, S., & Johnson, K. (2017). Evidence of substantial development of inhibitory control
28
29 and sustained attention between 6 and 8 years of age on an unpredictable go / no-go task. *Journal of*
30
31 *Experimental Child Psychology*, 157, 66-80.
- 32
33
34 Lewkowicz, D. (2000). Infants' perception of the audible, visible, and bimodal attributes of multimodal
35
36 syllables. *Child Development*, 71, 1241–1257.
- 37
38 Libertus K., Landa R., and Haworth J. (2017) Development of attention to faces during the first 3 years:
39
40 Influences of stimulus type. *Frontiers in Psychology*, 8:1976. <https://doi.org/10.3389/fpsyg.2017.01976>
- 41
42 Macken, W., Phelps, F., & Jones, D. (2009). What causes auditory distraction? *Psychonomic Bulletin & Review*,
43
44 16, 139-144.
- 45
46 Manly, T., Anderson, V., Nimmo-Smith, I., Turner, A., Watson, P., & Robertson, I. (2001). The differential
47
48 assessment of children's attention: *The Test of Everyday Attention for Children (TEA-Ch)*, normative sample
49
50 and ADHD performance. *Journal of Psychology & Psychiatry*, 42, 1065-1081.
- 51
52
53 Marks, L. (2004). Cross-modal interactions in speeded classification. In G. Calvert, C. Spence, & B. Stein(Eds.),
54
55 *The handbook of multisensory processes* (pp. 85–105). Cambridge, MA: MIT Press.
- 56
57
58
59
60

VISUAL SPEECH, DETECTION, & ATTENTION

33

- 1
2
3 Marslen-Wilson, W., & Zwitserlood, P. (1989). Accessing spoken words: The importance of word onsets.
4
5 *Journal of Experimental Psychology: Human Perception and Performance*, 15, 576-585.
6
7 Mavica, L., & Barenholtz, E. (2013). Matching voice and face identity from static images. *Journal of*
8
9 *Experimental Psychology: Human Perception and Performance*, 39, 307–312.
10
11 McClelland, J., & Elman, J. (1986). The TRACE model of speech perception. *Cognitive Psychology*, 18, 1-86.
12
13 McConachie, H., & Moore, V. (1994). Early expressive language of severely visually impaired children.
14
15 *Developmental Medicine & Child Neurology*, 36, 230-240.
16
17 Miller, J. (1988). A warning about median reaction time. *Journal of Experimental Psychology: Human*
18
19 *Perception & Performance*, 14, 539-543.
20
21 Miller, J., & Ulrich, R. (2003). Simple reaction time and statistical facilitation: A parallel grains model. *Cognitive*
22
23 *Psychology*, 46, 101-151.
24
25 Molholm, S., Ritter, W., Murray, M., Javitt, D., Schroeder, C., & Foxe, J. (2002). Multisensory auditory-visual
26
27 interactions during early sensory processing in humans: A high-density electrical mapping study. *Cognitive*
28
29 *Brain Research*, 14, 115-128.
30
31 Mordkoff, T., & Yantis, S. (1993). Dividing attention between color and shape: Evidence of coactivation.
32
33 *Perception & Psychophysics*, 53, 357-366.
34
35 Munhall K, Gribble P, Sacco L, & Ward M. (1996). Temporal constraints on the McGurk effect. *Perception*
36
37 *& Psychophysics*, 58:351–362.
38
39 Napolitano, A., & Sloutsky, V. (2004). Is a picture worth a thousand words? The flexible nature of modality
40
41 dominance in young children. *Child Development*, 75, 1850 – 1870.
42
43 O'Toole, A., Roak, D., & Abdi, H. (2002). Recognizing moving faces: A psychological and neural synthesis.
44
45 *Trends in Cognitive Sciences*, 6, 261-266.
46
47 Otsuka, Y., Konishi, Y., Kanazawa, S., Yamaguchi, M., Abdi, H., & O'Toole, A. (2009). Recognition of moving and
48
49 static faces by young infants. *Child Development*, 80, 1259-1271.
50
51 Rabbitt, P., & Goward, L. (1994). Age, information processing speed, and intelligence. *The Quarterly Journal of*
52
53
54
55
56
57
58
59
60

VISUAL SPEECH, DETECTION, & ATTENTION

34

1
2
3 *Experimental Psychology Section A*, 47, 741-760.

4
5 Rabbitt, P., Osman, P., Moore, B., & Stollery, B. (2001). There are stable individual differences in performance
6
7 variability, both from moment to moment and from day to day. *The Quarterly Journal of Experimental*
8
9 *Psychology Section A*, 54, 981-1003.

