

Published in final edited form as:

J Exp Child Psychol. 2014 October ; 126: 295–312. doi:10.1016/j.jecp.2014.05.003.

Children Use Visual Speech to Compensate for Non-Intact Auditory Speech

Susan Jerger, Ph.D.^{*}, Markus F. Damian, Ph.D.⁺, Nancy Tye-Murray, Ph.D.[^], and Hervé Abdi, Ph.D.^{*}

Susan Jerger: sjerger@utdallas.edu; Markus F. Damian: m.damian@bristol.ac.uk; Nancy Tye-Murray: nmurray@wustl.edu; Hervé Abdi: herve@utdallas.edu

^{*}School of Behavioral & Brain Sciences, Univ. of Texas at Dallas, 800 W. Campbell Rd. Richardson, TX 75080

⁺School of Experimental Psychology, University of Bristol, 12a Priory Road, Room 1D20, Bristol BS8 1TU, UK

[^]Dept. Otolaryng., Washington Univ. School of Medicine, Box 8115, 660 S. Euclid Ave., St. Louis, MO 63110

When adults engage in casual conversations in noisy environments, they typically understand each other without effort. Such skilled understanding in degraded soundscapes seems related to the inherently multimodal nature of speech perception as dramatically illustrated by the classic McGurk effect (McGurk & MacDonald, 1976). In this task, an audiovisual speech stimulus with mismatched auditory and visual onsets (e.g., hearing /ba/ while seeing /ga/) is presented to participants. Adults typically perceive a mixture of the auditory and visual inputs (i.e., /da/ or /8a/). Our ability to integrate auditory and visual speech helps us understand speech in noisy environments as well as unfamiliar content in clear environments (Arnold & Hill, 2001; MacLeod & Summerfield, 1987). Because visual speech is so useful to communication when the message is complex/degraded or the environment noisy (i.e., classrooms), it is paramount to investigate the development of multimodal speech perception during the preschool-elementary school years.

Development of Multimodal Speech Perception

Extant studies with multiple types of tasks report that—compared to adults—children from about 5 yrs to the pre-teen/teenage yrs show reduced sensitivity to visual speech (e.g., McGurk & MacDonald, 1976; Desjardins, Rogers, & Werker, 1997; Erdener & Burnham, 2013; Jerger, Damian, Spence, Tye-Murray, & Abdi, 2009; Ross et al., 2011; Sekiyama & Burnham, 2008; Tremblay et al., 2007; see Fort, Spinelli, Savariaux, & Kandel, 2012, for an exception re: vowel monitoring). For example, McGurk and MacDonald (1976) noted that

© 2014 Elsevier Inc. All rights reserved.

Corresponding author: Susan Jerger, Ph.D., School of Behavioral & Brain Sciences, University of Texas-Dallas, 800 W. Campbell Rd., Richardson, Tx 75080, telephone: 214-236-6972, FAX: 972-883-2491, sjerger@utdallas.edu.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final citable form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

children were influenced by visual speech only about half as often as adults. Poorer sensitivity to visual speech in children has been attributed to developmental differences in articulatory proficiency or speechreading skills (Desjardins et al., 1997; Erdener & Burnham, 2013), linguistic experiences and perceptual tuning into language-specific phonemes (Erdener & Burnham, 2013; Sekiyama & Burnham, 2008), or to differences in the perceptual weighting of visual speech cues (Green, 1998).

Supplementing these more specific theories, Jerger et al., (2009) —who observed a U-shaped developmental function with a significant influence of congruent visual speech on phonological priming in 4-yr-olds and 12-yr-olds, but not in 5–9-yr-olds—adopted a dynamic systems viewpoint (Smith & Thelan, 2003). Jerger et al. hypothesized that children's poorer sensitivity to visual speech from 5–9-yr-olds was reflecting a period of dynamic growth as relevant perceptual, linguistic, and cognitive skills were reorganizing in response to external and internal factors. Externally, reorganization may be due to literacy instruction at about 5 to 6 yrs during which time knowledge transmutes from phonemes as coarticulated nondistinct parts of speech into phonemes as separable distinct written and heard elements (Bryant, 1995; Burnham, 2003; Horlyck, Reid, & Burnham, 2012). Internally, reorganization may be due to phonological processes becoming sufficiently proficient (at around the same ages) to support the use of inner speech for learning, remembering, and problem solving (Conrad, 1971). The actuality of reorganization is confirmed by evoked potential studies indicating—during this age period—developmental restructuring of the lexical phonological system (Bonte & Blomert, 2004).

A dynamic systems viewpoint (Smith & Thelan, 2003) stresses two points motivating this research. First, periods of poorer sensitivity to visual speech may reflect a transition period—in contrast to a loss of ability—during which time relevant perceptual, linguistic, and cognitive resources are harder to access. Second, dynamic periods of reorganization and growth are characterized by less robust processing systems and decreases in processing efficiencies. So according to a dynamic system viewpoint, the influence of visual speech may vary as a function of task/stimulus demands for these softly reassembled resources. This suggests that tasks with less complex stimuli and/or lower task demands may not tax a child's limited processing resources as readily. This would make the harder to access resources more accessible, and therefore the child's performance might reveal greater sensitivity to visual speech. Below we introduce our new task to set up reviewing the related literature.

The New Visual Speech Fill-In Effect

Our approach assesses performance for words and nonwords with intact visual speech coupled to non-intact auditory speech (excised consonant onsets). The strategy was to insert visual speech into the gap created by the excised auditory onset to study the possibility of a visual speech fill-in effect—operationally defined by the difference in performance between the audiovisual (AV) mode and auditory only (AO) modes. As an example, stimuli for the word bag would be: 1) AV: intact visual (/b/ag) coupled to non-intact auditory (/–b/ag) and 2) AO: static face coupled to the same non-intact auditory (/–b/ag). The visual speech fill-in effect occurs when listeners experience hearing the *same* auditory stimulus as *different*

depending on the presence/absence of visual speech, such as hearing /bag/ in the AV mode but /ag/ in the AO mode. The AO mode also controls for the influence on performance of remaining coarticulatory cues in the stimulus. Below we review related studies that investigated the perception of AV or AO speech containing an excised segment-replaced with noise, however, instead of visual speech.

Previous Studies

Studies with adults indicate that listeners report hearing AO speech with an excised segment replaced with noise as intact, a phenomenon that has been called illusory filling-in, illusory continuity, auditory induction, and perceptual or phonemic restoration (Bashford & Warren, 1987; Samuel, 1981; Shahin, Bishop, & Miller, 2009; Warren, 1970). Of particular interest to the current approach—inserting visual speech into the auditory gap—three of the studies questioned whether the addition of visual speech enhances this type of illusory phenomenon. The studies investigated adults' ability to discriminate between AV speech with 1) an excised auditory phoneme replaced with noise vs 2) an intact auditory phoneme with superimposed noise. Two studies (Shahin & Miller, 2009; Shahin, Kerlin, Bhat, & Miller, 2012) varied the auditory gap duration to determine whether AV congruent speech lengthened the gap duration producing illusory continuity relative to auditory plus incongruent or static visual speech. Results showed that participants perceived the stimulus as continuous—rather than interrupted—over longer gaps when they saw and heard congruent speech compared to static or incongruent speech. These investigators reasoned that congruent visual speech or visual context may have made the noise sound more speech-like or may have helped restore the speech envelope during the interruption, thus enhancing the illusion of continuity. The other study (Trout & Poser, 1990) investigated whether—relative to AO speech-AV speech enhanced the ability to discriminate the noise-replacing vs -noise-superimposed types of stimuli when the place of articulation was easy to see (e.g., bilabial) vs difficult to see (e.g., alveolar and velar). Although results were difficult to interpret unequivocally, visual speech did not appear to influence the results. In short, the previous adult studies suggest opposing conclusions about whether visual speech influences this type of illusory phenomenon.

In contrast to the above adult studies, studies in children have focused on AO speech. Results have shown that-compared to adults-recognition of AO speech with excised phonemes replaced by noise is unusually disrupted in 5-yr-olds (Walley, 1988; Newman, 2004). For Walley, this indicated that children require more intact AO speech to identify even familiar words, although Walley also cautioned that the noise may have influenced performance. Newman (2004) extended her research paradigm and included a comparison of AO speech with excised phonemes replaced by noise vs silence. The children showed an adult-like advantage when noise filled the gap, and so Newman concluded that—even though children have remarkable difficulty understanding non-intact AO speech—children show adult-like perceptual restoration. Overall these results make it difficult to draw strong conclusions about this type of illusory phenomenon in children.

In overview, the literature does not provide clear-cut predictions about performance by children on the new visual speech fill-in task. The child literature on the development of AV

speech perception, however, clearly indicates that the influence of visual speech is not as consistent and robust in children as in adults. To account for this difference, Jerger et al., (2009) theorized that the influence of visual speech on performance in children might be modulated by the information processing requirements of tasks. The present research reports two projects that systematically manipulated aspects of information processing. Study 1 assessed the effects of age, salience of visual speech cues, and lexical status on the visual speech fill-in effect. Study 2 studied the effects of task demands and of child factors on visually influenced performance by comparing results on the visual speech fill-in and McGurk tasks. Below we address some relevant issues for selecting test stimuli.

