



Published in final edited form as:

Behav Res Methods. 2012 December ; 44(4): 1121–1128. doi:10.3758/s13428-012-0194-0.

Accuracy of Perceptual and Acoustic Methods for the Detection of Inspiratory Loci in Spontaneous Speech

Yu-Tsai Wang, Ph.D.,

School of Dentistry, National Yang-Ming University, Taipei, Taiwan

Ignatius S. B. Nip, Ph.D.,

School of Speech, Language, and Hearing Sciences, San Diego State University

Jordan R. Green, Ph.D.,

Department of Special Education and Communication Disorders, University of Nebraska-Lincoln

Ray D. Kent, Ph.D.,

Waisman Center, University of Wisconsin-Madison

Jane Finley Kent, M.S., and

Waisman Center, University of Wisconsin-Madison

Cara Ullman, M.S.

Department of Special Education and Communication Disorders, University of Nebraska-Lincoln

Abstract

The current study investigates the accuracy of perceptually and acoustically determined inspiratory loci in spontaneous speech for the purpose of identifying breath groups. Sixteen participants were asked to talk about simple topics in daily life at a comfortable speaking rate and loudness while connected to a pneumotach and audio microphone. The locations of inspiratory loci were determined based on the aerodynamic signal, which served as a reference for loci identified perceptually and acoustically. Signal detection theory was used to evaluate the accuracy of the methods. The results showed that the greatest accuracy in pause detection was achieved (1) perceptually based on the agreement between at least 2 of the 3 judges; (2) acoustically using a pause duration threshold of 300 ms. In general, the perceptually-based method was more accurate than was the acoustically-based method. Inconsistencies among perceptually-determined, acoustically-determined, and aerodynamically-determined inspiratory loci for spontaneous speech should be weighed in selecting a method of breath-group determination.

During speech, breathing patterns are constantly changing to balance the varying demands of an utterance with those of underlying homeostatic respiration. Among primates, humans appear unique in this refined and flexible capability for sound production (MacLarnon & Hewitt, 1999). During speech, the duration of inspiration typically only represents 9% to 19% of the full breath cycle (inspiration + expiration). (Loudon, Lee, & Holcomb, 1988). The characteristic respiratory pattern for speech (quick inspiration and a gradual and controlled expiration) inevitably imposes a breath-related structure on vocal output. This structure is commonly known as the breath group, a sequence of syllables or words produced on a single breath. Management of breath groups is one aspect of efficient and effective communication, for optimum vocal performance in both healthy and disordered speakers.

The location and degree of inspiration must be planned prior to the production of an utterance to ensure there is adequate aerodynamic power for conveying the linguistic properties of an utterance (Winkworth, Davis, Adams & Ellis, 1995).

Identification of breath groups is often a fundamental step in the analysis of recorded speech samples, especially for reading passages, dialogs and orations. Inspirations mark intervals of speech that can be subsequently examined for prosody and related variables. The usefulness of breath groups has been demonstrated in studies of (a) normal speech breathing (Hoit & Hixon, 1987; Mitchell, Hoit & Watson, 1996), (b) development of speech in infants (Nathani & Oller 2001), (c) design of speech technologies such as automatic speech recognition and text-to-speech synthesis (Ainsworth, 1973; Rieger, 2003), and (d) the assessment and treatment of speech disorders (Che, Wang, Hu, & Green, 2011; Huber & Darling, 2011; Yorkston, 1996). Common to these various applications is the need to identify groupings of syllables or words produced on a single breath, which is the inevitable respiratory imprint on spoken communication.

Breath groups have been determined in three ways: perceptually (by listening to the speech output), acoustically (usually by detecting pauses or silences that exceed a criterion threshold), and physiologically (typically by recording chest wall movements or the direction of airflow during speech). The physiologic method may be considered the gold standard; however, it is not always easily incorporated in studies of speech and cannot be used to analyze previously recorded samples (such as archival recordings) that did not employ physiological measures. Although many studies identify and evaluate breath groups perceptually and acoustically, the basic question about breath group studies using perceptual and acoustic methods is how well they correlate with physiologic analysis.

Perceptual determination is an indirect detection based on auditory judgments of speech features associated with the respiratory cycle (Bunton, Kent, & Rosenbek, 2000; Oller & Smith, 1977; Schlenck, Bettrich, & Willmes, 1993; Wang, Kent, Duffy & Thomas, 2005; Wozniak, Coelho, Duffy, & Liles, 1999). Both the perceptual and acoustic methods can be applied to previously recorded speech samples and can be accomplished with only modest investment in hardware or software. Although most studies have investigated breath groups using either of these indirect methods, the accuracy of these approaches has not been tested; the perceptual method is entirely subjective based on listeners' impressions; the acoustic method, requires the user to specify a minimum duration for an acceptable pauses. Therefore, silent portions in the speech signal, which may be pauses, but that do not exceed this criterion are not investigated (Campbell & Dollaghan, 1995; Green, Beukelman, & Ball, 2004; Walker, Archibald, Cherniak, & Fish, 1992; Yunusova, Weismer, Kent, & Rusche, 2005).

