

Articulatory tradeoffs reduce acoustic variability during American English /r/ production

Frank H. Guenther^{a)}

*Department of Cognitive and Neural Systems, Boston University, Boston, Massachusetts 02215
and Research Laboratory of Electronics, Massachusetts Institute of Technology, Cambridge,
Massachusetts 02139*

Carol Y. Espy-Wilson

*Department of Electrical and Computer Engineering, Boston University, Boston, Massachusetts 02215
and Research Laboratory of Electronics, Massachusetts Institute of Technology, Cambridge,
Massachusetts 02139*

Suzanne E. Boyce

*Department of Communication Sciences and Disorders, University of Cincinnati, Cincinnati, Ohio
and Research Laboratory of Electronics, Massachusetts Institute of Technology, Cambridge,
Massachusetts 02139*

Melanie L. Matthies

*Department of Communication Disorders, Boston University, Boston, Massachusetts 02215
and Research Laboratory of Electronics, Massachusetts Institute of Technology, Cambridge,
Massachusetts 02139*

Majid Zandipour

*Research Laboratory of Electronics, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139
and Department of Cognitive and Neural Systems, Boston University, Boston, Massachusetts 02215*

Joseph S. Perkell

*Research Laboratory of Electronics, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139
and Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge,
Massachusetts 02139*

(Received 30 December 1997; revised 8 September 1998; accepted 25 January 1999)

The American English phoneme /r/ has long been associated with large amounts of articulatory variability during production. This paper investigates the hypothesis that the articulatory variations used by a speaker to produce /r/ in different contexts exhibit systematic tradeoffs, or articulatory trading relations, that act to maintain a relatively stable acoustic signal despite the large variations in vocal tract shape. Acoustic and articulatory recordings were collected from seven speakers producing /r/ in five phonetic contexts. For every speaker, the different articulator configurations used to produce /r/ in the different phonetic contexts showed systematic tradeoffs, as evidenced by significant correlations between the positions of transducers mounted on the tongue. Analysis of acoustic and articulatory variabilities revealed that these tradeoffs act to reduce acoustic variability, thus allowing relatively large contextual variations in vocal tract shape for /r/ without seriously degrading the primary acoustic cue. Furthermore, some subjects appeared to use completely different articulatory gestures to produce /r/ in different phonetic contexts. When viewed in light of current models of speech movement control, these results appear to favor models that utilize an acoustic or auditory target for each phoneme over models that utilize a vocal tract shape target for each phoneme. © 1999 Acoustical Society of America. [S0001-4966(99)00205-2]

PACS numbers: 43.70.Aj, 43.70.Bk [AL]

INTRODUCTION

The American English phoneme /r/ has long been associated with relatively large amounts of articulatory variability (Alwan *et al.*, 1997; Delattre and Freeman, 1968; Espy-Wilson and Boyce, 1994; Hagiwara, 1994, 1995; Ong and Stone, 1998; Westbury *et al.*, 1995, 1998). In fact, the endpoints of the articulatory continuum for /r/ can be analyzed

as functionally different articulator configurations that use different primary articulators (tongue tip versus tongue dorsum). These endpoints have been characterized in the literature as “bunched” (using the tongue dorsum) and “retroflexed” (using the tongue blade/tip). Often, the same speaker will use different types of /r/ in different productions, e.g., in different phonetic contexts. At the same time, the primary acoustic cue for /r/ is relatively simple and stable: a deep dip in the trajectory of the third spectral energy peak of the acoustic waveform, or third formant frequency (F_3) (Boyce and Espy-Wilson, 1997; Delattre and Freeman, 1968; Westbury *et al.*, 1995, 1998). Furthermore, no consistent acoustic

^{a)}Address correspondence to: Prof. Frank H. Guenther, Boston University, Center for Adaptive Systems and Department of Cognitive and Neural Systems, 677 Beacon Street, Boston, MA 02215, Fax number: (617) 353-7755, Electronic mail: guenther@cns.bu.edu

difference between bunched and retroflexed /r/'s has been discovered.

