

A JOINT ACOUSTIC-ARTICULATORY STUDY OF NASAL SPECTRAL REDUCTION IN READ VERSUS SPONTANEOUS SPEAKING STYLES

Vikram Ramanarayanan^{*}, Dani Byrd[^], Louis Goldstein[^] and Shrikanth Narayanan^{*^}

^{*}Signal Analysis and Interpretation Laboratory, Ming Hsieh Department of Electrical Engineering,
[^]*Department of Linguistics*; University of Southern California, Los Angeles, CA-90089-0899

ABSTRACT

Data on nasal articulation obtained through real-time magnetic resonance imaging (MRI) are jointly used with acoustic analyses of the speech signal to analyze nasal production differences in read and spontaneous speech, especially focusing on the details of acoustic spectral reduction. In this exploratory study, vowel-nasal-vowel (VNV) segments from one speaker were examined and measures corresponding to the speed of the velum and spectral center-of-gravity of the nasal were extracted. It is observed that lower velum speeds in spontaneous nasal production could result in more vowel nasalization and hence a ‘lowering’ of the center-of-gravity of the acoustic spectrum. Such an analysis has implications for understanding speech planning and for informing design of automatic speech analysis algorithms.

Index Terms— speech production, real-time MRI, nasals, vocal tract, image motion analysis, read speech, spontaneous speech, spectral reduction.

1. INTRODUCTION

A joint consideration of direct articulatory and acoustic data of read and spontaneous speech styles can provide an improved understanding of the underlying differences in their production. Such investigations can be of potential use for efforts in both speech science and speech technology. This paper considers such a comparative study focusing on the details of the nasal consonant in read and spontaneous speaking styles. In terms of speech analysis, it can help us understand the nature of coupling between the oral tract and nasal tract such as (i) across durational differences of phones observed in the speech signal, (ii) when speech flow is sudden (unplanned) vs. smooth (planned), and (iii) during different phases of the human respiratory cycle. In terms of technology applications it can help inform improved speech representations such as for production inspired automatic speech recognition (ASR).

In this paper, we focus on investigating a specific aspect of read and spontaneous speech production differences--acoustic “spectral reduction”--which is characterized by

the reduction in spectral and durational distinctions between sounds as the speaking style becomes more informal, or the stress on the syllable is reduced. van son and Pols [1] report how intervocalic energy differences for nasal consonants are higher in spontaneous speech than read speech, reversing a trend which is generally seen for all other consonants; while Nakamura et al. [2] give explicit evidence for *acoustic* reduction of the MFCC spectrum for most consonants of English as the speech becomes more and more spontaneous causing a reduction in ASR performance. By investigating how these two observations are related using direct articulatory data in conjunction with the acoustic speech signal, we hope to obtain insights into this behavior from an from an articulatory perspective.

Varying degrees of coarticulation, velum speed and timing are some of the important factors that must be taken into consideration while examining the articulatory causes of spectral reduction in nasals. Byrd et al. [3] have examined the timing effects of syllable structure and stress on the segment-internal coordination of nasals and found the predominance of in-phase relations of oral/nasal gestures during onsets, and anti-phase during codas. Moll and Daniloff [4] also report extensive anticipatory coarticulation of the velum movement toward velopharyngeal opening in CVN and CVVN sequences. In this work, since we are looking for global differences between read and spontaneous speech styles, we will not consider explicitly the effects of syllable position.

The specific research questions we would like to pose with respect to spectral reduction phenomena here are: how is the timing and speed of the velum in read speech different from that in spontaneous speech, and further, how are their acoustic consequences different?

The recent advances in real-time mid-sagittal magnetic resonance imaging (MRI) offer an excellent tool to investigate solution this problem, since it allows us to capture the full vocal tract during speech production and quantify the ‘choreography’ of the articulators [5], making it an ideal technique to compare articulation during read and spontaneous speech production.

The paper is organized as follows: Section 2 details methods of MR data acquisition and reconstruction. Section 3 describes the analysis carried out on the acoustic signal as well as a method to characterize the midsagittal profile of the vocal tract. Finally, Section 4 discusses our results and summarizes directions for future work.