In contrast to the indirect methods, the physiologic determination directly detects inspiratory and expiratory events through either airflow (Wang, Green, Nip, Kent, Kent & Ullman, 2010b) or chestwall movements (Bunton, 2005; Forner & Hixon, 1977; Hammen & Yorkston, 1994; Hixon, Goldman & Mead, 1973; Hixon, Mead & Goldman, 1976; Hoit & Hixon, 1987; Hoit, Hixon, Watson, & Morgan, 1990; McFarland, 2001; Mitchell, Hoit, & Watson, 1996; Winkworth et al., 1995; Winkworth, Davis, Ellis, & Adams, 1994). Physiologic detection requires adequate instrumentation and may impose at least slight encumbrances on participants, such as the need to wear a face mask for oral airflow measures.

The present study is a follow-up to an earlier investigation of breath-group detection in a task of passage reading. Because accuracy of breath-group detection may be affected by the speaking task, it is necessary to examine the performance of different methods of detection

in at least spontaneous speech and passage reading, which have been primary tasks in the study of speech production. Studies have shown that these two speaking tasks are associated with somewhat different patterns in breath group structure (Wang, Green, Nip, Kent and Kent, 2010a).

Methods

Participants and Stimuli

Sixteen healthy adults (6 males, 10 females), ranging in age from 20 to 64 years ($m = 40$; $sd = 15$) participated in the study. All participants were native speakers of North American English and with no self-reported history of speech, language, or neurological disorders. Participants had normal or corrected hearing and vision. Participants were screened to ensure that they had adequate speech, language, and cognitive skills required to discuss simple topics regarding daily life. In addition to the 16 speakers, three individuals from the University of Wisconsin-Madison judged where inspiratory loci fell in each speaking sample based on auditory-perceptual cues in the audio recording.

Experimental protocol

Participants were seated and were instructed to hold a circumferentially-vented mask (Glottal Enterprises MA-1L) tightly against their faces. Expiratory and inspiratory airflows during the speaking tasks were recorded using a pneumotachograph (airflow) transducer (Biopac SS111A) that was coupled to the facemask. Previous research has demonstrated that facemasks do not significantly alter breathing patterns (Collyer, Davis, Collyer, & Davis, 2006). Although respiratory activity may be affected by the participants' use of facemasks in combination with the hand and arm muscle forces needed to hold the mask tightly against the face, participants in the current study were talking comfortably. Audio signals were recorded digitally at 48 kHz (16-bit quantization) using a professional microphone (Sennheiser), which was placed approximately 2 – 4 cm away from the vented mask. Participants were also video-recorded using a Canon XL-1s digital video recorder; however, only the audio signals were used for the analysis of breath group determination.

Participants were asked to talk about the following topics with a comfortable speaking rate and loudness in as much detail as possible: their family, activities in an average day, their favorite activities, what they do for enjoyment, and their plans for their future. The topics were presented on a large screen using an LCD projector. Participants were given time to formulate their responses to the topics before the recording was initiated to obtain reasonably organized and fluent spontaneous speech samples. Each response was required to be composed of at least 6 breath groups (as monitored by an airflow transducer).

Breath group determination

Aerodynamics—Data from the pneumotachometer and the simultaneous digital audio signal were recorded using Biopac Student Lab 3.6.7. The airflow signal was sampled at 1000 Hz and subsequently low-pass filtered ($F_{LP} = 500$ Hz). The resultant airflow signal was later used to visually identify actual inspiratory loci, represented by the upward peak in the airflow trace indicating inspiration (Figure 1) whereas a downward trend in the signal indicated expiration. On the rare occasions where there was uncertainty about the location of the inspiratory location, the first and the second authors examined the airflow traces in order to reach a consensus agreement on the inspiratory location.

Perception—Breath groups for the speech samples were determined perceptually by three judges at the University of Wisconsin-Madison. The judges were native English speakers trained on how to identify breath groups using known perceptual cues that signal the

production of inspiratory pauses. The judges for the determination of breath group were trained to learn how to determine the location of breath groups based on possible cues before performing their tasks. They were asked to listen to other conversation speech samples and to mark the points on their transcription sheets at which inspiration occurs. When the inspiration was not audible, the judges estimated the inhalation point based on auditory-perceptual impression and various acoustic cues, such as longer pause duration, f_0 declination, and longer phrase-final duration, which are fairly reliable indicators of pauses in normal speech and infant vocalization (Nathani & Oller, 2001; Oller & Lynch, 1992). The judges were also provided with a standard set of instructions explaining the task (see Appendix 1). In addition, the judges were allowed to listen to the speech samples repeatedly to ensure that they were confident in their determination on the breath group location. The procedures of breath group determination were as follows:

1. The speech samples were orthographically transcribed by a trained transcriptionist who did not serve as a judge in the determination of inspiratory loci.
2. Punctuations and upper-lower case distinctions (except for pronoun I, and proper names) were removed from the orthographic transcripts to prevent judges from analyzing breath groups based on punctuation and related visual cues in the transcript. Three spaces separated each word prevent judges from using word order to separate breath groups.
3. The speech samples prepared for the judges for the task of breath group determination were randomized for order of speaker using a table of random numbers.
4. The judges listened to the speech samples at normal loudness and marked perceived inspiratory loci on the transcripts. The judges were asked to make a best guess of the inhalation location based on their auditory-perceptual impressions when inspirations were not obvious. Therefore, these judgments could be based on multiple cues available to listeners, such as longer pause duration, f_0 declination, and longer phrase-final word or syllable duration. Judges were allowed to listen to the digitized speech samples repeatedly until they were satisfied with their determination of the breath group location.
5. The perceptual judgments of inspiratory loci were compared across each possible pairing of the three judges and across all the three judges to gauge the inter-judge reliability. Measurement reliability was defined as the number of points that the judges agreed upon an inspiratory location divided by the total number of perceptually determined inspiratory loci by the three judges.

Acoustics—A custom Matlab algorithm called Speech Pause Analysis or SPA (Green, Beukelman, & Ball, 2004) determined the acoustically-identified locations of the breath groups for the speech samples. The software required that a section of pausing to be identified manually to specify the minimum amplitude threshold for speech. The software also required specification of durational threshold values for the minimum pause and speech segment durations. For the current study, five pause duration thresholds were tested: 150 ms, 200 ms, 250 ms, 300 ms, and 350 ms. These were selected to cover the range of pause duration thresholds typically used in previous studies, for example, inspiratory loci have been defined as pauses greater than 150 msec (Yunusova, Weismer, Kent, & Rusche, 2005), 250 msec (Walker, Archibald, Chemiak, & Fich, 1992), or 300 msec (Campbell & Dollaghan, 1995). The minimum threshold for speech segment duration was held constant at 25 ms. Once these parameters were set, the acoustic waveform was rectified and then identified signal boundaries based on portions of the recording that fell below the signal amplitude threshold and above the specified minimum pause duration (e.g., 250 ms).

Portions that exceeded the minimum amplitude threshold were identified as speech. Adjacent speech regions were considered to be a single region if a pause region was less than the minimum pause duration. Finally, all the speech and pause regions in the speech samples were calculated by the algorithm.

Accuracy

The loci of inspiration determined by the airflow signal for all speakers were marked first. Inspiratory loci in the aerodynamic signal were taken as the true inspiratory events because they reflected the physiologic events. The aerodynamically-determined inspiratory loci were set to determine the accuracy of the perceptually-determined loci and acoustically-determined loci. Once inspiratory loci were determined using each of the three methods, the loci between conditions were compared. First, the number of perceptually or acoustically judged inspiratory loci was totaled. These loci were then compared to those identified using the aerodynamic signal. Loci identified perceptually and acoustically were then coded as a “true positive” when loci identified by the judges matched an inspiration identified by the aerodynamic method. Loci in which judges perceived an inspiration but were not indicated in the aerodynamic signal were coded as a “false positive”. Aerodynamically-determined loci that were not identified by judges were coded as a “miss”.

Statistical analysis

Signal detection analysis (MacMillan & Creelman, 1991) was used to evaluate which perceptually-based method and which pause threshold used in acoustic analysis yielded the most accurate results. Specifically, sensitivity as indicated by the true positive rate (TPR), the false positive rate (FPR) (1 - specificity), accuracy, and d-prime values were determined for each perceptual judgment and for each pause threshold.

Results

Accuracy

The total number of inspirations determined from the airflow signal for all speakers was 1106. The number of pauses greater than 150 ms detected by the SPA algorithm was 2281, which was considered the total number of potential inspirations for judges to make their decisions.

Perception—The total number of inspiratory loci determined perceptually by the three judges was 1177. The number of inspiratory locations determined individually by judge 1 (J1), judge 2 (J2), and judge 3 (J3) were 1088, 1094, and 1054, respectively. The number of consistent judgments between at least two of the three judges (i.e., J1J2, J1J3, J2J3, or J1J2J3), was 1080. The number of consistent judgments across all three judges was 979. The highest inter-judge reliability between two of the three judges was 0.92 (1080/1177). The inter-judge reliability across all the three judges was 0.83 (979/1177).