How is it that a speaker can produce perceptually acceptable /r/'s despite using such variable vocal tract shapes? One possible answer to this question is that the variations in vocal tract shape for /r/ are not haphazard, but are instead systematically related in a way that maintains a relatively stable acoustic signal across productions despite large variations in vocal tract shape across productions. In other words, the different vocal tract shapes used to produce /r/ by a particular subject might involve articulatory tradeoffs, or *trading relations*. The concept of articulatory trading relations is illustrated by the following example. Assume that narrowing either of two constrictions at different locations along the vocal tract (call them location 1 and location 2) has the same effect on an important acoustic cue for a phoneme. Assume further that narrowing either constriction causes a reduction in *F3*. In such a case, one could use different combinations of the two constrictions to achieve the same acoustic effect. For example, to achieve a particular value of *F3*, one might form a very narrow constriction at location 1 and a less narrow constriction at location 2, or one might alternatively form a very narrow constriction at location 2 and a less narrow constriction at location 1. If a speaker used one of these options in one phonetic context and the other option in a second phonetic context, a negative covariance between the sizes of these two constrictions would be seen across phonetic contexts.

The primary purpose of the current study is to investigate the issue of whether the various vocal tract shapes used by an individual to produce /r/ in different phonetic contexts exhibit articulatory trading relations that act to maintain a relatively stable acoustic signal. As discussed at the end of this article, this issue has important implications for theories of speech motor control and speech production. Largely for this reason, several recent experiments have investigated the trading relations issue for phonemes other than /r/ (e.g., de Jong, 1997; Perkell *et al.*, 1993, 1994; Savariaux *et al.*, 1995a), but the results have not been uniform across subjects: Although most subjects exhibit expected articulatory trading relations, some others do not. A possible reason for this ambiguity is that these studies have primarily concentrated on one hypothesized trading relationship, and subjects who do not exhibit this trading relation may exhibit other, unanalyzed trading relations that act to reduce acoustic variability. For example, Perkell *et al.* (1993) investigated an hypothesized trading relation between lip rounding and tongue body raising for the vowel /u/. Three of four subjects showed weak trading relations, but the fourth subject showed the opposite pattern. This fourth subject may have been using other trading relations that overrode the effect of the lip rounding/tongue body raising relationship. In the current study, we employ analysis procedures that allow us to assess the combined effects of multiple articulatory covariances on the variability of the acoustic signal. Furthermore, American English /r/ was chosen¹ because the large amount of articulatory variability associated with /r/ productions should make it easier to detect trading relations if they are present.

I. METHODS

A. Data collection

An electromagnetic midsagittal articulometer (EMMA) system (Perkell *et al.*, 1992) was used to track the movements of six small (5 mm long×2.5 mm diameter) transducer coils. The coils were attached in the midsagittal plane to the tongue (three coils), lips (two coils), and lower incisor (one coil) with bio-compatible adhesive. Transducers were also placed on the upper incisor and the bridge of the nose, for defining a coordinate system with a maxillary frame of reference. A directional microphone was suspended 14 in. from the subject's mouth and the acoustic signal was recorded simultaneously with the EMMA signals. Standard EMMA calibration protocols were completed prior to each experiment (cf. Perkell *et al.*, 1992 for details). The current study focused on the positions of the three tongue transducers, which were located approximately 1, 2.5, and 5 cm back from the tongue tip (with the tongue in a neutral configuration).

The seven subjects were young adults, two females (subjects 2 and 3) and five males. They had no history of speech, language, or hearing deficits or pronounced regional dialects. Each of the seven subjects produced 4–6 repetitions of the carrier phrase "Say ____ for me" for each of the five test utterances; /warav/, /wabrav/, /wadrav/, /wagrav/, and /wavrav/. The articulatory and acoustic data from these utterances were time aligned to allow direct comparison between the two data types.