2. DATA

One native speaker (female) of American English was engaged in a simple dialog on topics of general nature (e.g.,: “what music do you listen to...”, “tell me more about your favorite cuisine ...”, etc.) while she was lying supine inside the MR scanner. For each of the speech “turn”, audio responses and MRI videos of speech during the vocal tract articulation were recorded for 30 seconds and time-synchronized. The same speaker was also recorded/imaged while reading some of the TIMIT sentences and the rainbow passage during a separate scan session. Further details regarding the recording and imaging setup can be found in [5,6]. Midsagittal real-time MR images of the vocal tract were acquired with a repetition time of TR=6.5ms on a GE Signa 1.5T scanner with a 13 interleaf spiral gradient echo pulse sequence. The slice thickness was approximately 3mm. A sliding window reconstruction at a rate of 22.4 frames per second was employed. The field of view was adjusted depending on the subject’s head size.

3. ANALYSES

It is important to note that all analysis considered here in the articulatory domain were carried out based on phenomena observed in the acoustic-domain, since most baseline scientific analyses have been carried out on the more widely-available speech signal. In this section, we first describe the acoustic domain pre-processing and analysis, such as the segmentation of the speech signal into phones so that context-specific (VCV) analyses can be performed and extraction of a spectral centroid measure. Then we explain how relevant features were extracted from the MRI videos using automatically-determined air-tissue boundary information.

The SONIC speech recognizer [7] was used to perform a first-pass automatic forced alignment of the recorded audio data to the phone sequence. However, the background noise especially due to the cryogenic pump in the MRI scan room caused misalignment of some phones/groups of phones, which warranted a second-pass manual correction of these alignments.

Spectral centre-of-gravity (CoG) values were computed for each vowel-nasal-vowel (VNV) sequence extracted automatically from the segmented speech data band-limited to 2.5 kHz (for both read and spontaneous

speech). This helps us approximately obtain the frequency about which maximum spectral energy is distributed without unwanted averaging effects due to noisy sample values at higher frequencies up to 10kHz (since the audio is sampled at 20kHz). This was computed as follows [1]:

$$CoG = \frac{\sum_i f_i E_i}{\sum_i E_i}$$

where f_i is the center frequency of each FFT band of frequencies, and E_i is the spectral power corresponding to each frequency.

3.1. Contour Extraction

The air-tissue boundary of the articulatory structures was automatically extracted using an algorithm that hierarchically optimizes the observed image data fit to an anatomically informed object model using a gradient descent procedure [8]. The object model is chosen such that different regions of interest such as the palate, tongue, velum etc. are each defined by a dedicated region (see Figure 1).

3.2. Velum Speed Measure

The velum contour *alone* was used to create a ‘mask’ image corresponding to each image of the MRI video sequence, with all pixels inside the contour rendered white and the rest, black. Then by taking the absolute value of the difference between successive masks, we can obtain a measure of how fast the velum is moving. For a detailed description of this process, see [9]. In addition, the velum opening for each MR image was computed as the minimum distance between the velum and pharyngeal wall contours.

3.3. Vocal Tract Area Descriptors (VTADs) Extraction

In the following paragraph, the extraction of vocal tract variables, including the lip aperture (LA), tongue tip constriction degree (TTCD), tongue dorsum constriction degree (TDCD), tongue root constriction degree (TRCD) and velic aperture (VEL), is described. For each image in the MRI video sequence, LA is computed as the minimum distance between the upper lip and lower lip contour segments. VEL is computed as the minimum distance between the velum and pharyngeal wall contours. In order to extract the tongue-related tract variables (TTCD, TDCD and TRCD), especially for those frames where the articulator in question is not a critical one, the main problem is defining the point on the palate with respect to which we can measure the constriction degree for that articulator. This problem can be alleviated by using frames in which an articulator is critical in order to define a set of possible ‘palate constriction locations’, which can then be in computing the constriction degrees for that articulator for all other frames.