Study 1

Stimuli Construction Issues

Speech in Noise—To reduce the ceiling effect for AO speech, previous AV studies have studied AO speech in noise. The ability to identify AO speech in noise, however, does not reach adult-like performance until about 11–14 yrs for consonants and 10 yrs for vowels (Johnson, 2000). Thus we reasoned that noise might have diminished the effect of visual speech on previous tasks because the children had difficulty inhibiting task-irrelevant input and resisting interference (Bjorklund & Harnishfeger, 1995). We therefore chose to reduce the ceiling effect for AO performance with non-intact speech.

Salience of Speech—Our task focuses on onsets to evaluate whether children's performance can benefit from visual speech. Relative to the other parts of an utterance, onsets are easier to speechread, more reliable with less articulatory variability, and more stressed (Gow, Melvold, & Manuel, 1996). Phonemes are also more easily processed in the onset position; onsets reveal better accuracy on phonological awareness tasks in children (Treiman & Zukowski, 1996) and nonword tasks in adults (Gupta, Lipinski, Abbs, & Lin, 2005). Finally, onsets are important because speech perception proceeds incrementally even in infants (Fernald, Swingly, & Pinto, 2001). These effects are compatible with theories proposing that speech input activates phonological and lexical-semantic information as it unfolds, according to the match between the evolving input and representations in memory (e.g., Marslen-Wilson & Zwitserlood, 1989).

Synchrony Between Visual and Auditory Speech Onsets—Our task presents a visual consonant + vowel onset coupled to an auditory silent gap + vowel onset (our methodological criterion created a silent gap of about 50 ms for words and 65 ms for nonwords, see Methods). A question is whether this change in the normal synchrony relation between the visual and auditory speech onsets impacts AV integration. The literature suggests not. Listeners normally have access to visual cues before auditory cues (Bell-Berti & Harris, 1981). Adults synthesize visual and auditory cues (without any detection of asynchrony or any effect on intelligibility) when visual speech leads auditory speech by as much as 200 ms (Grant, van Wassenhove, & Poeppel, 2004). Visual speech can lead auditory speech by as much as 180 ms without altering the McGurk effect (Munhall, Gribble, Sacco, & Ward, 1996). This literature suggests that cross-modal synchrony is probably not the basis of AV integration (Munhall & Tohkura, 1998).

Integrity Between Speech Cues—This issue has been studied with the Garner task (1974) in which participants selectively attend to features of stimuli (e.g., consonants vs vowels). Participants are asked to classify the stimuli on the basis of one of the features (e.g., /b/ vs /g/) while the other feature 1) remains constant (control condition, /ba/ vs /ga/) or 2) varies irrelevantly (/ba/, /bi/ vs /ga/, /gi/). Results have shown that irrelevant variation in vowels interferes with classifying consonants and vice versa (see, e.g., Tomiak, Mullennix, & Sawusch, 1987). Such a pattern of interaction indicates that AO speech cues are perceived in a mutually interdependent manner. A question is whether this tight coupling between AO speech cues generalizes to AV speech cues? Such a generalization is supported by a study with the Garner task (Green & Kuhl, 1989) showing that speech cues (e.g., voice onset time and place of articulation) are perceived in an interdependent manner not only when both of the cues are specified by AO speech but also when the voice onset time is specified by AO speech and the place of articulation is specified by a combination of AV speech (i.e., McGurk stimuli). These results imply that the auditory and visual speech cues of this research should be processed in an interdependent manner. Even though most participants listening to our stimuli in the AO mode report hearing a vowel onset (see Methods/Results), the vowel still contains some lawful variation from being produced in a consonant-vowel-consonant context rather than in isolation. To the extent that the speech perceptual system utilizes lawful variation (Tomiak et al., 1987), our visual speech onset cues may be more easily grafted onto the remaining compatible visual and auditory vowel-consonant cues to yield a unified AV percept. These results imply that the auditory and visual speech cues of this study should be processed in an interdependent manner. From this viewpoint, children's performance may be more sensitive to visual speech on our task than, for example, on the McGurk effect with its non-compatible auditory and visual content. This possibility will be addressed directly in Study 2. Below we predict results based on the effects of age, salience of the visual speech cues, and lexical status.

Predicted Results

Age—The evidence reviewed above on the development of AV speech perception predicts that children from about 5 to the pre-teen/teenage years will show reduced benefit from visual speech. To the extent that our new task is more sensitive to the influence of visual speech, other evidence predicts other possible age-related differences. Overall, however, the predictions are inconsistent.

a. Processing Skills: Younger children process AO speech cues less efficiently. As an example, gating studies document that younger children require more AO input to recognize words than teenagers (Elliott, Hammer, & Evan, 1987). Younger children also have less detailed and harder to access phonological representations (Snowling & Hulme, 1994). All this suggests that younger children may rely more on visual speech to supplement their less efficient processing of AO speech. Visual speech may also enhance because it 1) facilitates the detection of AO phonemes (Fort et. al., 2012), 2) provides extra phonetic information that facilitates the extraction of AO cues and reduces uncertainty (Campbell, 1988; Dodd, 1977), and 3) acts as a type of alerting mechanism that benefits younger children's immature attentional skills (Campbell, 2006). Finally, younger children with less mature articulatory proficiency tend to observe visual speech more-perhaps in order to

cement their knowledge of the acoustic consequences of articulatory gestures (Desjardins et al., 1997; Dodd, McIntosh, Erdener, & Burnham, 2008). These results suggest that the processing weights assigned to the auditory and visual modalities may shift with age, and younger children may show a greater influence of visual speech than older children.

b. Speechreading Skills: If any visual speech fill-in effect depends on visual speechreading, then older participants—being more proficient (See Table 1)—should show a greater visual speech fill-in effect than younger children.

Visual Speech Cues/Phonotactic Probabilities—The bilabial /b/ is more accurately speechread and is more common as an onset in English than the velar /g/ (Tye-Murray, 2009; Vitevitch & Luce, 2004; Storkel & Hoover, 2010; see Appendix, Table 1A for phonotactic probabilities). These differences predict that the non-intact /b/ onset will be more readily restored than the non-intact /g/ onset.

Lexical Status—To determine whether the availability of a lexical representation influences the visual speech fill-in effect, we compared performance for words vs nonwords (e.g., bag vs baz). As noted previously, theories propose that speech automatically activates its lexical representation and that activation of this knowledge 1) reduces the demand for processing resources and 2) facilitates the detection of phonemes (Bouton, Cole, & Serniclaes, 2012; Fort, Spinelli, Savariaux, & Kandel, 2010; Newman & Twieg, 2001; Rubin, Turvey, & van Gelder, 1976). Lexical status also affects visually influenced performance on the McGurk task, with the McGurk effect more prevalent for words than nonwords (Barutchu, Crewther, Kiely, Murphy, & Crewther, 2008) and for stimuli in which the visual stimulus forms a word and the auditory stimulus forms a nonword (Brancazio, 2004). To the extent that these results generalize to our task, we may see a greater visual speech fill-in effect for words than nonwords. An exception to this prediction is raised, however, by the finding that lexical-semantic access requires more attentional resources in 4–5-yr-olds than in older children (Jerger et al., 2013). This outcome predicts that the word stimuli may disproportionately drain the younger children's processing resources and reduce sensitivity to visual speech for words.

In addition to predicting results from the literature, we can predict results from theories of speech perception. A particularly relevant model for predicting the effects of lexical status is the TRACE theory of auditory speech perception with its interactive activation architecture consisting of acoustic feature, phoneme, and lexical levels (McClelland & Elman, 1986). In this model, information flows forwards (from feature to lexical level) and backwards (from lexical to feature level). Thus, the model proposes that the activation level of a phoneme is determined by activation from both the feature and lexical levels. For AV speech, Campbell (1988) proposes that visual speech adds visual features that feed forward to activate their associated phonemes. Thus for our AV stimuli, the visual speech features corresponding to the onset would be activated and would feed forward to activate the phoneme. With regard to the words vs nonwords, the response to nonwords in the TRACE model would be based on activation only at the feature and phoneme levels in the purest sense. In this case, performance would reflect the feature and phoneme pattern of activation. If, however, the nonwords partially activate a similar lexical item that in turn supports

phonological processing (see Gathercole et al., 1991), then performance might be driven by the lexical level as well. Our definition of the visual speech fill-in effect as the difference between performance for AV - AO modes should control for any lexical influences on performance that are not unique to visual speech.