B. *F3* extraction and alignment

The minimum measured *F3* value during /r/ production, which can be thought of as the acoustic "center" of /r/, served as a landmark for time alignment of the data across utterances for each speaker. Formant tracks were computed for all utterances using the ESPS/WAVES formant tracker and a 51.2-ms window and 3.2-ms frame rate. The *F3* minimum was detected using an automatic procedure that first identified all sonorant regions, then located the point of minimal *F3* from the relevant sonorant regions. *F3* values and transducer positions within a 140-ms time window centered at the *F3* minimum were extracted. Extracted *F3* traces for some utterances were corrupted due to technical difficulties in automatically tracking low-amplitude and low-frequency values of *F3* after stop consonants. Therefore, utterances whose *F3* tracks changed by more than 200 Hz in a 3.2-ms time step were eliminated from the study, leaving 12 to 27 analyzed utterances per subject. After this elimination process, the tongue shapes at the *F3* minimum of the remaining utterances were visually inspected, and two additional utterances (one each for subjects 1 and 4) were identified as having articulations that were incorrectly labeled as /r/ by the automatic extraction process. These two utterances were also eliminated from the study.

C. Effects of vocal tract shape parameters on *F3*

The vocal tract shape for /r/ involves a palatal constriction formed by the tongue in the anterior third of the tract.

Roughly speaking, the third formant frequency ($F3$) of /r/ is the resonance resulting from the cavities anterior to the palatal constriction (e.g., Alwan *et al.*, 1997; Espy-Wilson *et al.*, 1997; Stevens, 1998). This part of the vocal tract consists of an acoustic compliance due to a large front cavity volume and two parallel acoustic masses due to natural tapering by the teeth/lips and the palatal constriction behind the front cavity. The resulting resonance is inversely proportional to the product of the total acoustic mass and the acoustic compliance. Because it is difficult to accurately infer lip aperture from EMMA data, we focus on the effects of the acoustic mass due to the size and location of the palatal constriction. From these considerations, we conclude that the frequency of $F3$ can be decreased by tongue movements that lengthen the front cavity (thereby increasing the acoustic compliance of the front cavity), lengthen the constriction (thereby increasing the acoustic mass of the constriction behind the front cavity), or decrease the area of the constriction (thereby increasing the acoustic mass of the constriction).²

The predicted effects of these movements on $F3$ were confirmed using vocal tract area functions derived from structural MRI scans of a speaker producing /r/.³ Two area functions were derived: one representing a “bunched” /r/ configuration, and one representing a “retroflexed” /r/ configuration. Three manipulations were carried out on each area function to test the effects on $F3$ predicted from acoustic theory: (i) the palatal constriction was extended backward by narrowing the vocal tract area immediately behind the constriction; (ii) the front cavity was lengthened by displacing the palatal constriction backward; and (iii) the vocal tract area at the palatal constriction was decreased. For all three manipulations, an acoustic signal was synthesized (using S. Maeda’s VTCALCS program; Maeda, 1990) and compared to the signal synthesized from the original area function. Each manipulation resulted in a lower $F3$ in both the bunched and retroflexed /r/ cases, as expected from the acoustic theory analysis.

Because all three manipulations act to lower $F3$, subjects could maintain a relatively stable $F3$ despite vocal tract shape variations across contexts if these variations involved tradeoffs between the different manipulations. When looking at the different vocal tract shapes for /r/ across contexts, these tradeoffs would be manifested by correlations between constriction length, front cavity length, and constriction area. Specifically, the following three correlations would be expected to aid in maintaining a relatively stable $F3$ across utterances while allowing variations in vocal tract shape:

- (1) a negative correlation between constriction length and front cavity length, since increases in constriction length and front cavity length both act to reduce $F3$;
- (2) a positive correlation between constriction length and constriction area, since increases in constriction length reduce $F3$ and decreases in constriction area reduce $F3$; and
- (3) a positive correlation between front cavity length and constriction area, since increases in front cavity length reduce $F3$ and decreases in constriction area reduce $F3$.

D. Predicted articulatory covariances

To determine whether a subject uses any of the three trading relations hypothesized above, we must first describe the trading relations in terms of the x and y coordinates of the tongue transducers. For tongue configurations during /r/ production, a forward movement of the tongue front transducer generally corresponds to a shortening of the front cavity, an upward movement of the tongue front transducer generally corresponds to a decrease in the area of the palatal constriction for /r/, and, since the point of maximal constriction for /r/ is typically anterior to the tongue back transducer, an upward movement of the tongue back transducer generally corresponds to a lengthening of the palatal constriction and possibly a decrease in the area of the constriction. When determining the signs of the transducer coordinate correlations corresponding to the trading relations delineated above, we must take into account that increasing values of the tongue front horizontal position correspond to decreases in front cavity length, and increasing values of the tongue front vertical position correspond to decreases in constriction area. From these considerations, we can surmise that the three trading relation strategies described above should be evidenced by the following articulatory correlations:

- (1) a positive correlation between tongue back height and tongue front horizontal position;
- (2) a negative correlation between tongue back height and tongue front height; and
- (3) a positive correlation between tongue front horizontal position and tongue front height.