Referenced to the 1106 actual inspiratory loci, J1 correctly identified 1066, missed 42, and added 22 (false alarm). J2 correctly identified 1065, missed 43, and added 29. J3 correctly identified 1010, missed 98, and added 44. The loci that were consistent between at least two of the three judges were 1068 correctly identified, 40 missed, and 12 added. The loci that were consistent across all three judges were 976 correctly identified, 132 missed, and 3 added.

Table 1 shows the sensitivity, specificity, accuracy, and d-prime data for the perceptual judgments. Inspiratory locations were perceived correctly (true positive rate, TPR) about 95% of the time on average, and the false alarm rate (false positive rate, FPR) varied among

the three judges. J1 had the highest sensitivity, specificity, accuracy, and d-prime. When the decision was based on the agreement across all 3 judges, the specificity was increased substantially but the sensitivity and accuracy were decreased to 88%. However, when the decision was based on agreement between at least 2 of the 3 judges, the specificity was near 99%, and the sensitivity, accuracy, and d-prime were all at their highest. Overall, the best discrimination of the perceptual judgment of inspiratory loci in spontaneous speech was based on the consistency between at least 2 of the 3 judges. However, as shown in the Receiver Operating Characteristic (ROC) curve of Figure 2, the separate results for the 3 judges are clustered rather tightly.

Acoustics—The number of pauses acoustically determined by the SPA algorithm is given in parentheses in the following summary for the five different pause thresholds: 150 ms (2281), 200 ms (1864), 250 ms (1657), 300 ms (1513), and 350 ms (1406). Table 2 shows the sensitivity, specificity, accuracy, and d-prime data for the SPA algorithm results. Figure 2 shows the ROC for the combined perceptual and acoustic results. The TPR (sensitivity) values of the five different pause thresholds were all above 98%, but the FPR differs greatly among different threshold values, with smaller thresholds resulting in greater FPRs. The smaller thresholds had near perfect sensitivity but very poor specificity; consequently, lower accuracy. Thus, in terms of the d-prime value, the SPA acoustically determined inspiratory loci of 300 ms threshold had the best performance.

Compared with the actual inspiratory locations determined by the aerodynamic signal, the perceptually determined method with the best performance had smaller TPR and FPR, but larger accuracy and d-prime than those of the acoustically determined method for this spontaneous speech task (Table 2). Moreover, the sensitivity values of the five different pause thresholds were all higher than those of perceptual judgments, but the specificity values were much larger and varied widely (Table 2). Consequently, based on accuracy and d-prime analysis, the performance of the perceptually-based breath determination of breath groups is judged to be better than that of the acoustic method of pause detection.

Discussion

The current study indicates that (1) The greatest accuracy in the perceptual detection of inspiratory loci was achieved with agreement between 2 of the 3 judges; (2) The most accurate pause duration threshold used for the acoustic detection of inspiratory loci was 300 ms; and (3) The perceptual method of breath group determination was more accurate than the acoustically-based determination of pause duration.

For the perceptual approach, the criterion of agreement between 2 of the 3 judges yielded the highest TPR, accuracy (0.977), and d-prime (4.116). This approach had approximately 1.75 % (40/2281) false negatives and 0.53 % (12/2281) false positives. Apparently, the more stringent criterion of consistency across all 3 judges led to an increase of false negatives that was much larger than the decrease of false positives, thereby reducing both accuracy and d-prime. In contrast, the most accurate approach for detecting inspiratory loci based on listening in a reading task (Wang et al., 2010b) was agreement across all three judges, which achieved an accuracy of 0.902, a d-prime of 4.140, and a small number of both false negatives (approximately 10 %) and false positives (0%). The accuracy of the perceptual approach was better for spontaneous speech in the present study than it was for passage reading in the study by Wang et al. (2010b). The differences between spontaneous speech and reading are likely explained by differences in breath group structure, as discussed in Wang, Green, Nip, Kent and Kent (2010a). Breath groups had longer durations for spontaneous speech as compared with reading. In addition, inspiratory pauses for

spontaneous speech are more likely to fall in grammatically-inappropriate locations, potentially making the inspirations to be more perceptually salient to the judges.

Using acoustic algorithms to identify inspiratory loci, the optimal threshold of pause detection in the present study was 300 ms, which achieved an accuracy of 0.817 and a d-prime value of 2.994. With this threshold, the false negative rate is 0.2% (5/2281), but the false positive rate is much higher, approximately 18% (412/2281). Wang, Green, Nip, Kent, Kent and Ullman (2010b) reported that the most accurate pause duration threshold for detecting inspiratory loci in the reading task was 250, which achieved an accuracy of 0.895, a d-prime of 3.561, a zero rate of false negatives and an approximately 10% rate of false positives. Task effects between reading and spontaneous speech occurred for the acoustic method, much as they did for the perceptual method. The accuracy and d-prime values in spontaneous speech were lower than those in reading. Furthermore, the false negative rate and false positive rate in spontaneous speech were both raised when compared with reading. Consequently, the acoustically-determined method in spontaneous speech performed more poorly than for reading, which is likely related to the task differences in the breath group structure and perhaps in cognitive-linguistic load.