Another relevant theory is the hierarchical model of speech segmentation (Mattys, White, & Melhorn, 2005) which proposes that in optimal situations, listeners assign the greatest weight to lexical-semantic content. If the lexical-semantic content is compromised, listeners switch and assign the greatest weight to phonetic-phonological content. If both the lexical-semantic and phonetic-phonological content are compromised, listeners switch and assign the greatest weight to acoustic-temporal (prosodic) content. It is also the case that monosyllabic words such as ours may activate their lexical representations without requiring phonological decomposition whereas nonwords require phonological decomposition (Mattys, 2014). If this model generalizes to our task, children in both the AO and AV modes should assign the greatest weight to lexical-semantic content for words and to phonetic-phonological content for nonwords. To the extent a greater weight on phonetics-phonology for nonwords increases children's attention to visual speech cues, we predict a greater visual speech fill-in effect for the nonwords than the words.

Method

Participants

The children were 92 native English speakers ranging in age (years;months) from 4;2 to 14;5 (53% boys). The racial distribution was 87% White, 7% Asian, and 4% Black, with 9% of participants reporting Hispanic ethnicity. The children were divided into four age groups: 4–5-yr-olds ($M = 4;11$), 6–7-yr-olds ($M = 6;11$), 8–9-yr-olds ($M = 8;10$), and 10–14-yr-olds ($M = 11;7$). Each group consisted of 22 children excepting the 10–14-yr group which contained 26 children. Visual perception, articulatory proficiency, and hearing and vision sensitivity were within normal limits for chronological age. Other demographic characteristics are detailed in Study 2.

Materials and Instrumentation

Stimuli—The stimuli were monosyllabic words and nonwords beginning with the consonants /b/ or /g/ coupled with the vowels /i/, /ae/, /a/, or /o/ (Appendix, Table 1A). The stimuli were recorded at the Audiovisual Recording Lab, Washington University School of Medicine. The talker was an 11-yr-old trained boy actor with clearly intelligible speech without pubertal characteristics (f_0 of 203 Hz). His full facial image and upper chest were recorded. He started and ended each utterance with a neutral face/closed mouth. The color video signal was digitized at 30 frames/s with 24-bit resolution at a 720×480 pixel size. The auditory signal was digitized at a 48 kHz sampling rate with 16-bit amplitude resolution. The utterances were adjusted to equivalent A-weighted root mean square sound levels.

Editing the Auditory Onsets—To edit the auditory track, we located the /b/ or /g/ onsets visually and auditorily with Adobe Premiere Pro and Soundbooth (Adobe Systems Inc., San

Jose, CA) and loudspeakers. We applied a perceptual criterion to operationally define a non-intact onset. We excised the waveform in 1 ms steps from the identified auditory onset to the point in the adjacent vowel for which at least 4 of 5 trained listeners (AO mode) heard the vowel as the onset. Splice points were always at zero axis crossings. Using this perceptual criterion, we excised on average 52 ms (/b/) and 50 ms (/g/) from the word onsets and 63 ms (/b/) and 72 ms (/g/) from the nonword onsets. Performance by young untrained adults for words (N=10) and nonwords (N=10) did not differ from the results presented herein for the 10–14-yr-old group. The visual track of the words and nonwords was also edited to form AV (dynamic face) vs AO (static face) modes of presentation.

AV vs AO Modes—All stimuli were presented as Quicktime movie files. The AV mode consisted of the talker's still neutral face and upper chest, followed by an AV utterance of a word or nonword, followed by the talker's still neutral face and upper chest. The AO mode consisted of the same auditory track but the visual track contained the talker's still neutral face and upper chest for the entire trial. The video track was routed to a high resolution computer monitor and the auditory track was routed through a speech audiometer to a loudspeaker.

Final Set of Items—The AV and AO modes of the word (or nonword) test items with intact and non-intact /b/ and /g/ auditory onsets were randomly intermixed and formed into lists, which also contained 14 filler items presented in the AV and AO modes. The filler items consisted of words (or nonwords) with intact not /b/ or /g/ consonant and vowel /i/, /ae/, /a/, or /o/ onsets. Illustrative filler items are the word/nonword pairs of eagle/eeble, apple/apper, cheese/cheeg, and table/tavel. Thus listeners heard trials randomly alternating between intact vs non-intact auditory onsets, AV vs AO modes, and test vs filler items. Each test item (intact and non-intact) was presented twice in each mode. These 64 test trials were intermixed with 48 filler trials, yielding 57% test trials. The set of 112 trials was divided into four lists (presented forward or backward for eight variations). The items comprising a list varied randomly under the constraints that 1) no onset could repeat, 2) the intact and non-intact pairs [e.g., bag and (-b)ag] could not occur without at least two intervening items, 3) a non-intact onset must be followed by an intact onset, 4) the mode must alternate after three repetitions, and 5) all types of onsets (intact /b/ and /g/, non-intact /b/ and /g/, vowels, and not /b/ or /g/) had to be dispersed uniformly throughout the lists. The number of intervening items between the intact vs non-intact pairs averaged 12 items. The intensity level of the stimuli was approximately 70 dB SPL. The responses of the participants were digitally recorded.

Procedure

The tester sat at a computer workstation and initiated each trial by pressing a touch pad (out of child's sight). The children, with a co-tester alongside, sat at a distance of 71 cm directly in front of an adjustable height table containing the computer monitor and loudspeaker. Their view of the talker's face subtended a visual angle of 7.17° vertically (eyebrow - chin) and 10.71° horizontally (eye level). The children completed the word/nonword repetition tasks along with other procedures in three sessions, scheduled approximately 10 days apart. In the first session, the children completed three of the word (or nonword) lists in separated

listening conditions; in the second session, the children completed the fourth word (or nonword) list and the first nonword (or word) list in separated conditions; and in the third session, the children completed the remaining three nonword (or word) lists in separated conditions. The order of presentation of the words vs nonwords was counterbalanced across participants in each age group. The analyses below were collapsed across the counterbalancing conditions.

The children were instructed to *repeat exactly* what the talker said. We told younger children that the task was a copy-cat game. For the words, participants were told that they might hear words or nonwords. For the nonwords, they were told that none of the utterances would be words. Due to the multiple procedures of our protocol, the children had heard the words and nonwords previously as distractors for the multimodal picture-word task (Jerger et al., 2009).

The children's utterances were transcribed independently by the tester and co-tester. For the utterances with non-intact onsets, the transcribers disagreed on 1.83% of word responses and 2.68% of nonword responses. For responses that were in disagreement, another trained listener independently transcribed the recorded utterances. Her transcription, which always agreed with one of the other transcribers, was recorded as the response. The criteria for scoring responses to the non-intact onsets were as follows.

1. Correct vowel onsets ["ean" for "(-b)ean"] were scored as an auditory-based response for both modes.
2. Correct consonant onsets ["bag" for "(-b)ag"] were scored as a visual-based response for the AV mode and as a coarticulatory/lexical-based response for the AO mode. Visemes (visually indistinguishable phonemes) of a consonant were counted as correct for the AV mode. Viseme alternatives represented less than 1% of correct responses for both words and nonwords.
3. Incorrect vowel or consonant onsets ["dear" for "(-g)ear"] were scored as errors.

Determination of Word Knowledge—Each parent identified each word that the child knew. For the remaining words, the word was considered known if the child could identify the word from a set of six alternatives and tell us about it. The number of unknown words identified by this approach averaged from 0.91 in the 4–5-yr-olds to 0.00 in the 10–14-yr-olds. All unidentified words were taught to the children. The results below did not differ for taught vs previously known words.

Results

Accuracy for Words and Nonwords with Intact or Non-Intact Onsets

The accuracy of repeating the *intact* words and nonwords in the two modes for all groups was near ceiling (>96%) for the onsets and the offsets (i.e., the remainder of the utterance after the onset). The accuracy of repeating the offsets of the *non-intact* stimuli in the AO vs AV modes respectively averaged 98.57% vs 98.78% (words) and 96.40% vs 94.77% (nonwords). Below we analyze the onset responses for the non-intact stimuli. The overall

proportion of onset errors for the AO vs AV modes respectively averaged 5.16% vs 2.38% (non-intact words) and 7.47% vs 2.79% (non-intact nonwords). To ensure that the responses in the children who made errors contributed equally to the group averages, the number of correct vowel and consonant onsets was normalized such that the sum of the correct responses always equaled 8. For example, if a child had 5 correct vowel responses, 2 correct consonant responses, and 1 error, her normalized data were 5.71 vowel responses and 2.29 consonant responses. Our initial analysis addressed whether performance for the AO baselines differed as a function of age, lexical status, or onset.

Stimuli with Non-Intact Onsets: AO Baselines

Figure 1 displays the proportion of correct consonant onset responses for the /b/ vs /g/ onsets of the words and nonwords in the AO mode. Results were analyzed with an analysis of variance with one between-participants factor (Group: ages 4–5, 6–7, 8–9, 10–14) and two within-participants factors (Stimulus: words vs nonwords; Onset: /b/ vs /g/). The overall collapsed results showed a significant age-related change with more correct consonant onset responses in the 4–5-yr-olds than in the 10–14-yr-olds (50% vs 28%), $F(3,88) = 5.16$, $MSE = 12.532$, $p = .002$, $\text{partial } \eta^2 = .150$. The overall proportion of correct consonant onset responses was also significantly greater for the words than nonwords (51% vs 23%) and for the /b/ than /g/ onsets (39% vs 34%), stimulus: $F(1,88) = 127.66$, $MSE = 9.593$, $p < .0001$, $\text{partial } \eta^2 = .592$; onset: $F(1,88) = 7.41$, $MSE = 0.758$, $p = .008$, $\text{partial } \eta^2 = .078$. Finally the /b/ onsets showed significantly more correct consonant onset responses than the /g/ onsets for the words (58% vs 44%) but not for the nonwords (21% vs 24%), with a significant stimulus \times onset interaction, $F(1,88) = 32.11$, $MSE = 1.480$, $p < .0001$, $\text{partial } \eta^2 = .267$. No other significant effects or interactions were observed.