10
11 Ratcliff, R. (1979). Group reaction time distributions and an analysis of distribution statistics. *Psychological*
12
13 *Bulletin*, 86, 446-461.

14
15 Reinvang, I. (1998). Validation of reaction time in continuous performance tasks as an index of attention by
16
17 electrophysiological measures. *Journal of Clinical and Experimental Neuropsychology*, 20, 885-897.

18
19 Reisberg, D., McLean, J., & Goldfield, A. (1987). Easy to hear but hard to understand: A lip-reading advantage
20
21 with intact auditory stimuli. In B. Dodd & R. Campbell (Ed.s), *Hearing by eye: The psychology of lip-reading*
22
23 (pp. 97-111). Hillsdale, NJ: Lawrence Erlbaum Associates, Inc.

24
25 Robertson, I., Manly, T., Andrade, J., Baddeley, B., & Yiend, J. (1997). 'Oops!': Performance correlates of
26
27 everyday attentional failures in traumatic brain injured & normal subjects. *Neuropsychologia*, 35, 747-758.

28
29 Rosenthal, R., Rosnow, R., & Rubin, D. (2000). *Contrasts and effect sizes in behavioral research. A*
30
31 *correlational approach*. New York: Cambridge University Press.

32
33
34
35 Ross, M., & Lerman, J. (1971). *Word Intelligibility by Picture Identification*. Pittsburgh: Stanwix House, Inc.

36
37 Schroger, E., & Widmann, A. (1998). Speeded responses to audiovisual signal changes result from bimodal
38
39 integration. *Psychophysiology*, 35, 755-759.

40
41 Schwartz, J., Berthommier, F., & Savariaux, C. (2004). Seeing to hear better: Evidence for early audio-visual
42
43 interactions in speech identification. *Cognition*, 93, B69-78.

44
45 Sloutsky, V., & Napolitano, A. (2003). Is a picture worth a thousand words? Preference for auditory modality
46
47 in young children. *Child Development*, 74, 822-833.

48
49 Smith, H., Dunn, A., Baguley, T., & Stacey, P. (2016a). Concordant cues in faces and voices: Testing the back-
50
51 up signal hypothesis. *Evolutionary Psychology*, 1-10, <https://doi.org/10.1177/1474704916630317>.

52
53
54 Smith, H., Dunn, A., Baguley, T., & Stacey, P. (2016b). Matching novel face and voice identity using static and
55
56
57
58
59
60

VISUAL SPEECH, DETECTION, & ATTENTION

35

dynamic facial images. *Attention, Perception & Psychophysics*, 78, 868-879.

Smith, L., & Thelen, E. (2003). Development as a dynamic system. *Trends in Cognitive Sciences*, 7, 343-348.

Stevenson, R., Ghose, D., Fister, J., Sarko, D., Altieri, N., Nidiffer, A.,...Wallace, M. (2014). Identifying and quantifying multisensory integration: A tutorial review. *Brain Topography*, 27, 707-730.

Stevenson, R., Sheffield, S., Butera, I., Gifford, R. & Wallace, M. (2017). Multisensory integration in cochlear implant recipients. *Ear and Hearing*, 38, 521–538

Teinonen, T., Aslin, R., Alku, P., & Csibra, G. (2008). Visual speech contributes to phonetic learning in 6-month-old infants. *Cognition*, 108, 850-855.

Thillay, A., Roux, S., Gissot, V., Carteau-Martin, I., Knight, R., Bonnet-Brilhault, F., & Bidet-Caulet, A. (2015). Sustained attention and prediction: Distinct brain maturation trajectories during adolescence. *Frontiers in Human Neuroscience*, 9(Article 519), <https://doi.org/10.3389/fnhum.2015.00519>.

Tjan, B., Chao, E., & Bernstein, L. (2013). A visual or tactile signal makes auditory speech detection more efficient by reducing uncertainty. *European Journal of Neuroscience*, 39, 1323 - 1331.

Tse, C., Balota, D., Yap, M., Duchek, J., & McCabe, D. (2010). Effects of healthy aging and early stage dementia of the Alzheimer's type on components of response time distributions in three attention tasks. *Neuropsychology*, 24, 300-315.

Tye-Murray, N., Spehar, B., Myerson, J., Sommers, M., & Hale, S. (2011). Cross-modal enhancement of speech detection in young and older adults: Does signal content matter? *Ear & Hearing*, 32, 650-655.

van Wassenhove, V., Grant, K., & Poeppel, D. (2007). Temporal window of integration in auditory-visual speech perception. *Neuropsychologia*, 45, 598-607

van Wassenhove, V., Grant, K., & Poeppel, D. (2005). Visual speech speeds up the neural processing of auditory speech. *Proceedings National Academy of Sciences of United States of America*, 102, 1181-1186.