Because the minimum inter-breath-group pause in reading for healthy speakers is 250 ms (Wang et al., 2010a), the 150 ms and 200 ms thresholds produced no false negatives but many false positives, which lowered their accuracy. In contrast, with thresholds above 200 ms, the decrease in the number of false positive was substantially more than the increase of the number of false negatives, which increases the accuracy. Generally speaking, the false positive rate differed among different pause thresholds, indicating that the selection of the pause threshold is very sensitive to the detection of false positives in spontaneous speech. Because the spontaneous speech samples in the current study were produced fluently by healthy adults who were familiar with the topics to be addressed, there was negligible occurrence of prolonged cognitive hesitations or articulatory or speech errors. Therefore, the current findings may not apply to speech produced by talkers with neurological or other impairments, whose speech might be characterized by either faster or slower speaking rate and with more pauses of long durations unrelated to inspiration. A threshold of 300 ms might potentially be either too short for individuals who speak significantly slower or too long for speakers with faster than typical speaking rates.

Taking together the present results and those of Wang, Green, Nip, Kent, Kent and Ullman (2010b), it can be concluded that for both spontaneous speech and passage reading, the perceptual method of breath group determination is more accurate than the acoustic method based on pause duration. The ability of listeners to identify breath groups is no doubt aided by their knowledge that speech is typically produced on a prolonged expiratory phase. Simple acoustic measurements of pauses are naive to this expectation, which is one reason by perceptual assessment can be more accurate than acoustic pause detection. The larger d-prime obtained for the perceptual approach may indicate that listeners are sensitive to many cues beyond pause duration. Factors related to physiologic needs, cognitive demands, and linguistic accommodations to affect the locations of inspirations and the durations of inter-breath-group pauses are possibly perceptible by human ears. Perceptual cues for inspiration include the occurrence of pauses at a major constituent boundary, anacrusis, final syllable lengthening, and final syllable pitch movement (Wozniak, et al., 1999). Some of these factors could be included in an elaborated acoustic method that relies on more than just pause duration.

The choice of method for breath-group determination should be based on a consideration of the risk-benefit ratio. If errors cannot be tolerated, then physiologic methods are preferred if not mandatory. But if this is not possible (as in the analysis of archived audio signals), then

the choice between perceptual and acoustic methods should weigh the risk of greater errors (likely to occur with the acoustic method) against the relative costs (in terms of both analysis time and technology). As shown in Figure 2, the results for any one judge in the perceptual method were more accurate than any of the pause duration thresholds used in the acoustic study. Perceptual determination appears to be a better choice based on accuracy alone. Of course, these findings pertain to studies interested in identifying only inspiratory pauses, and not those located at phrase and word boundaries; the high false positive rates obtained for the acoustic methods suggest that this approach may be well suited for this purpose although additional research is needed. If it is desired to examine the relationship between breath groups and linguistic structures, then preparation of a transcript is necessary for any method of breath-group determination. Finally, it should be recognized that the present results and those of Wang, Green, Nip, Kent, Kent and Ullman (2010b) pertain to healthy adult speakers. Generalization of the results to younger or older speakers or to speakers with disorders should be done with caution.

Acknowledgments

This work was supported in part by Research Grant number 5 R01 DC00319, R01 DC009890, and R01 DC006463 from the National Institute on Deafness and Other Communication Disorders (NIDCD-NIH), and NSC 100-2410-H-010-005-MY2 from National Science Council, Taiwan. Additional support was provided by the Barkley Trust, University of Nebraska-Lincoln, Department of Special Education and Communication Disorders. Some of the data were presented in a poster session at the 5th International Conference on Speech Motor Control, Nijmegen, 2006. We would like to acknowledge Hsiu-Jung Lu and Yi-Chin Lu for data processing.