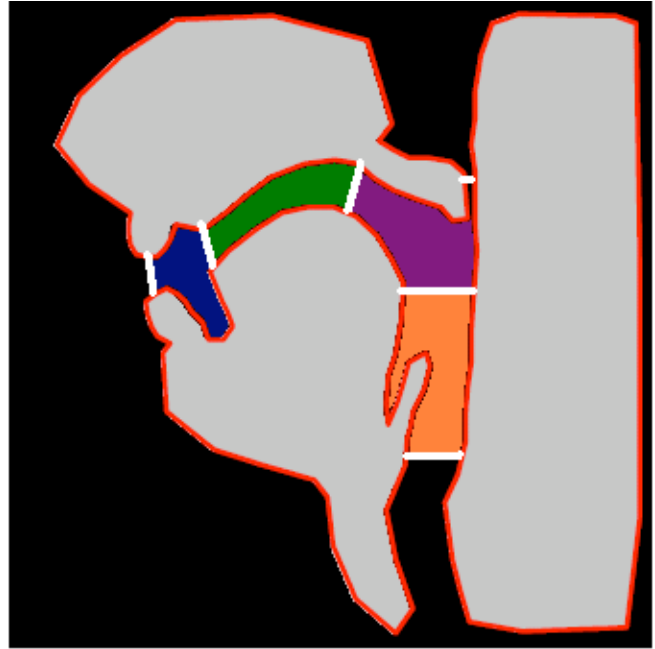
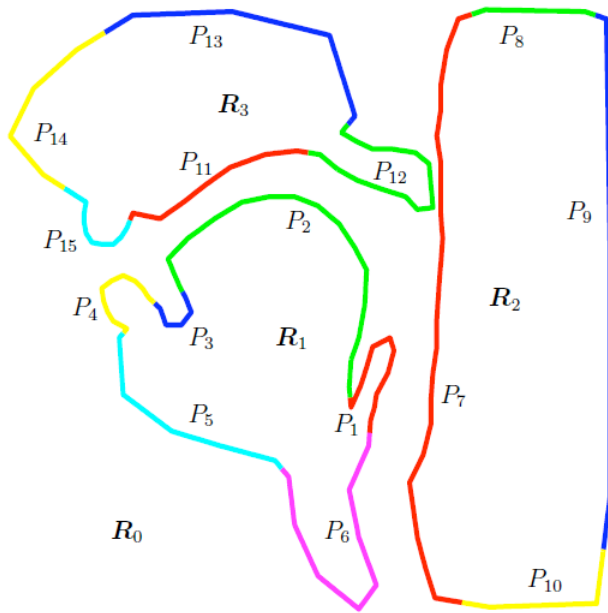


Figure 1 (Left) Contour outlines extracted for each image of the vocal tract. Note the template definition such that each articulator is described by a separate contour. (Right) A schematic depicting the concept of vocal tract area descriptors (adapted from [8]).

For example, in order to compute the TTCD for a vowel /a/, in which the tongue tip is not critical, we use the constriction location (on the palate) of the tongue tip constriction for all /t/, /d/ frames, where the tongue tip is a critical articulator and use the mean of this point cloud as the point on the palate from which to measure minimum distance to the tongue contour. We find that choosing /t/ and /d/ frames, /k/ and /g/ frames, and /a/ and /r/ frames as critical frames for the tongue tip, tongue dorsum and tongue root respectively works well in practice. Finally, the lowermost boundary of the vocal tract area for our purposes is computed as the minimum distance between the root of the epiglottis and pharyngeal wall contour (please see Figure 1). However, due to poor signal-to-noise ratio of images in this region, this is not always robust.

Once these tract variables are computed, we can then use them to partition the vocal tract midsagittal cross-sectional area, into the area between the LA and TTCD (which we call A1 (blue in Fig)), the area between the TTCD and TDCD (or A2 (green)), the area between the TDCD and TRCD (or A3 (purple)), and the area below the TRCD as A4 (orange) (Note that we are not using A4 for conclusive analyses, due to the reason described above). Once these areas are obtained, we can formalize the differences in vocal tract shaping more concretely.

4. RESULTS AND DISCUSSION

4.1. Results for all nasals in general

A one-way parametric¹ analysis of variance (ANOVA) with post-hoc Tukey test was used to test the hypothesis that the means of the z-scores of spectral CoG values calculated for each VNV instance in the read and spontaneous speech utterances were the same. All nasals (/m/,/n/,/ng/) in spontaneous speech were found to have a significantly lower ($p \leq 0.05$) spectral centroid as compared to those in read speech, which agrees with the van son and Pols study [1] which asserts that spectral slope and therefore CoG frequency is determined by the speech effort, which is generally higher in the case of read speech. It is, however, difficult to tease apart how much of the observed effect is individually due to the source and filter.

One-way parametric ANOVAs with post-hoc Tukey tests were also used to test the hypothesis that the means of the z-scores of the areas defined in Section 3.3 for each VNV instance in the read and spontaneous speech utterances were the same. The average area of the A2 vocal tract region (defined between TTCD and TDCD) for each nasal instance was found to be significantly *higher* ($p=0$) for read nasals as opposed to spontaneously produced nasals. In addition, the A1 region (defined between the LA and TTCD) was found to be significantly *lower* ($p \leq 0.05$) for

¹ P-P plots and Kolmogorov-Smirnov tests ($p=0.05$) were used to ascertain parametricity wherever required.

read nasals as opposed to spontaneously produced nasals. However, these observations could also be due to a more anterior constriction location (No significant differences were however found). Vowel coarticulation effects are also a potential confound if the /i/.../i/ vowel contexts force the consonant to be shaped more /i/-like. Repeating these experiments for different vowel contexts showed that the reported phenomenon was still significant irrespective of the underlying vowel context. This provides another possible reason for the concentration of more energy around lower frequencies (lower CoG) in spontaneous speech beyond those arising from a potential reduced glottal effort.