Venker, C., Goodwin, J., Roe, D., Kaemingk, K., Mulvaney, S., & Quan, S. (2007). Normative psychomotor vigilance task performance in children ages 6 to 11—*The Tucson Children's Assessment of Sleep Apnea (TuCASA)*. *Sleep Breath*, 11, 217-224.

VISUAL SPEECH, DETECTION, & ATTENTION

36

- 1
2
3 Vickers, J. (2007). *Perception, cognition, and decision training: The quiet eye in action* (pp. 47-64). Champaign,
4 IL: Human Kinetics.
- 5
6
7 Watkins, S., Dalton, P., Lavie, N., & Rees, G. (2007). Brain mechanisms mediating auditory attentional capture
8 in humans. *Cerebral Cortex*, 17, 1694-1700
- 9
10
11 Weissman, D., Roberts, K., Visscher, K., & Woldorff, M. (2006). The neural bases of momentary lapses in
12 attention. *Nature Neuroscience*, 9, 971-978.
- 13
14
15
16 Whelan, R. (2008). Effective analysis of reaction time data. *The Psychological Record*, 58, 475-482.
- 17
18 Wickens, C. (1974). Temporal limits of human information processing: A developmental study. *Psychological*
19 *Bulletin*, 81, 739-755.
- 20
21
22 Woodworth, R., & Schlosberg, H. (1954). *Experimental psychology*. New York: Henry Holt and Company.
- 23
24 Zhou, B., & Krott, A. (2016). Bilingualism enhances attentional control in non-verbal conflict tasks – evidence
25 from ex-Gaussian analyses. *Bilingualism: Language & Cognition*,
26
27 <https://doi.org/10.1017/S1366728916000869>.
- 28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Footnotes

Footnote 1. A motor (key-press) component is also involved in the task, but it is assumed to be approximately constant within individuals and is not considered (e.g., Miller & Ulrich, 2003).

Footnote 2. We thank one of the reviewers for recommending this analysis.

For Peer Review

Figure Legends

Figure 1. Mean response times in the A, V, and AV modes for the static and dynamic faces in the four age groups and in all participants. The error bars are ± 1 standard error of the mean.

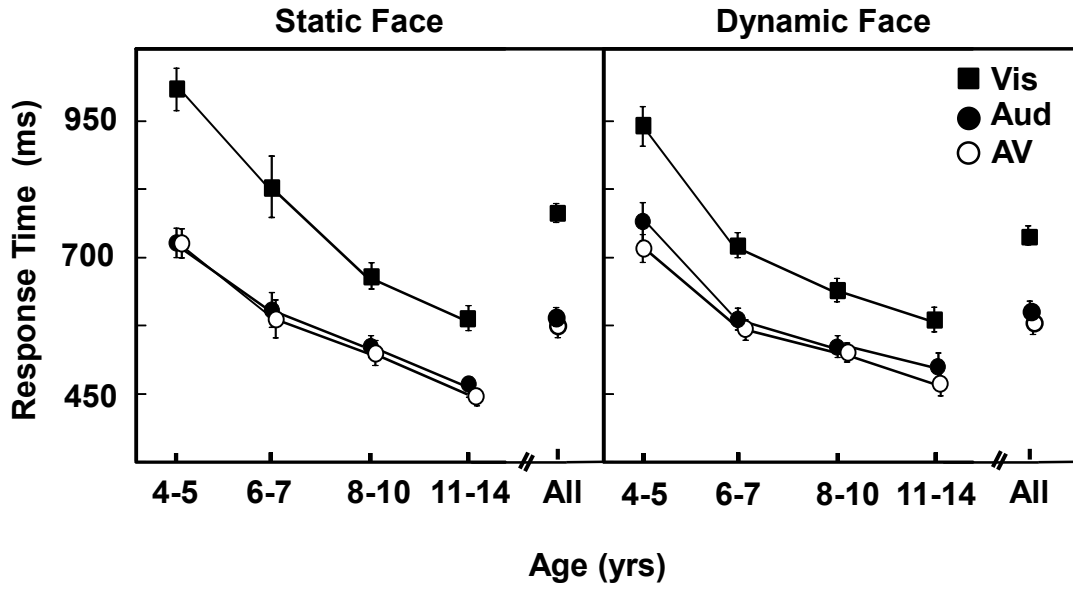
Figure 2. The mean difference scores (V response times – A response times) in the age groups and in all participants for the static and dynamic faces at the 1st/faster and 3rd/slower quartiles of the CDFs. The error bars are ± 1 standard error of the mean. Every data point showed a significant difference for the V vs. A modes. An asterisk indicates the data points showing a significant difference for the static vs. dynamic silent faces.