References

- Ainsworth W. A system for converting English text into speech. *IEEE Transactions on Audio and Electroacoustics*. 1973; 21:288–290.
- Bunton K. Patterns of lung volume use during an extemporaneous speech task in persons with Parkinson disease. *Journal of Communication Disorders*. 2005; 38:331–348. [PubMed: 15963334]
- Bunton K, Kent RD, Rosenbek JC. Perceptuo-acoustic assessment of prosodic impairment in dysarthria. *Clinical Linguistics and Phonetics*. 2000; 14:13–24. [PubMed: 22091695]
- Campbell TF, Dollaghan CA. Speaking rate, articulatory speed, and linguistic processing in children and adolescents with severe traumatic brain injury. *Journal of Speech and Hearing Research*. 1995; 38:864–875. [PubMed: 7474979]
- Che WC, Wang YT, Lu HJ, Green JR. Respiratory changes during reading in Mandarin-speaking adolescents with prelingual hearing impairment. *Folia Phoniatica et Logopaedica*. 2011; 63:275–280. [PubMed: 21372590]
- Collyer S, Davis PJ, Collyer S, Davis PJ. Effect of facemask use on respiratory patterns of women in speech and singing. *Journal of Speech Language and Hearing Research*. 2006; 49:412–423.
- Forner LL, Hixon TJ. Respiratory kinematics in profoundly hearing-impaired speakers. *Journal of Speech and Hearing Research*. 1977; 20:373–408. [PubMed: 895106]
- Green JR, Beukelman DR, Ball LJ. Algorithmic estimation of pauses in extended speech samples of dysarthric and typical speech. *Journal of Medical Speech-Language Pathology*. 2004; 12:149–154. [PubMed: 20628555]
- Hammen VL, Yorkston KM. Respiratory patterning and variability in dysarthric speech. *Journal of Medical Speech-Language Pathology*. 1994; 2:253–261.
- Hixon TJ, Goldman MD, Mead J. Kinematics of the chest wall during speech production: volume displacements of the rib cage, abdomen, and lung. *Journal of Speech and Hearing Research*. 1973; 16:78–115. [PubMed: 4267384]
- Hixon TJ, Mead J, Goldman MD. Dynamics of the chest wall during speech production: function of the thorax, rib cage, diaphragm, and abdomen. *Journal of Speech and Hearing Research*. 1976; 19:297–356. [PubMed: 135885]
- Hoit JD, Hixon TJ. Age and speech breathing. *Journal of Speech and Hearing Research*. 1987; 30:351–366. [PubMed: 3669642]

- Hoit JD, Hixon TJ, Watson PJ, Morgan WJ. Speech breathing in children and adolescents. *Journal of Speech and Hearing Research*. 1990; 33:51–69. [PubMed: 2314085]
- Huber JE, Darling M. Effect of Parkinson's disease on the production of structured and unstructured speaking tasks: respiratory physiologic and linguistic considerations. *Journal of Speech, Language, and Hearing Research*. 2011; 54:33–46.
- Loudon RG, Lee L, Holcomb BJ. Volumes and breathing patterns during speech in healthy and asthmatic subjects. *Journal of Speech and Hearing Research*. 1988; 31:219–227. [PubMed: 3294505]
- MacLarnon AM, Hewitt GP. The evolution of human speech: the role of enhanced breathing control. *American Journal of Physical Anthropology*. 1999; 109:341–363. [PubMed: 10407464]
- Macmillan, NA.; Creelman, CD. *Detection Theory: A user's guide*. New York: Cambridge University Press; 1991.
- McFarland DH. Respiratory markers of conversational interaction. *Journal of Speech Language and Hearing Research*. 2001; 44:128–143.
- Mitchell HL, Hoit JD, Watson PJ. Cognitive-linguistic demands and speech breathing. *Journal of Speech and Hearing Research*. 1996; 39:93–104. [PubMed: 8820701]
- Nathani S, Oller DK. Beyond ba-ba and gu-gu: Challenges and strategies in coding infant vocalizations. *Behavior Research Methods, Instruments, and Computers*. 2001; 33:321–330.
- Oller, DK.; Lynch, MP. Infant vocalizations and innovations in infraphonology: Toward a broader theory of development and disorders. In: Ferguson, C.; Menn, L.; Stoel-Gammon, C., editors. *Phonological development: Models, research, implications*. Parkton, MD: York Press; 1992. p. 509-536.
- Oller DK, Smith BL. Effect of final-syllable position on vowel duration in infant babbling. *Journal of the Acoustical Society of America*. 1977; 62:994–997. [PubMed: 908792]
- Rieger JM. The effect of automatic speech recognition systems on speaking workload and task efficiency. *Disability and Rehabilitation*. 2003; 25:224–235. [PubMed: 12623631]
- Schlenck KJ, Bettrich R, Willmes K. Aspects of disturbed prosody in dysarthria. *Clinical Linguistics & Phonetics*. 1993; 7:119–128.
- Walker JF, Archibald LM, Cherniak SR, Fish VG. Articulation rate in 3- and 5-year-old children. *Journal of Speech & Hearing Research*. 1992; 35:4–13. [PubMed: 1735975]
- Wang YT, Green JR, Nip ISB, Kent RD, Kent JF. Breath Group Analysis for Reading and Spontaneous Speech in Healthy Adults. *Folia Phoniatrica et Logopaedica*. 2010a; 62:297–302. [PubMed: 20588052]
- Wang YT, Green JR, Nip ISB, Kent RD, Kent JF, Ullman C. Accuracy of perceptually based and acoustically based inspiratory loci in reading. *Behavior Research Methods*. 2010b; 42:791–797. [PubMed: 20805602]
- Wang YT, Kent RD, Duffy JR, Thomas JE. Dysarthria in traumatic brain injury: a breath group and intonational analysis. *Folia Phoniatrica et Logopaedica*. 2005; 57:59–89.
- Winkworth AL, Davis PJ, Adams RD, Ellis E. Breathing patterns during spontaneous speech. *Journal of Speech and Hearing Research*. 1995; 38:124–144. [PubMed: 7731204]
- Winkworth AL, Davis PJ, Ellis E, Adams RD. Variability and consistency in speech breathing during reading: lung volumes, speech intensity, and linguistic factors. *Journal of Speech and Hearing Research*. 1994; 37:535–556. [PubMed: 8084185]
- Wozniak RJ, Coelho CA, Duffy RJ, Liles BZ. Intonation unit analysis of conversational discourse in closed head injury. *Brain Injury*. 1999; 13:191–203. [PubMed: 10081600]
- Yorkston K. Treatment efficacy: dysarthria. *Journal of Speech and Hearing Research*. 1996; 39:546–557. [PubMed: 8783133]
- Yunusova Y, Weismer G, Kent RD, Rusche NM. Breath-group intelligibility in dysarthria: Characteristics and underlying correlates. *Journal of Speech, Language, and Hearing Research*. 2005; 48:1294–1310.