In short, performance for the AO baselines showed more coarticulatory or lexically-based responses for the words than the nonwords and for the younger children than the older children. To probe whether this outcome was reflecting remaining coarticulatory cues or a lexical effect, we evaluated performance for the words in the children receiving the opposing counterbalancing conditions (words first vs nonwords first). Results supported a change in the weighting of the lexical-semantic information. When only children in the nonwords first condition were considered, the proportion of consonant onset responses for the word baselines dropped to 34% (instead of 51%). If performance had been reflecting remaining coarticulatory evidence in the waveform, results *for the same auditory waveform* should not have changed across the counterbalancing conditions. Finally, it is possible that results for the words may have differed depending on whether the non-intact word did or did not form a new vowel-onset word [(-g)ear vs (-g)uts]. To address this possibility, we carried out an item analysis. Performance did not differ as a function of whether the non-intact words did or did not form a new vowel-onset word. Overall results in Figure 1 indicate that the AO baselines for both the words and nonwords are sufficiently below ceiling to allow us to evaluate the difference in the proportion of correct responses for the AV vs AO modes (visual speech fill-in effect).

Visual Speech Fill-In Effect

Figure 2 shows the difference (AV - AO) in the proportion of correct consonant onset responses for the non-intact words and nonwords as a function of the onset. Results were analyzed with an analysis of variance consisting of one between-participants factor (Group: ages 4–5, 6–7, 8–9, 10–14) and two within-participants factors (Stimulus: words vs nonwords; Onset: /b/ vs /g/). The findings indicated a significant age-related increase in the visual speech fill-in effect, with an overall magnitude of only 6% in the 4–5-yr-olds but 25% in the 10–14-yr-olds, $F(3,88) = 16.13$, $MSE = 2.795$, $p < .0001$, partial $\eta^2 = .355$. The visual speech fill-in effect was also significantly larger for nonwords than words and for /b/ than /g/ onsets, with an overall magnitude of 25% for nonwords but only 12% for words and 35% for /b/ but only 2% for /g/, stimulus: $F(1,88) = 41.52$, $MSE = 2.435$, $p < .0001$, partial $\eta^2 = .321$; onset: $F(1,88) = 190.81$, $MSE = 3.189$, $p < .0001$, partial $\eta^2 = .684$.

With regard to the interactions, the /b/ onsets showed a significantly larger visual speech fill-in effect for nonwords than words (47% vs 22%) whereas the /g/ onsets did not (2% for both types of stimuli), producing a significant onset \times stimulus interaction, $F(1,88) = 42.77$, $MSE = 2.235$, $p < .0001$, partial $\eta^2 = .327$. Finally, the visual speech fill-in effect showed a significantly smaller difference between the /b/ and /g/ onsets in the 4–5-yr-olds (21%) than in the older children (about 36%), producing a significant onset \times age group interaction, $F(3,88) = 3.20$, $MSE = 3.189$, $p = .027$, partial $\eta^2 = .098$. No other significant effects or interactions were observed.

To determine whether each age group showed a significant visual speech fill-in effect, we conducted multiple t-tests assessing whether each difference score differed from zero. The multiple comparison problem was controlled with the false discovery rate (FDR) procedure (Benjamini & Hochberg, 1995). Results for the nonwords showed a significant visual speech fill-in effect in all age groups for the /b/ onsets and in the 8–9-yr-olds for the /g/ onsets. Results for the words showed a significant visual speech fill-in effect for the 6–7-, 8–9-, and 10–14-yr-olds for the /b/ onsets. FDR results for the /g/ onsets approached significance in the 10–14-yr-olds for both the words and nonwords.

To address whether the smaller visual speech fill-in effect for the words relative to the nonwords—with /b/ onsets—was associated with differences in the AO baselines (Figure 1), we evaluated performance only in the children in the nonwords first condition who had lower AO baseline performance (34%). Results for the /b/ onsets continued to show a significantly smaller visual speech fill-in for the words than nonwords, $F(1,42) = 13.65$, $MSE = 3.533$, $p = .0006$, partial $\eta^2 = .245$. Thus the smaller visual speech fill-in effect for the words with /b/ onsets was not reflecting any limitation on performance associated with the slightly higher AO baseline seen in Figure 1. In short, Study 1 indicated that age, lexical status, and speechreadability/phonotactic probability influenced the visual speech fill-in effect. In the next study, we investigated the effects of task demands and child factors on visually influenced performance by comparing results on the visual speech fill-in and McGurk tasks.

Study 2

As discussed earlier, a dynamic systems viewpoint (Smith & Thelan, 2003) suggests that the influence of visual speech in children may vary as a function of task demands. To probe the effect of task demands and to identify the child factors that influence performance, we compared the influence of visual speech on performance for the nonwords of our new task (naturally compatible intact visual and non-intact auditory /b/ onsets) vs the McGurk task (conflicting auditory and visual /b/ and /g/ onsets: auditory /bΛ/ and visual /gΛ/) (McGurk & MacDonald, 1976). Below we predict results from our theories and the literature.

Predicted Results

Relative Weighting of Auditory vs Visual Speech—The literature shows a shift in the relative weights of the auditory and visual modes as the quality of the input shifts. For examples, children with *normal* hearing, listening to McGurk stimuli with lower fidelity (spectrally reduced) auditory speech, respond more on the basis of the intact visual input (Huyse, Berthommier, & Leybaert, 2013); but when the visual input is also artificially degraded, the children shift and respond more on the basis of the lower fidelity auditory input. Children with normal hearing or mild-moderate hearing loss and good auditory word recognition-listening to conflicting auditory and visual inputs such as auditory /meat/ coupled with visual /street/-respond on the basis of the auditory input (Seewald, Ross, Giolas, & Yonovitz, 1985). In contrast, children whose hearing loss is more severe-and whose perceived auditory input is more degraded-respond more on the basis of the visual input. Finally Japanese adults and children barely show a McGurk effect for intact auditory input but show a significant McGurk effect when the auditory input is degraded by noise or an unfamiliar (e.g., foreign) accent (Sekiyama & Burnham, 2008; Sekiyama & Tohkura, 1991). These results suggest that children exploit the relative quality of stimulus attributes to modulate the relative weighting of auditory and visual speech. We hypothesize that children's relative weighting of the auditory vs visual speech inputs will depend on the quality of the auditory input. With intact auditory input—such as McGurk stimuli—speech perception in children will be more auditory bound; with non-intact auditory onsets coupled to intact visual onsets, however, the children's performance will depend more on visual speech. If so, performance for our task will be more sensitive to the influence of visual speech than the McGurk stimuli.

Compatible vs Incompatible Audiovisual Input—A compatible AV utterance whose onsets are within the time window producing a perception of synchrony (Grant et al., 2004; Munhall et al., 1996) is more likely to be treated as a single multisensory event (Vatakis & Spence, 2007). For example, Vatakis and Spence manipulated the temporal onsets of auditory and visual inputs that were matching or not. Listeners were significantly less sensitive to temporal differences between the matched onsets. We view the silent period characterizing our non-intact auditory onsets as more harmonious with the concurrent visual speech onsets than a McGurk stimulus that has conflicting intact onsets. Finally, another seemingly relevant consideration is that AO and AV consonant-vowel stimuli are processed in an interdependent manner on the Garner task (1974) by adults (Tomiak et al., 1987) and by children (Jerger et al., 1993). The latter study assessed performance for other types of

speech cues with the Garner task in individuals from 3 yrs to 79 yrs and found interdependent processing at all ages. To the extent these results generalize to our task, results on the Garner task further suggest that our auditory and visual onsets should be processed interdependently. Overall these data suggest that listeners will more likely show a greater influence of visual speech on our task than the McGurk task.

Child Characteristics Underpinning Task Performance

Speechreadability—Study 1 indicates that children have difficulty speechreading the /g/ onset, a finding that agrees with the literature (Tye-Murray, 2009). This suggests that listeners will likely show a greater influence of visual speech on our task, particularly for the /b/ visual onset, than on the McGurk task with its /g/ visual onset. To the extent that performance is reflecting speechreading, older children with better speechreading skills should show a greater sensitivity to visual speech.

Language—With relation to the TRACE model discussed previously, older children may benefit more from visual speech because they have more robust, detailed phonological and lexical representations (Snowling & Hulme, 1994) that are more easily activated by sensory input. Further the older children having stronger language skills may be able to utilize the activation pattern across the feature, phoneme, and word levels better than the younger children. Thus children with more mature language skills may show greater sensitivity to the influence of visual speech.