Note that the use of all three trading relations by a single subject is unlikely given that they impose competing constraints; i.e., if tongue back height and tongue front horizontal position are positively correlated as in relation (1), and tongue front horizontal position and tongue front height are positively correlated as in relation (3), it is very likely that tongue back height and tongue front height will also be positively correlated, thus violating relation (2).

E. Analysis of articulatory and acoustic variances

To quantify the combined effects of articulatory covariances on $F3$ variability, an analysis was performed using both acoustic and articulatory data to estimate $F3$ variance as a function of articulatory variances. The relationship between transducer coordinates and $F3$ during /r/ can be written for each speaker as follows:

$$F3 = A_0 + \sum_{i=1}^N A_i c_i + E, \quad (1)$$

where the A_i are constants, the c_i are the transducer coordinates, N is the number of transducer coordinates considered in the analysis, and E is a residual term that accounts for the effects on $F3$ due to all other sources, including articulators not included in the analysis, measurement errors, and nonlinearities in the relationship between $F3$ and the transducer coordinates. The equation relating $F3$ variance to articulatory variances at each point in time is then:

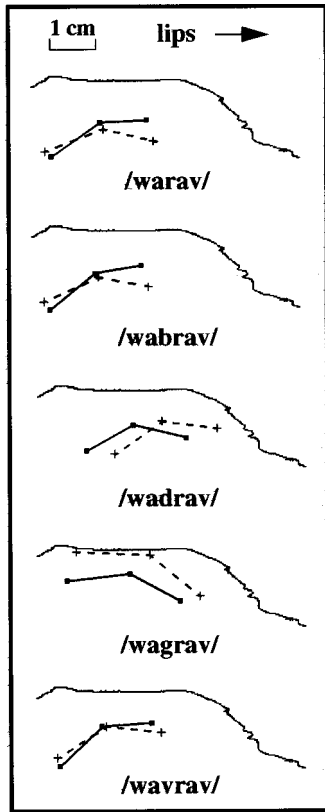


FIG. 1. Sample lingual articulations used by subject 1 to produce /r/ in the five phonetic contexts. For each context, two schematized tongues shapes and a palatal trace are shown. Each tongue shape schematic was formed by connecting the three tongue transducers with straight lines. The tongue shape at the F_3 minimum for /r/ is drawn with solid lines. The tongue shape 70 ms prior to the F_3 minimum is drawn with dashed lines.

$$\begin{aligned} \text{Var}(F_3) = & \sum_i A_i^2 \text{Var}(c_i) + \text{Var}(E) \\ & + 2 \sum_{i < j} \sum A_i A_j \text{Cov}(c_i, c_j) \\ & + 2 \sum_i A_i \text{Cov}(c_i, E). \end{aligned} \quad (2)$$

To determine the effects of articulatory covariances on F_3 variability, we can compare the variance estimate of Eq. (2) to the following variance estimate that excludes the covariances between the analyzed transducer coordinates:

$$\begin{aligned} \text{Var}(F_3) = & \sum_i A_i^2 \text{Var}(c_i) + \text{Var}(E) \\ & + 2 \sum_i A_i \text{Cov}(c_i, E). \end{aligned} \quad (3)$$

If the F_3 variance estimate in the absence of articulatory covariances [Eq. (3)] is significantly larger than the variance estimate including the articulatory covariances [Eq. (2)], we conclude that the primary effect of the articulatory covariances is a reduction in the variance of F_3 .