Appendix 1: The instruction of breath group determination for conversational speech samples

You will be provided with a transcription of the conversational speech samples without punctuations for each speaker in the present study. The task is to mark the points at which speakers stop for a breath. When you identify this point, place a mark on the corresponding location on the transcript. Make your best guess as to where the speaker stops to take a breath. Sometimes you can hear an expiration and/or inspiration, but in other cases you may have to make the judgment based on other cues, such as longer pause duration, f_0 declination, and longer phrase-final duration. In this task, you can listen to the sound files repeatedly before you are confident in your determination on the breath group location. Do you have any questions?

\$watermark-text

\$watermark-text

\$watermark-text

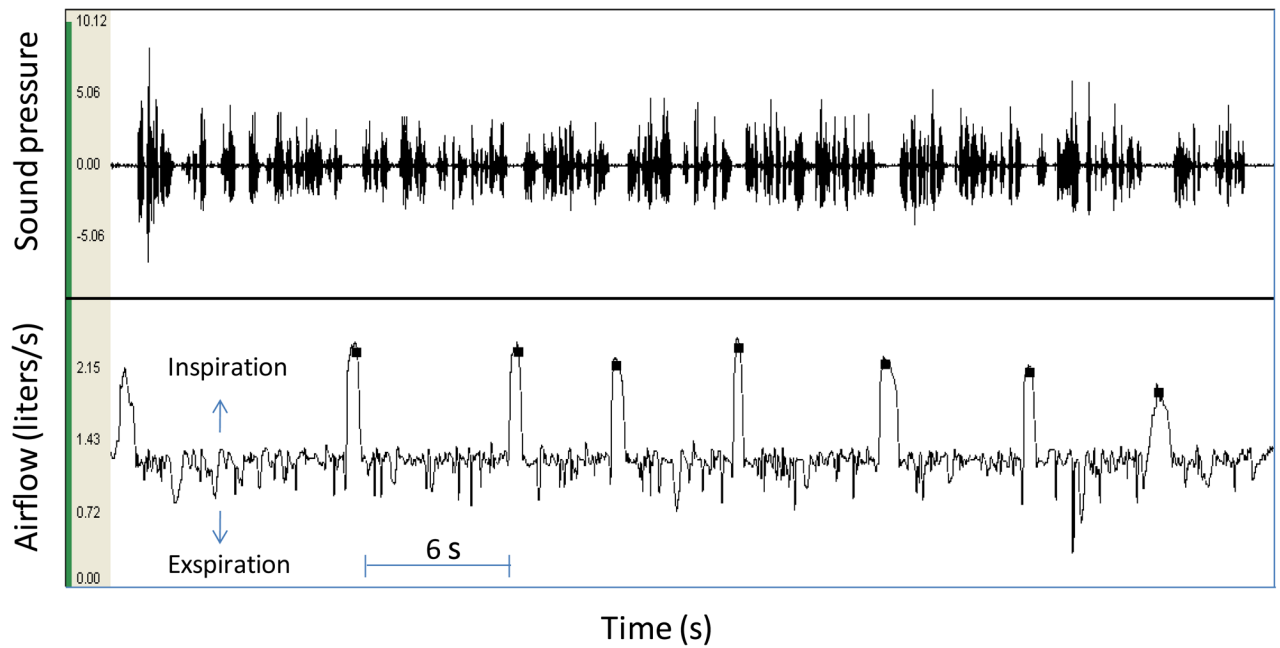


Figure 1.

A demonstration of the locations of inspiration indicated by the dots for the recorded speech sample based on the aerodynamic signal (the lower panel). The upper panel is the corresponding sound pressure of the acoustic signal. The arrows indicate the direction of airflow.