Method

Participants

Participants were the 92 children of Study 1. Table 1 summarizes their average ages, vocabulary skills and visual speechreading skills. Vocabulary skills were within normal limits for all groups. To quantify visual speechreading ability in the following regression analyses, we used the results scored by word onsets.

Stimuli and Procedure

Receptive vocabulary skills were estimated with the Peabody Picture Vocabulary Test-IV (Dunn & Dunn, 2007); the children heard a word and pointed to the picture-out of four alternatives-illustrating that word's meaning. Speechreading skills were estimated with the Children's Audio-Visual Enhancement Test (Tye-Murray & Geers, 2001); children repeated words presented in the AO and visual only (VO) modes. The stimuli for the visual speech fill-in task were the nonwords with /b/ onsets. The stimuli for the McGurk task (McGurk & MacDonald, 1976) were /bΛ/ and /gΛ/ utterances recorded at the same Audiovisual Recording Lab by the same talker described above. The auditory track of /bΛ/ was combined with the visual track of /gΛ/ with Adobe Premiere Pro and Soundbooth. The auditory and visual utterances were aligned during the release of the consonant. The final McGurk stimulus (auditory /bΛ/ - visual /gΛ/) was presented as a Quicktime movie with the talker's still neutral face and upper chest, followed by the audiovisual utterance of the stimulus, followed by the talker's still neutral face and upper chest e appended the McGurk stimulus to the end of each of the nonword lists described above. The influence of visual speech was

quantified 1) by the difference in the number of correct consonant onset responses for the AV vs AO modes for the visual speech fill-in task ($N=8$) and 2) by the absolute number of visually influenced responses for the McGurk task ($N=4$). These derived measures were tallied as reflecting the influence of visual speech if results showed 1) a visual speech fill-in effect of at least 25% (difference score ≥ 2) or 2) at least 25% visually based responses for the McGurk trials ($N = 1$).

Results

Comparison of Visual Speech Fill-In vs McGurk Effects

Visual speech influenced performance in 82% of children for the visual speech fill-in task and 52% of children for the McGurk task. The visually influenced McGurk responses in the children consisted of /dΛ/, /ðΛ/, and /gΛ/. Thirty-nine percent (39%) of children showed only a visual speech fill-in effect and 10% of children showed only a McGurk effect. We conducted a multiple regression analysis for each task with predictor variables of age, vocabulary skill, and visual speechreading skill and a criterion variable of the quantified effect of visual speech. The intercorrelations among this set were .526 (age & visual speechreading), .105 (age & vocabulary), and -.058 (vocabulary & visual speechreading).

Table 2 summarizes the regression results, with the multiple correlation coefficients and omnibus F statistics for all of the variables considered simultaneously followed by the part (also called semi-partial) correlation coefficients and the partial F statistics evaluating the variation in performance uniquely accounted for (after removing the influence of the other variables) by each individual variable (Abdi et al., 2009). The set of variables significantly predicted the influence of visual speech on performance, with the children's ages, vocabulary skills, and visual speechreading skills together predicting approximately 17% to 18% of the variance in performance for each task. Part correlations (expressing the unique influence of each variable) indicated that the visual speech fill-in effect varied significantly as a function of the children's ages and vocabulary skills, each uniquely accounting for about 7% to 8% of the variance in performance. In contrast, the influence of visual speech on McGurk performance varied significantly only as a function of the children's speechreading skills, which uniquely accounted for about 13% of the variance in performance.

In short, the influence of visual speech on the visual speech fill-in and McGurk tasks appears to be underpinned by dissociable independent variables. A possible caveat is that results for the visual speech fill-in and McGurk tasks may have been influenced by the different baseline levels of performance. Thus we analyzed the visual speech fill-in effect for the nonwords with /g/ onsets (for which visual speech also exerted a lesser influence on performance, Study 1). The set of variables significantly predicted the visual speech fill-in effect and accounted for 9% of the variance in performance, $R = .346$, $F(3,88) = 3.89$, $p = .012$. Part correlations indicated that the visual speech fill-in effect varied significantly as a unique function of the children's age, $F(1,88) = 5.69$, $p = .019$, with neither their vocabulary ($p = .12$) nor their speechreading ($p = .639$) skills achieving significance. These data support the conclusion that age and vocabulary skills underlie the visual speech fill-in effect but speechreading skills underlie the McGurk effect.

Discussion

In this research, Study 1 assessed the effects of age, salience of visual speech cues, and lexical status on the new visual speech fill-in effect, and Study 2 assessed the effects of task/stimulus demands and child factors by comparing results on the visual speech fill-in and McGurk tasks. Results of Study 1 showed—contrary to previous findings—that children from 4 to 14 years of age significantly benefited from visual speech. However, this benefit critically depended on task/stimulus properties and individual abilities. With regard to age, results for nonwords showed a significant visual speech fill-in effect in all age groups (4–14-yrs) for easy visual speech cues (/b/) and in the 8–9-yr-olds for difficult visual speech cues (/g/). Results for the words showed a significant visual speech fill-in effect in all age groups except the 4–5-yr-olds for the easy visual speech cues. Results for both the nonwords and words approached significance for the difficult visual speech cues (/g/) in the 10–14-yr-olds when controlling for multiple comparisons. These results disagree with predictions that younger children with less mature linguistic and processing skills will rely on visual speech to a greater extent than older children. The results do, however, support the prediction that word stimuli disproportionately drain processing resources in 4–5-yr-olds and reduce sensitivity to visual speech because lexical-semantic access requires more attentional resources in these younger children than in older children (Jerger et al., 2013).

Our finding of an age-related increase in children's sensitivity to visual speech agrees with the literature in general; however, our results document an influence of visual speech at much younger ages than previously observed (e.g., Desjardins et al., 1997; Erdener & Burnham, 2013; Ross et al., 2011; Sekiyama & Burnham, 2008; Tremblay et al., 2007). The current developmental functions are not, however, consistent with the U-shaped developmental functions observed on the multimodal picture-word (phonological priming) task (Jerger et al., 2009), which showed a significant influence of congruent visual speech in 4-yr-olds and 12-yr-olds, but not in 5–9-yr-olds. Clearly sensitivity to visual speech varies dramatically as a function of task/stimulus demands.

As implied in the developmental patterns above, other critical task/stimulus properties were the salience of the visual speech cues and semantic content. The visual speech fill-in effect was significantly larger for easy than difficult speech cues and for nonwords than words. The latter results imply that children's sensitivity to visual speech may vary depending on their relative weighting and decomposition of the phonetic-phonological content in agreement with the proposals of the hierarchical model of speech segmentation (Mattys et al., 2005). Our results also suggest that compatible AV nonwords are more optimal stimuli than words when research goals are concerned with assessing children's sensitivity to visual speech. That said, words are clearly the building blocks of language and remain essential stimuli for studying the contributions of visual speech to children's communicative abilities.

Results of Study 2 showed that the McGurk task significantly underestimates children's sensitivity to visual speech compared to the visual speech fill-in task. Visually influenced performance was revealed in about 80% of children on the visual speech fill-in task but in only about 50% of children on the McGurk task. Thirty-nine percent (39%) of children showed only a visual speech fill-in effect. The relatively greater sensitivity to visual speech

of our task may be reflecting previous findings indicating that a compatible AV utterance is more likely to be viewed as a single multisensory event and more likely to be synthesized and thus processed in an interdependent manner (Tomiak et al., 1987; Vatakis & Spence, 2007). The current outcome is also consistent with the research reporting that children's relative weighting of auditory and visual inputs depends on the quality of the input-with speech perception for intact auditory input such as McGurk stimuli more auditory bound (e.g., Huyse et al., 2013; Seewald et al., 1985; Sekiyama & Burnham, 2008). Importantly Study 2 also documented that sensitivity to visual speech varies *in the same children* depending on the task/stimulus demands.

Finally and perhaps most importantly is the finding that age and vocabulary uniquely determine visually influenced performance for the visual speech fill-in task whereas speechreading skills uniquely determine visually influenced performance on the McGurk task. Findings for the McGurk task agree with previous investigators who proposed that greater sensitivity to visual speech reflects greater speechreading skills (Erdener & Burnham, 2013; see Jerger et al., 2009, for opposing evidence re: phonological priming). Findings for the visual speech fill-in effect may agree with previous investigators who proposed that greater sensitivity to visual speech is due to a heightened focus on the phonemic distinctions of one's native language (e.g., Erdener & Burnham, 2013). To the extent that phonological skills are critical to how well children learn vocabulary (e.g., Gathercole & Baddeley, 1989), our vocabulary association may be compatible with Erdener and Burnham's native phonology association, with both effects representing skills that are contingent on more basic phonological mechanisms.