In order to see whether the width of the velum opening contributed in any significant way to our observations, z-scores of the maximum width of velum opening for each VNV instance were computed for both read and spontaneous speech. There were no significant differences ($p=0.05$) observed between the means of these 2 samples, suggesting that the extent to which the velum can open is similar in both read and spontaneous speech for our speaker. The length of nasal duration for each VNV sequence was also computed for read and spontaneous speech utterances and compared. The z-scores of the two samples showed no significant differences when tested using a one-way ANOVA ($p=0.05$), although nasals in spontaneous speech were observed to be slightly longer. A one-way ANOVA was also used to compare the z-scores of the average as well as the maximum speed of velum movement extracted for each VNV sequence in the read and spontaneous speech samples. Again, no significant differences ($p=0.05$) were found between read and spontaneous samples in either case, which suggests that the extent and speed of velum opening have a minimal effect on spectral reduction, at least for our speaker. However, the variance of the speed measure was much higher in the case of spontaneous nasals, as would be expected.

4.2. Alveolar nasals

Another interesting finding was that the normalized short-term acoustic energy of the nasal /n/ in read speech was much higher than in spontaneous speech ($p<=0.05$). However, no significant differences were found between the normalized short-term energies of the other nasals (/n/,/ng/) for the two speaking styles.

In addition, the palate constriction location coordinates were found for each alveolar nasal when the constriction degree was minimum. The variance of these values was found to be much higher for spontaneous alveolar nasals as compared to the read nasals. This observation agrees with studies suggesting that spectral reduction is characterized by a blurring in constriction location position. In addition, the absolute position of this location

was farther down the vocal tract from the lips in spontaneous speech as compared to read speech, confirming our earlier observation in Section 4.1 with the cavity areas. Rates of change of *absolute* areas (between LA & TTCD, TTCD & TDCD, and TDCD & TRCD) were not found to be significantly different for read and spontaneous nasals.

5. ACKNOWLEDGEMENTS

Work described in this paper was supported by NIH Grant DC007124, the USC Imaging Sciences Center, and the USC Center for High Performance Computing and Communications (HPCC).

6. REFERENCES

- [1] Van Son, R. J. J. H. and Pols, L. C. W. (1999): An acoustic description of consonant reduction. *Speech Communication*. 28: 125-140.
- [2] Nakamura, M., Iwano, K., and Furui, S. (2008). Differences between acoustic characteristics of spontaneous and read speech and their effects on speech recognition performance. *Computer Speech and Language*. 22: 171– 184.
- [3] Byrd, D., Tobin, S., Bresch, E. and Narayanan, S. (2009). Timing effects of syllable structure and stress on nasals: A real-time MRI examination. *Journal of Phonetics*. 37(1): 97-110.
- [4] Moll, K. & Daniloff, R. (1971). Investigation of the timing of velar movement during speech. *Journal of the Acoustical Society of America*. 50: 678-684.
- [5] Narayanan, S., Nayak, K., Lee, S., Sethy, A., and Byrd, D. (2004). An approach to real-time magnetic resonance imaging for speech production. *Journal of the Acoustical Society of America*. 115(4): 1771-1776.
- [6] Bresch, E., Nielsen, J., Nayak, K., and Narayanan, S.. (2006). Synchronized and noise-robust audio recordings during realtime MRI scans. *Journal of the Acoustical Society of America*. 120(4): 1791-1794.
- [7] Bryan Pellom, "SONIC: The University of Colorado Continuous Speech Recognizer", University of Colorado, #TRCSLR- 2001-01, Boulder, Colorado, March, 2001
- [8] Bresch, E., and Narayanan, S. (2009). Region segmentation in the frequency domain applied to upper airway real-time magnetic resonance images. *IEEE Transactions on Medical Imaging*. 28(3): 323-338.
- [9] Ramanarayanan, V., Bresch, E., Byrd, D., Goldstein, L., Narayanan, S. (2009). Analysis of pausing behavior in spontaneous speech using real-time magnetic resonance imaging of articulation. *Journal of the Acoustical Society of America*. 126 (5): EL160-EL165.