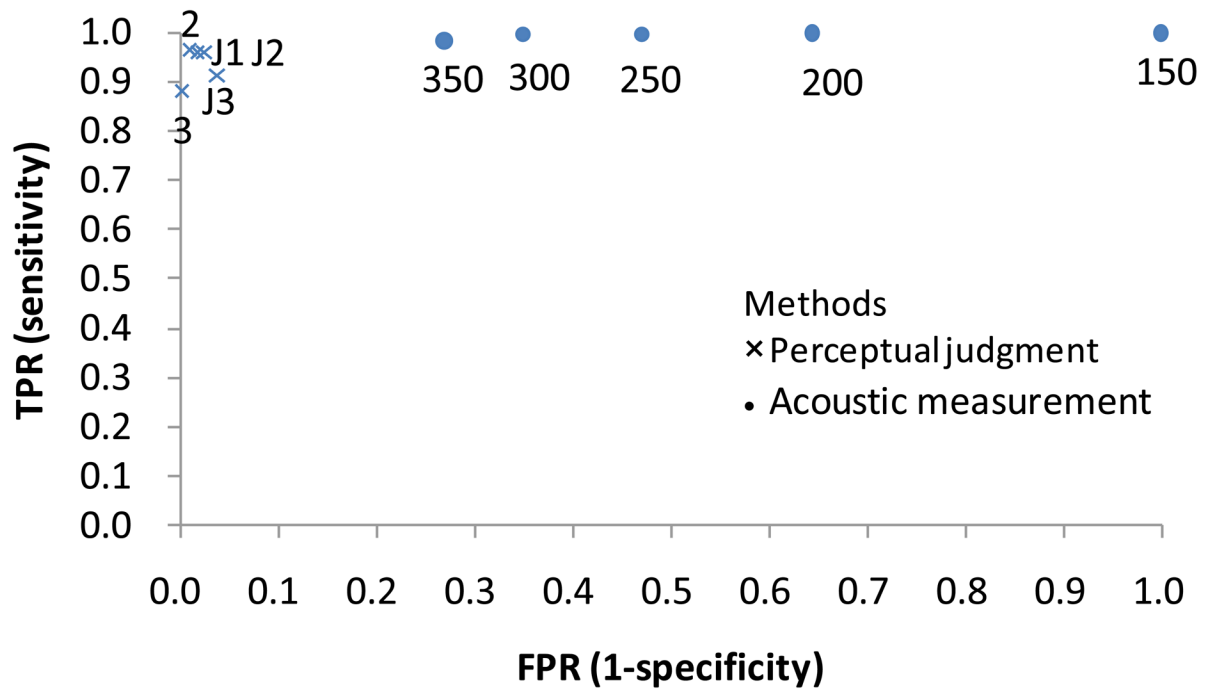


Figure 2. Receiving operator characteristic (ROC) curve for the perceptual and acoustic methods of breath-group determination. Perceptual results are shown for each judge and agreements between 2 judges (2) and 3 judges (3). Acoustic results are shown for various thresholds of pause duration.

Table 1

The sensitivity, specificity, accuracy, and d-prime data of perceptual judgments determined by judge (J1, J2, J3), by the consistency of at least 2 of the 3 judges, and by the consistency of all the 3 judges. True positive rate (TPR) refers to sensitivity, whereas false positive rate (FPR) refers to 1 - specificity.

Judge(s)	Inspiratory location		TPR	FPR	Accuracy	d-prime	beta (ratio)
	Yes	No					
J1	Yes	1066	0.9621	0.0188	0.972	3.856	1.799
	No	42	1151				
J2	Yes	1065	0.9612	0.0247	0.968	3.729	1.452
	No	43	1144				
J3	Yes	1010	0.9116	0.0375	0.938	3.131	1.960
	No	98	1129				
2	Yes	1068	0.9639	0.0102	0.977	4.116	2.915
	No	40	1161				
3	Yes	976	0.8809	0.0026	0.941	3.979	25.122
	No	132	1170				

Table 2

The sensitivity, specificity, accuracy, and d-prime data of acoustically determined by SPA algorithm. True positive rate (TPR) refers to sensitivity, whereas false positive rate (FPR) refers to 1 - specificity.

Threshold (ms)	Inspiratory location		TPR	FPR	Accuracy	d-prime	beta (ratio)
	Yes	No					
150	Yes	1106	0.9995	0.9996	0.485	-0.017	1.056
	No	0					
200	Yes	1106	0.9995	0.6451	0.668	2.947	0.004
	No	0	417				
250	Yes	1104	0.9982	0.4706	0.757	2.983	0.015
	No	2	622				
300	Yes	1101	0.9955	0.3506	0.817	2.994	0.036
	No	5	763				
350	Yes	1089	0.9846	0.2698	0.854	2.774	0.117
	No	17	858				