The association between language skills and the visual speech fill-in effect may also explain why younger children benefit less from visual speech. Again, the TRACE theory of AO speech perception proposes that input activates acoustic feature, phoneme, and lexical levels (McClelland & Elman, 1986). For AV speech, Campbell (1988) proposes that visual speech activates visual features and their associated phonemes and so older children may benefit more from visual speech because they have more robust, detailed phonological and lexical representations (Snowling & Hulme, 1994) that are more easily activated by sensory input. Further the older children-having stronger language skills-may be able to utilize the activation pattern across the feature, phoneme, and word levels better than the younger children.

Finally, Study 2 documents that benefit from visual speech is not a unitary phenomenon. The McGurk effect was affected by one child factor (speechreading) but not by the other two (age and vocabulary) whereas the visual speech fill-in effect was affected by two of the child factors (age and vocabulary) but not by the other one (speechreading). This dissociation implies that benefit from visual speech is a complicated multi-faceted phenomenon underpinned by heterogeneous abilities.

In conclusion, these results show that-under some conditions-pre/school and elementary school-aged children benefit from visual speech during multimodal speech perception. Importantly, our new task extends the range of measures that can be used to assess visual speech processing by children and provides results critical to integrating visual speech into

our theories of speech perceptual development. We conceptualize the new visual speech fill-in effect as tapping into a perceptual process that may enhance the accuracy of speech perception in everyday listening conditions. These results emphasize that children experience hearing a speaker's utterance rather than the auditory stimulus per se. In children as in adults, there is more to speech perception than meets the ear.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

This research was supported by the National Institute on Deafness and Other Communication Disorders, grant DC-00421. Sincere appreciation to the speech science colleagues who advised us to adopt a perceptual criterion for editing the non-intact stimuli. We thank Dr. Brent Spar for recording the audiovisual stimuli. We thank the children and parents who participated and the research staff who assisted, namely Aisha Aguilera, Carissa Dees, Nina Dinh, Nadia Dunkerton, Alycia Elkins, Brittany Hernandez, Cassandra Karl, Demi Krieger, Michelle McNeal, Jeffrey Okonye, Rachel Parra, and Kimberly Periman of UT-Dallas (data collection, analysis, presentation), and Derek Hammons and Scott Hawkins of UT-Dallas and Brent Spehar of CID-Washington University (computer programming).

Source of Funding. National Institute on Deafness and Other Communication Disorders, grant DC-00421.

References

References

- Arnold P, Hill F. Bisensory augmentation: A speechreading advantage when speech is clearly audible and intact. *British Journal of Psychology*. 2001; 92:339–355.
- Barutcu A, Crewther S, Kiely P, Murphy M, Crewther D. When /b/ill with /g/ill becomes /d/ill: Evidence for a lexical effect in audiovisual speech perception. *European Journal of Cognitive Psychology*. 2008; 20:1–11.
- Bashford J, Warren R. Multiple phonemic restorations follow the rules for auditory induction. *Perception & Psychophysics*. 1987; 42:114–121. [PubMed: 3627931]
- Bell-Berti F, Harris K. A temporal model of speech production. *Phonetica*. 1981; 38:9–20. [PubMed: 7267724]
- Benjamini Y, Hochberg Y. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society Series B (Methodological)*. 1995; 57:289–300.
- Bjorklund, D.; Harnishfeger, K. The evolution of inhibition mechanisms and their role in human cognition and behavior. In: Dempster, F.; Brainerd, C., editors. *Interference and inhibition in cognition*. San Diego: Academic Press; 1995. p. 141–173.
- Bonte M, Blomert L. Developmental changes in ERP correlates of spoken word recognition during early school years: A phonological priming study. *Clinical Neurophysiology*. 2004; 115:409–423. [PubMed: 14744584]
- Bouton S, Cole P, Serniclaes W. The influence of lexical knowledge on phoneme discrimination in deaf children with cochlear implants. *Speech Communication*. 2012; 54(2):189–198.
- Brancazio L. Lexical influences in audiovisual speech perception. *Journal of Experimental Psychology: Human Perception and Performance*. 2004; 30:445–463. [PubMed: 15161378]
- Bryant, P. Phonological and grammatical skills in learning to read. In: deGelder, B.; Morais, J., editors. *Speech and reading: A comparative approach*. Erlbaum (UK), Taylor & Francis: Hove, East Sussex; 1995. p. 249–266.
- Burnham D. Language specific speech perception and the onset of reading. *Reading and Writing: An inter disciplinary Journal*. 2003; 16:573–609.

- Campbell, R. Audio-visual speech processing. In: Brown, K.; Anderson, A.; Bauer, L.; Berns, M.; Hirst, G.; Miller, J., editors. The encyclopedia of language and linguistics. Amsterdam: Elsevier; 2006. p. 562-569.
- Campbell R. Tracing lip movements: Making speech visible. *Visible Language*. 1988; 22:32–57.
- Conrad R. The chronology of the development of covert speech in children. *Developmental Psychology*. 1971; 5:398–405.
- Desjardins R, Rogers J, Werker J. An exploration of why preschoolers perform differently than do adults in audiovisual speech perception tasks. *Journal of Experimental Child Psychology*. 1997; 66:85–110. [PubMed: 9226935]
- Dodd B. The role of vision in the perception of speech. *Perception*. 1977; 6:31–40. [PubMed: 840618]
- Dodd B, McIntosh B, Erdener D, Burnham D. Perception of the auditory-visual illusion in speech perception by children with phonological disorders. *Clinical Linguistics & Phonetics*. 2008; 22:69–82. [PubMed: 18092221]
- Dunn, L.; Dunn, D. The peabody picture vocabulary test-IV. Fourth ed. Minneapolis, MN: NCS Pearson, Inc; 2007.
- Elliott L, Hammer M, Evan K. Perception of gated, highly familiar spoken monosyllabic nouns by children, teenagers and older adults. *Perception and Psychophysics*. 1987; 42:150–157. [PubMed: 3627935]
- Erdener D, Burnham D. The relationship between auditory-visual speech perception and language-specific speech perception at the onset of reading instruction in English-speaking children. *Journal of Experimental Child Psychology*. 2013; 114:120–138. [PubMed: 23773915]
- Fernald A, Swingley D, Pinto J. When half a word is enough: Infants can recognize spoken words using partial phonetic information. *Child Development*. 2001; 72:1003–1015. [PubMed: 11480931]
- Fort M, Spinelli E, Savariaux C, Kandel S. Audiovisual vowel monitoring and the word superiority effect in children. *International Journal of Behavioral Development*. 2012; 36(6):457–467.
- Fort M, Spinelli E, Savariaux C, Kandel S. The word superiority effect in audiovisual speech perception. *Speech Communication*. 2010; 52(6):525–532.
- Garner, W. The processing of information and structure. Potomax, MD: Erlbaum; 1974.
- Gathercole S, Baddeley A. Evaluation of the role of phonological STM in the development of vocabulary in children: A longitudinal study. *Journal of Memory and Language*. 1989; 28:200–213.
- Gathercole S, Willis C, Emslie H, Baddeley A. The influence of number of syllables and wordlikeness on children's repetition of nonwords. *Applied Psycholinguistics*. 1991; 12:349–367.
- Gow D, Melvold J, Manuel S. How word onsets drive lexical access and segmentation: Evidence from acoustics, phonology, and processing. *Spoken Language ICSLP Proceedings of the Fourth International Conference*. 1996; 1:66–69.
- Grant K, vanWassenhove V, Poeppel D. Detection of auditory (cross-spectral) and auditory-visual (cross-modal) synchrony. *Speech Communication*. 2004; 44:43–53.
- Green, K. The use of auditory and visual information during phonetic processing: Implications for theories of speech perception. In: Campbell, R.; Dodd, B.; Burnham, D., editors. *Hearing by eye II. Advances in the psychology of speechreading and auditory-visual speech*. Hove, UK: Taylor & Francis; 1998.
- Green K, Kuhl P. The role of visual information in the processing of place and manner features in speech perception. *Perception & Psychophysics*. 1989; 45:34–42. [PubMed: 2913568]
- Gupta P, Lipinski J, Abbs B, Lin PH. Serial position effects in nonword repetition. *Journal of Memory and Language*. 2005; 53:141–162.
- Horlyck S, Reid A, Burnham D. The relationship between learning to read and language-specific speech perception: maturation versus experience *Scientific Studies of. Reading*. 2012; 16(3):218–239.
- Huyse A, Berthommier F, Leybaert J. Degradation of labial information modifies audiovisual speech perception in cochlear-implanted children. *Ear and Hearing*. 2013; 34(1):110–121. [PubMed: 23059850]