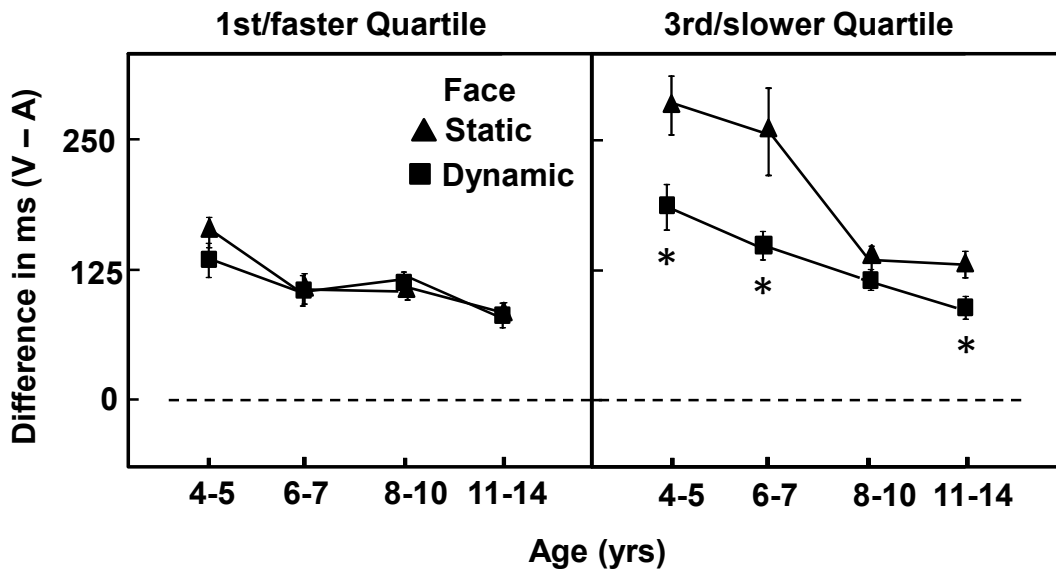
Figure 3. The mean difference scores (AV response times – A response times) in the age groups and in all participants for the static and dynamic faces at the 1st/faster and 3rd/slower quartiles of the CDFs. The error bars are ± 1 standard error of the mean. A star indicates the data points showing a significant difference for the AV vs. A modes; an asterisk indicates the data point showing a significant difference for the static vs. dynamic faces.

Figure Legends: Appendix

Figure 1App. The cumulative distribution functions (CDFs) for the A, AV, and V modes in the static (1App_a) and dynamic (1App_b) facial conditions for all age groups.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56