- Jerger S, Damian M, Mills C, Bartlett J, Tye-Murray N, Abdi H. Effect of perceptual load on semantic access by speech in children. *Journal of Speech Language & Hearing Research*. 2013; 56:388–403.
- Jerger S, Damian MF, Spence MJ, Tye-Murray N, Abdi H. Developmental shifts in children's sensitivity to visual speech: A new multimodal picture-word task. *Journal of Experimental Child Psychology*. 2009; 102:40–59. [PubMed: 18829049]
- Jerger S, Pirozzolo F, Jerger J, Elizondo R, Desai S, Wright E, Reynosa R. Developmental trends in the interaction between auditory and linguistic processing. *Perception & Psychophysics*. 1993; 54:310–320. [PubMed: 8414890]
- Johnson CE. Children's phoneme identification in reverberation and noise. *Journal of Speech Language and Hearing Research*. 2000; 43:144–157.
- MacLeod A, Summerfield Q. Quantifying the contribution of vision to speech perception in noise. *British Journal of Audiology*. 1987; 21:131–141. [PubMed: 3594015]
- Marslen-Wilson W, Zwitserlood P. Accessing spoken words: The importance of word onsets. *Journal of Experimental Psychology: Human Perception and Performance*. 1989; 15:576–585.
- Mattys, S. Speech perception. In: Reisberg, D., editor. *The Oxford handbook of cognitive psychology*. Oxford, UK: Oxford University Press; 2014. p. 391–412.
- Mattys S, White L, Melhorn J. Integration of multiple speech segmentation cues: A hierarchical framework. *Journal of Experimental Psychology-General*. 2005; 134:477–500. [PubMed: 16316287]
- McClelland J, Elman J. The TRACE model of speech perception. *Cognitive Psychology*. 1986; 18:1–86. [PubMed: 3753912]
- McGurk H, MacDonald M. Hearing lips and seeing voices. *Nature*. 1976; 264:746–748. [PubMed: 1012311]
- Munhall K, Gribble P, Sacco L, Ward M. Temporal constraints on the McGurk effect. *Perception & Psychophysics*. 1996; 58:351–362. [PubMed: 8935896]
- Munhall K, Tohkura Y. Audiovisual gating and the time course of speech perception. *Journal of the Acoustical Society of America*. 1998; 104:530–539. [PubMed: 9670544]
- Newman R. Perceptual restoration in children versus adults. *Applied Psycholinguistics*. 2004; 25(4): 481–493.
- Newman S, Twieg D. Differences in auditory processing of words and pseudowords: An fMRI study. *Human Brain Mapping*. 2001; 14:39–47. [PubMed: 11500989]
- Ross L, Molholm S, Blanco D, Gomez-Ramirez M, Saint-Amour D, Foxe J. The development of multisensory speech perception continues into the late childhood years. *European Journal of Neuroscience*. 2011; 33:2329–2337.
- Rubin P, Turvey MT, Vangelder P. Initial phonemes are detected faster in spoken words than in spoken nonwords. *Perception & Psychophysics*. 1976; 19:394–398.
- Samuel AG. Phonemic restoration: Insights from a new methodology. *Journal of Experimental Psychology-General*. 1981; 110(4):474–494. [PubMed: 6459403]
- Seewald RC, Ross M, Giolas TG, Yonovitz A. Primary modality for speech perception in children with normal and impaired hearing. *Journal of Speech and Hearing Research*. 1985; 28(1):36–46. [PubMed: 3981996]
- Sekiyama K, Burnham D. Impact of language on development of auditory-visual speech perception. *Developmental Science*. 2008; 11(2):306–320. [PubMed: 18333984]
- Sekiyama K, Tohkura Y. McGurk effect in non-English listeners: Few visual effects for Japanese subjects hearing Japanese syllables of high auditory intelligibility. *Journal of the Acoustical Society of America*. 1991; 90:1797–1805. [PubMed: 1960275]
- Shahin A, Bishop C, Miller L. Neural mechanisms for illusory filling-in of degraded speech. *Neuroimage*. 2009; 44:1133–1143. [PubMed: 18977448]
- Shahin A, Kerlin J, Bhat J, Miller L. Neural restoration of degraded audiovisual speech. *Neuroimage*. 2012; 60:530–538. [PubMed: 22178454]
- Shahin A, Miller L. Multisensory integration enhances phonemic restoration. *Journal of the Acoustical Society of America*. 2009; 125:1744–1750. [PubMed: 19275331]

- Smith L, Thelen E. Development as a dynamic system. *Trends in Cognitive Sciences*. 2003; 7:343–348. [PubMed: 12907229]
- Snowling M, Hulme C. The development of phonological skills. *Philosophical Transactions of the Royal Society of London Series B*. 1994; 346:21–27. [PubMed: 7886149]
- Storkel H, Hoover J. An on-line calculator to compute phonotactic probability and neighborhood density based on child corpora of spoken American English. *Behavior Research ' Methods*. 2010; 42(2):497–506. [PubMed: 20479181]
- Tomiak G, Mullennix J, Sawusch J. Integral processing of phonemes: Evidence for a phonetic mode of perception. *Journal of the Acoustical Society of America*. 1987; 81:755–764. [PubMed: 3584684]
- Treiman R, Zukowski A. Children's sensitivity to syllables, onsets, rimes, and phonemes (vol 61, pg 433, 1996). *Journal of Experimental Child Psychology*. 1996; 62(3):432–455. [PubMed: 8812054]
- Tremblay C, Champoux R, Voss P, Bacon B, Lepore F, Theoret H. Speech and non-speech audio-visual illusions: A developmental study. *PLoS One*. 2007; 2(8):e742. [PubMed: 17710142]
- Trout JD, Poser WJ. Auditory and visual influences on phonemic restoration. *Language and Speech*. 1990; 33:121–135. [PubMed: 2283923]
- Tye-Murray, N. *Foundations of aural rehabilitation: Children, adults, and their family members*. 3rd ed. Singular Publishing Group: San Diego; 2009.
- Tye-Murray, N.; Geers, A. *Children's audio-visual enhancement test*. Central Institute for the Deaf: St. Louis, MO; 2001.
- Vatakis A, Spence C. Crossmodal binding: Evaluating the “unity assumption” using audiovisual speech stimuli. *Perception & Psychophysics*. 2007; 69(5):744–756. [PubMed: 17929697]
- Vitevitch M, Luce P. A web-based interface to calculate phonotactic probability for words and nonwords in English. *Behavior Research Methods, Instruments, & Computers*. 2004; 36:481–487.
- Walley AC. Spoken word recognition by young children and adults. *Cognitive Development*. 1988; 3(2):137–165.
- Warren RM. Perceptual restoration of missing speech sounds. *Science*. 1970; 167(3917):392–393. [PubMed: 5409744]

Appendix

Table 1A
The word and nonword items, which were constructed to have as comparable phonotactic probabilities as possible.
a. Adult values (Vitevitch & Luce, 2004)

	Segment Frequency		Biphone Frequency		
	Word	Nonword	Word	Nonword	
bag	baz	0.1486	0.1507	0.0087	0.0066
bean	beece	0.1791	0.1619	0.0053	0.0040
bone	bohs	0.1996	0.1793	0.0045	0.0047
bus	buhl	0.1692	0.1641	0.0073	0.0080
	average	0.1741	0.1640	0.0064	0.0058
gap	gak	0.1425	0.1589	0.0081	0.0086
gear	geen	0.1361	0.1538	0.0001	0.0031
gold	gothd	0.1181	0.1187	0.0015	0.0016
guts	guks	0.1813	0.1687	0.0048	0.0084
	average	0.1445	0.1500	0.0036	.0054
overall	average	0.1593	0.1570	0.0050	0.0056

b. Child values (Storkel & Hoover, 2010)

Segment Frequency			Biphone Frequency		
Word	Nonword	Word	Nonword	Word	Nonword
bag	baz	0.1776	0.1827	0.0108	0.0089
ban	beece	0.2244	0.1813	0.0082	0.0063
bone	bohs	0.2391	0.1961	0.0080	0.0073
bus	buv	0.2003	0.1601	0.0110	0.0089
	average	0.2103	0.1801	0.0095	0.0079
gap	gak	0.1481	0.1763	0.0063	0.0088
gear	geen	0.1642	0.1816	0.0008	0.0044
gold	gothd	0.1580	0.1568	0.0031	0.0030
guts	guks	0.2170	0.2094	0.0082	0.0120
	average	0.1718	0.1810	0.0046	0.0054
overall	average	0.1911	0.1805	0.0071	0.0074

Note: The positional segment frequency is the sum of the likelihoods of occurrence of each phoneme in its position within the utterance; the biphone frequency is the sum of the likelihoods of co-occurrence of each two adjacent phonemes.