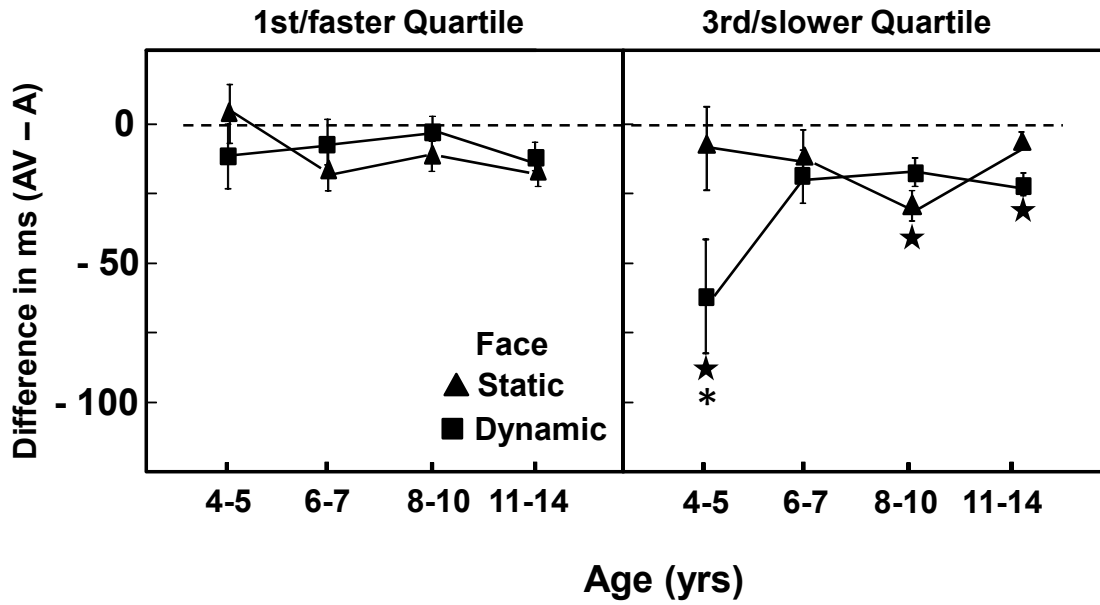


Table 1

Results of mixed-design analysis of variance (ANOVA)

<i>Factors</i>	<i>Mean Square Error</i>	<i>F value</i>	<i>p value</i>	<i>partial η^2</i>
Age Group	.040	34.80	< .0001	.485
Mode	.002	524.76	< .0001	.825
Facial Condition	.005	1.11	ns	.011
Mode x Age Group	.002	1.95	ns	.051
Facial Condition x Age Group	.005	0.89	ns	.023
Mode x Facial Condition	.001	15.11	< .0001	.121
Mode x Facial Condition x Age Group	.001	2.02	ns	.050

Note: The ANOVA contained one between-participant factor (Age Group: 4–5-yrs, 6–7-yrs, 8–10-yrs, and 11–14-yrs) and two within-participant factors (Mode: V, A, and AV; Facial Condition: static vs. dynamic).

The dependent variable was the log transformed response times. The degrees of freedom were 3,111 for age group and facial condition x age group; 1, 111 for facial condition; 2,222 for mode and facial condition x mode; and 6,222 for mode x age group and facial condition x mode x age group.

ns = not significant.

Initially we conducted this analysis with Gender as a factor, but Gender did not influenced the results.

Thus, Gender was eliminated.

Table 2

Results of mixed-design analysis of variance (ANOVA)

<i>Factors</i>	<i>Mean Square Error</i>	<i>F value</i>	<i>p value</i>	<i>partial η^2</i>
Age Group	.028	31.14	< .0001	.462
Mode	.001	29.51	< .0001	.210
Facial Condition	.004	0.47	ns	.005
Mode x Age Group	.001	0.59	ns	.012
Facial Condition x Age Group	.004	0.68	ns	.018
Mode x Facial Condition	.001	3.85	ns	.034
Mode x Facial Condition x Age Group	.001	3.47	.018	.086

Note: The ANOVA contained one between-participant factor (Age Group: 4–5-yrs, 6–7-yrs, 8–10-yrs, and 11–14-yrs) and two within-participant factors (Mode: A vs. AV; Facial Condition: static vs. dynamic). The dependent variable was the log transformed reaction times at the 3rd/slower quartile. The degrees of freedoms were 3,111 for age group, mode x age group, facial condition x age group, and mode x facial condition x age group; and 1, 111 for mode, facial condition, and mode x facial condition.

ns = not significant.

Table 3

Results of paired *t*-tests in each age group for each facial condition

<i>Facial Condition</i>	<i>t value</i>	<i>p value</i>	<i>partial η^2</i>
4–5-yrs			
Static Face	0.38	ns	.001
Dynamic Face	3.19	.003	.246
6–7-yrs			
Static Face	1.42	ns	.053
Dynamic Face	1.50	ns	.091
8–10-yrs			
Static Face	4.69	<.0001	.421
Dynamic Face	2.44	ns	.154
11–14-yrs			
Static Face	0.63	ns	.015
Dynamic Face	3.28	.003	.294

Note. The dependent variable was the log transformed response times at the 3rd/slower quartile for the AV vs. A modes.

ns = not significant.

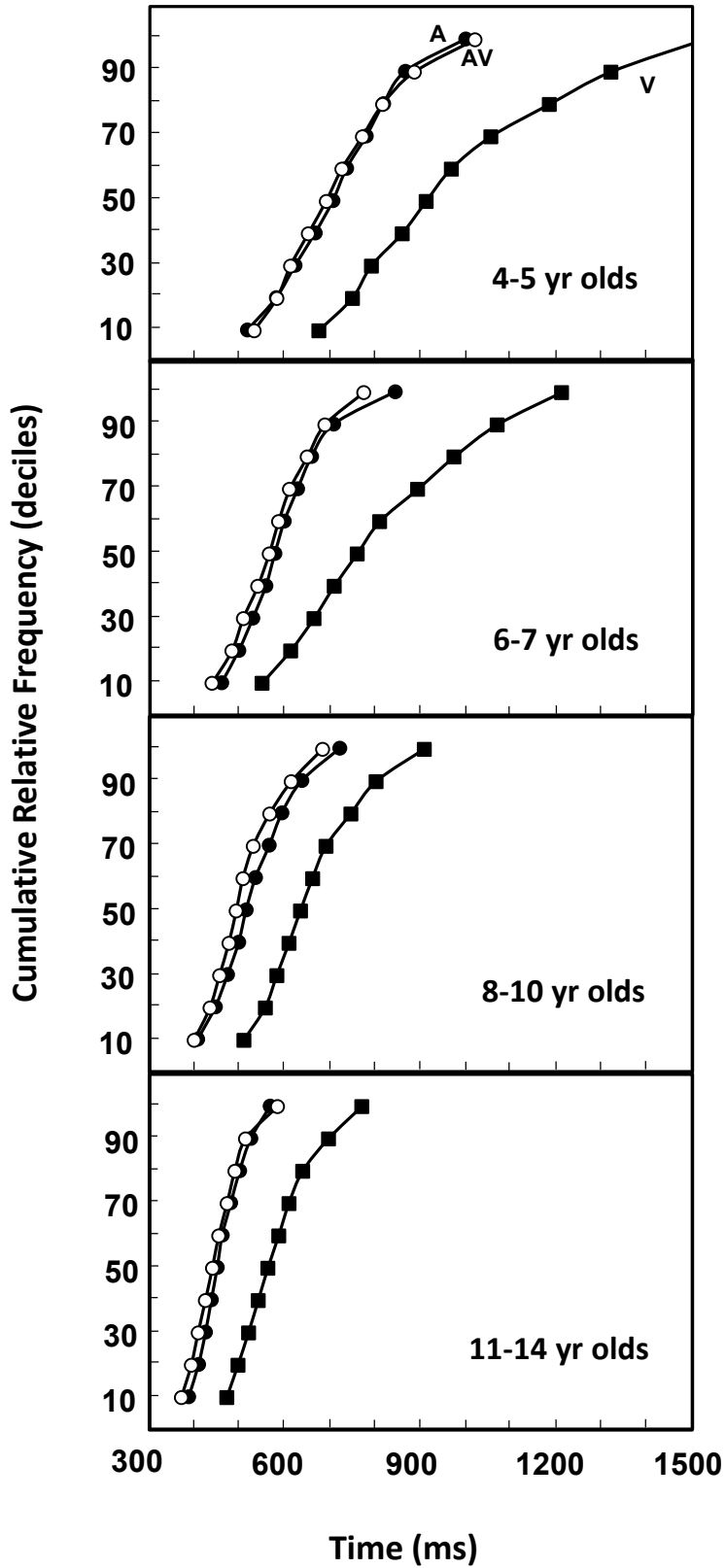


Figure 1App_a. Static Face

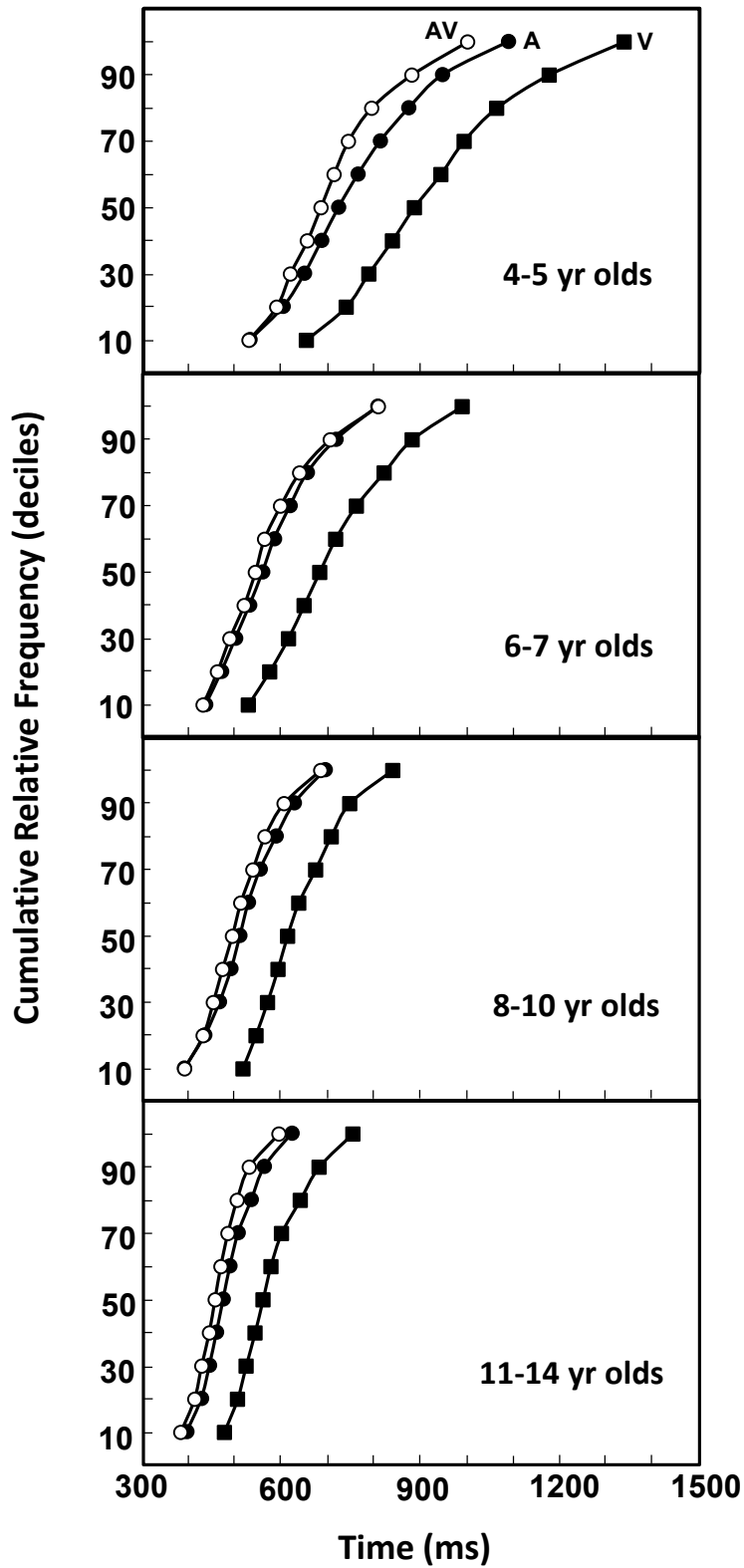


Figure 1App_b. Dynamic Face

Appendix

Table A1

F_{contrast} analyses to determine whether the V vs. A response times differ at each quartile for each facial condition in the age groups.

Quartile Facial Condition	Mode		F contrast	p value	partial η^2
	V	A			
4–5-yrs					
1st (fast) Quartile					
Static Face	726	561	354.23	<.0001	.761
Dynamic Face	711	580	228.82	<.0001	.673
3rd (slow) Quartile					
Static Face	1046	765	507.09	<.0001	.820
Dynamic Face	983	800	233.88	<.0001	.678
6–7-yrs					
1st (fast) Quartile					
Static Face	594	485	228.82	<.0001	.673
Dynamic Face	562	460	208.76	<.0001	.653
3rd (slow) Quartile					
Static Face	871	617	537.67	<.0001	.829
Dynamic Face	753	608	254.91	<.0001	.697
8–10-yrs					
1st (fast) Quartile					
Static Face	541	434	276.72	<.0001	.714
Dynamic Face	537	421	341.76	<.0001	.755
3rd (slow) Quartile					
Static Face	687	554	244.20	<.0001	.688
Dynamic Face	664	548	204.08	<.0001	.648
11–14-yrs					
1st (fast) Quartile					
Static Face	489	401	218.69	<.0001	.663
Dynamic Face	496	415	199.22	<.0001	.642
3rd (slow) Quartile					
Static Face	602	474	305.35	<.0001	.733
Dynamic Face	596	502	180.72	<.0001	.619

Note: Results were based on a mixed-design analysis of variance with one between-participant factor (Age Group: 4–5-yrs, 6–7-yrs, 8–10-yrs, and 11–14-yrs) and three within-participant factors (Mode: V vs. A; Facial Condition: static vs. dynamic; Quartile: 1st vs. 3rd). Although mean response times are presented to ease understanding, the dependent variable for analyses was the log transformed response times. For all $F_{\text{contrasts}}$, the mean square error = .0005 and the degrees of freedom = 1,111

Appendix

Table A2

$F_{contrast}$ analyses to determine whether the $V - A$ difference scores for the static vs. dynamic facial conditions differ at each quartile in the age groups.