Highlights

1. New child test to assess whether visual speech fills-in non-intact auditory speech
2. Children 4–14 years old benefited significantly from visual speech in some conditions
3. Visual speech fill-in effect shows greater sensitivity to visual speech than McGurk
4. Age and vocabulary underpin visual speech fill-in but lipreading underpins McGurk
5. Benefit from visual speech is multi-faceted phenomenon based on different skills

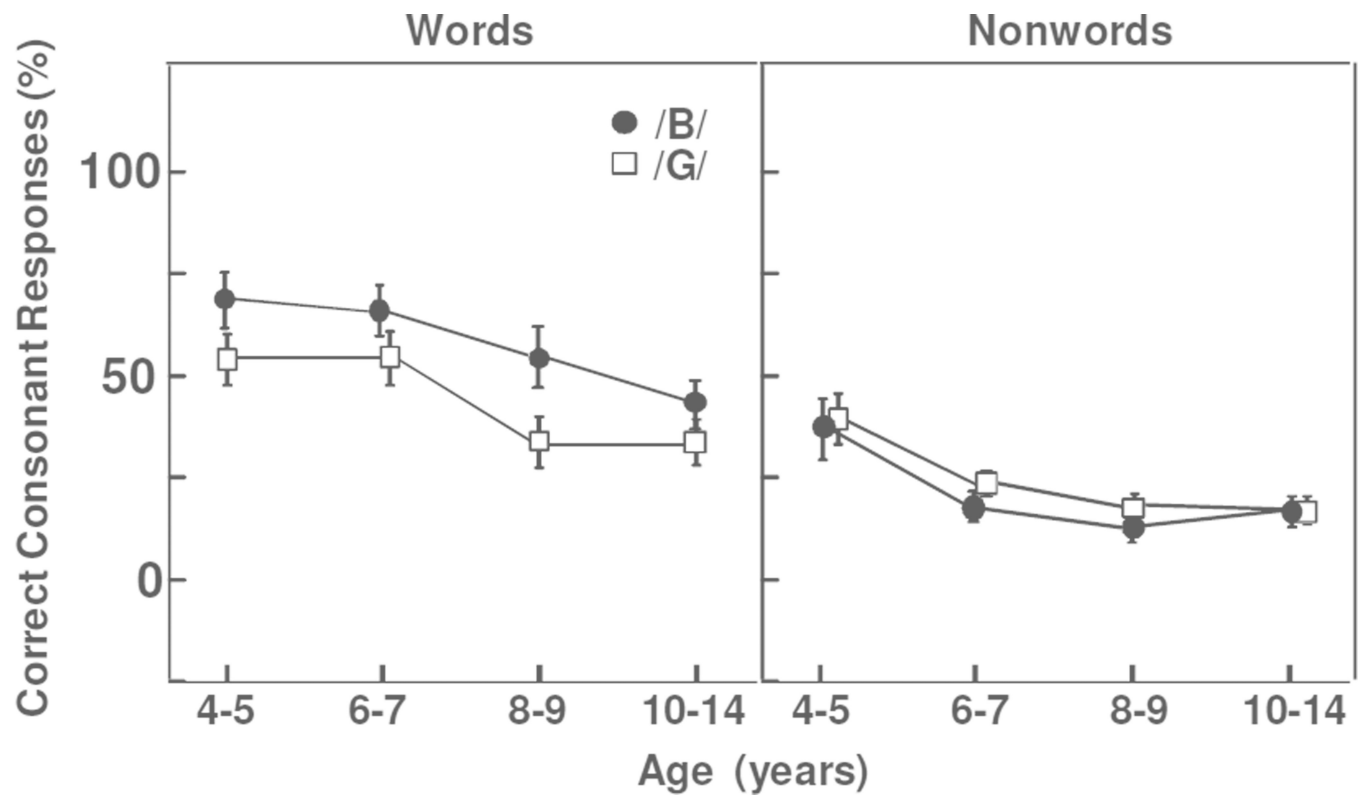


Figure 1.

Comparison of AO baselines for words and nonwords in the four age groups. A consonant onset response indicates that remaining coarticulatory information in the vowel or lexical effects formed the basis of the response.

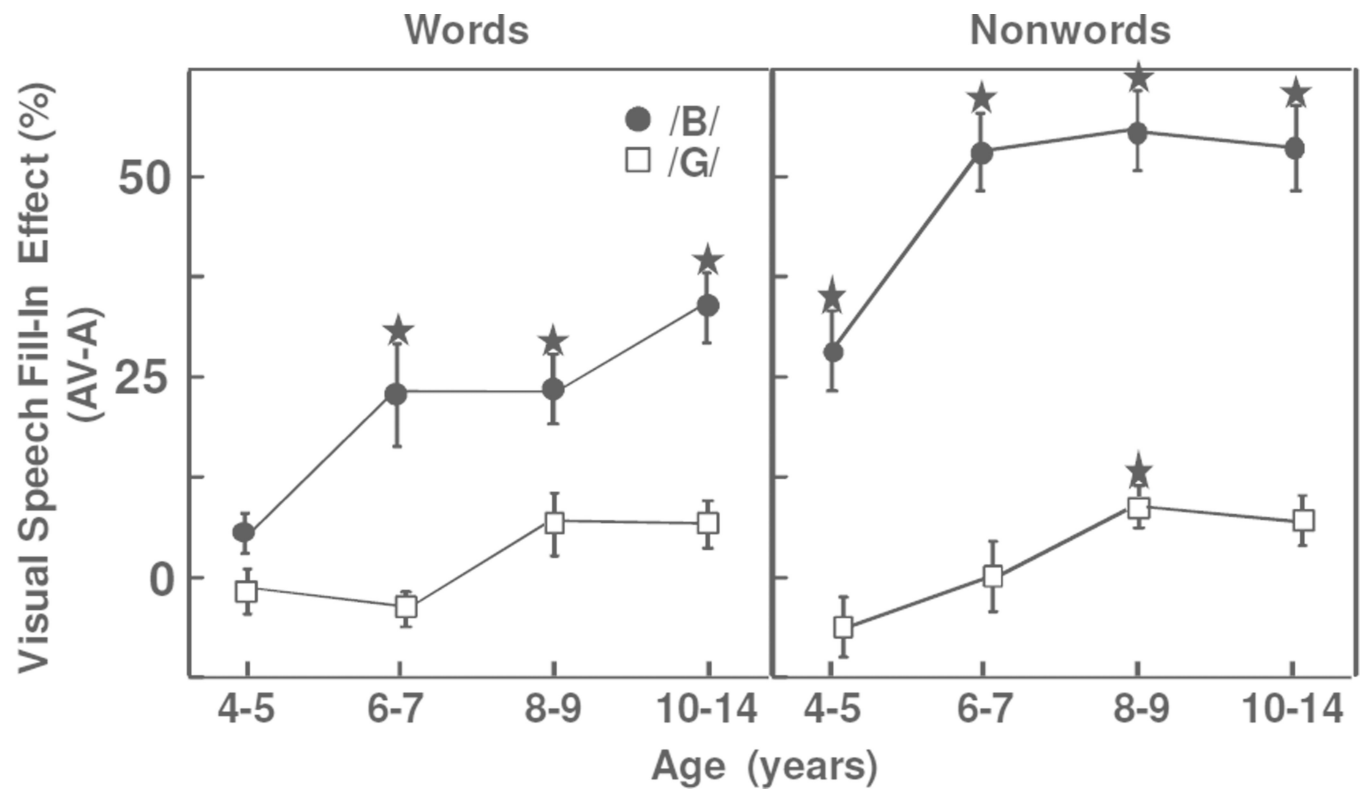


Figure 2.

Comparison of visual speech fill-in effect (difference in number of correct consonant onset responses for AV minus AO modes) for words and nonwords in the four age groups. The star indicates results that significantly differed from zero when controlling for multiple comparisons.

Table 1

Averages ages, vocabulary skills, and speechreading abilities (standard deviations in parentheses) in the four groups of children (N=92)

Measures	Age Groups (yrs)			
	4-5	6-7	8-9	10-14
Age	4;11	6;11	8;10	11;70
(years;months)	(0.73)	(0.67)	(0.68)	(1.38)
Vocabulary	120.48	116.75	121.82	121.35
(standard score)	(10.29)	(13.02)	(12.11)	(11.72)
Speechreading: Visual only [#]				
(percent correct)				
scored by words	4.86	9.97	17.72	24.59
	(5.35)	(8.09)	(13.72)	(12.43)
scored by word onsets [*]	46.06	53.73	66.14	72.86
	(21.30)	(20.53)	(14.43)	(13.07)

[#] Note. # AO results were at ceiling

^{*} Onsets were scored with visemes counted as correct (e.g., pat for bat)

Table 2

Multiple correlation coefficient and omnibus F for all of the variables considered simultaneously followed by the part correlation coefficients and the partial F statistics evaluating the variation in performance uniquely accounted for (after removing the influence of the other variables) by 1) age, vocabulary skills, or visual speechreading skills for onsets. The influence of visual speech was quantified by the 1) difference in the number of correct consonant onset responses for the AV - AO modes for the visual speech fill-in effect and 2) number of visually influenced responses for the McGurk effect.

Variables	Visual Speech Fill-In Effect /B/ Onsets			McGurk Effect		
	Multiple	Omnibus		Multiple	Omnibus	
	R	F	p	R	F	p
	.421*	6.16	.001	.407*	5.71	.001
	Part	Partial		Part	Partial	
	r	F	p	r	F	p
Age	.288*	8.64	.004	.000	0.05	ns
Vocabulary Skills	.261*	7.13	.009	.045	0.17	ns
Visual Speechreading (Onsets)	.032	0.03	ns	.355*	13.03	.001

Note. ns = not significant; df s = 3,88 for omnibus F and 1,88 for partial F