Quartile	Facial Condition		F contrast	p value	partial η^2
	Stat	Dynam			
4–5-yrs					
1st (fast) Quartile	165	131	6.04	ns	.052
3rd (slow) Quartile	281	183	26.39	<.0001	.192
6–7-yrs					
1st (fast) Quartile	109	102	0.10	ns	.001
3rd (slow) Quartile	254	145	25.22	<.0001	.185
8–10-yrs					
1st (fast) Quartile	107	116	1.36	ns	.012
3rd (slow) Quartile	133	116	0.88	ns	.008
11–14-yrs					
1st (fast) Quartile	88	81	0.29	ns	.003
3rd (slow) Quartile	128	94	199.22	.006	.066

Note: Results were based on a mixed-design analysis of variance with one between-participant factor (Age Group: 4–5-yrs, 6–7-yrs, 8–10-yrs, and 11–14-yrs) and two within-participant factors (Facial Condition: static vs. dynamic; Quartile: 1st vs. 3rd). Although mean difference scores ($V - A$) are presented to ease understanding, the dependent variable for analyses was always the log transformed difference scores. For all $F_{contrasts}$, the mean square error = .0010 and the degrees of freedom = 1,111 ns = not significant; Stat = static; Dynam = dynamic

Appendix

Table A3

F_{contrast} analyses to determine whether the AV vs. A response times differ at each quartile for each facial condition in age groups.

<i>Quartile Facial Condition</i>	<i>Mode</i>		<i>F contrast</i>	<i>p value</i>	<i>partial η^2</i>
	<i>AV</i>	<i>A</i>			
4–5-yrs					
1st (fast) Quartile					
Static Face	566	561	0.20	ns	.002
Dynamic Face	568	580	1.01	ns	.009
3rd (slow) Quartile					
Static Face	758	765	0.40	ns	.004
Dynamic Face	737	800	26.46	<.0001	.192
6–7-yrs					
1st (fast) Quartile					
Static Face	468	485	7.20	ns	.061
Dynamic Face	452	460	1.41	ns	.012
3rd (slow) Quartile					
Static Face	605	617	3.63	ns	.032
Dynamic Face	589	608	4.24	ns	.037
8–10-yrs					
1st (fast) Quartile					
Static Face	423	434	3.63	ns	.032
Dynamic Face	418	421	2.83	ns	.025
3rd (slow) Quartile					
Static Face	525	554	15.55	.0001	.123
Dynamic Face	531	548	4.24	ns	.095
11–14-yrs					
1st (fast) Quartile					
Static Face	386	401	6.66	ns	.057
Dynamic Face	402	415	4.24	ns	.037
3rd (slow) Quartile					
Static Face	466	474	0.40	ns	.004
Dynamic Face	480	502	10.51	.002	.086

Note: Results were based on a mixed-design analysis of variance (ANOVA) with one between-participant factor (Age Group: 4–5-yrs, 6–7-yrs, 8–10-yrs, and 11–14-yrs) and three within-participant factors (Mode: AV vs. A; Facial Condition: static vs. dynamic; Quartile: 1st vs. 3rd). Although mean response times are presented to ease understanding, the dependent variable for analyses was the log transformed response times. For all $F_{\text{contrasts}}$, the mean square error = .0005 and the degrees of freedom = 1,111. ns = not significant

Appendix

Table A4

$F_{contrast}$ analyses to determine whether the AV – A difference scores for the static vs. dynamic facial conditions differ at each quartile in the age groups.

Quartile	Facial Condition		F contrast	p value	partial η^2
	Stat	Dynam			
4–5-yrs					
1st (fast) Quartile	05	– 12	0.92	ns	.008
3rd (slow) Quartile	– 07	– 63	9.76	.002	.081
6–7-yrs					
1st (fast) Quartile	– 17	– 08	0.91	ns	.008
3rd (slow) Quartile	– 12	– 19	0.01	ns	.000
8–10-yrs					
1st (fast) Quartile	– 11	– 03	0.91	ns	.008
3rd (slow) Quartile	– 28	– 17	1.72	ns	.015
11–14-yrs					
1st (fast) Quartile	– 15	– 13	0.40	ns	.004
3rd (slow) Quartile	– 08	– 22	2.83	ns	.025

Note: Results were based on a mixed-design analysis of variance with one between-participant factor (Age Group: 4–5-yrs, 6–7-yrs, 8–10-yrs, and 11–14-yrs) and two within-participant factors (Facial Condition: static vs. dynamic; Quartile: 1st vs. 3rd). Although mean difference scores (AV – A) are presented to ease understanding, the dependent variable for analyses was always the log transformed difference scores. For all $F_{contrasts}$, the mean square error = .0010 and the degrees of freedom = 1,111. ns = not significant; Stat = static; Dynam = dynamic