Strictly speaking, a comparison of the F_3 variance estimates in Eqs. (2) and (3) tells us only about the effects of the covariances of the *linear* component of each transducer's

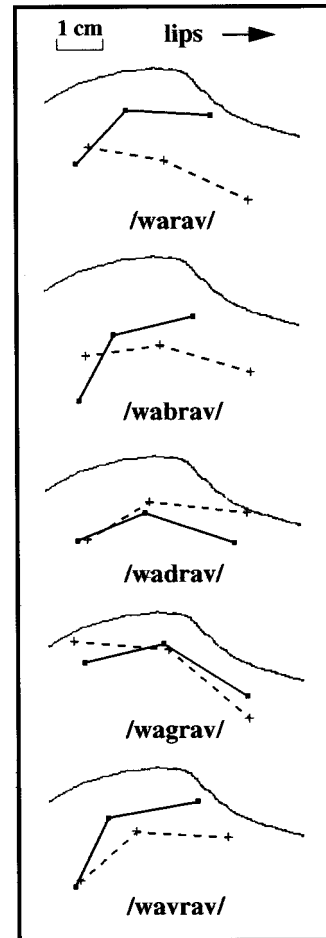


FIG. 2. Sample lingual articulations used by subject 2 to produce /r/ in the five phonetic contexts. For each context, two schematized tongues shapes and a palatal trace are shown. Each tongue shape schematic was formed by connecting the three tongue transducers with straight lines. The tongue shape at the F_3 minimum for /r/ is drawn with solid lines. The tongue shape 70 ms prior to the F_3 minimum is drawn with dashed lines.

relation to F_3 . However, the relationship between F_3 and transducer coordinates should be linear near a particular configuration of the vocal tract, since F_3 is presumably a continuous nonlinear function of the vocal tract area function, and such functions are locally linear. One would further expect that the relationship is still approximately linear for the relatively limited range of vocal tract configurations utilized by a particular subject for /r/. The linear approximations reported below captured approximately 80% of the variance when using only three pellet coordinates, providing support for the assertion that the primary effect of articulatory covariances on F_3 variance can be captured by considering only the linear component of each transducer's relationship to F_3 . Furthermore, the sign (positive or negative) of an articulatory covariance's contribution to F_3 variance depends only on the sign of the corresponding A_i terms, and we are primarily interested in the sign of the combined effects of articulatory covariances on F_3 variance. The expected signs of the A_i for tongue back height, tongue front horizontal position, and tongue front height can be deduced from acoustic theory considerations (Secs. IC and ID). A_i values were estimated for each subject using multiple linear regression on the acoustic and articulatory data. As discussed in Sec. IID, all

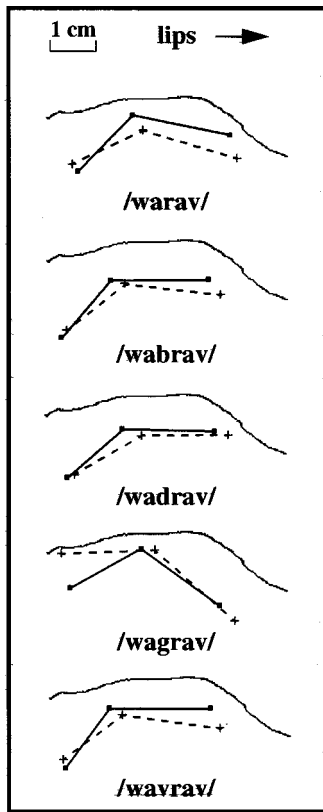


FIG. 3. Sample lingual articulations used by subject 3 to produce /r/ in the five phonetic contexts. For each context, two schematized tongue shapes and a palatal trace are shown. Each tongue shape schematic was formed by connecting the three tongue transducers with straight lines. The tongue shape at the F_3 minimum for /r/ is drawn with solid lines. The tongue shape 70 ms prior to the F_3 minimum is drawn with dashed lines.

21 estimated A_i values (three values for each of seven subjects) were of the sign expected from these acoustic theory considerations.

II. RESULTS

A. Temporal progression of tongue shapes

Figures 1–7 show sample lingual articulations used to produce /r/ in the five contexts by the seven subjects. For each context, two schematized tongue shapes and a palatal trace⁴ are shown. The tongue shape schematics were formed by connecting the three tongue transducers with straight lines. The solid tongue shape corresponds to the point in time at which F_3 reached its minimum value. The tongue shape 70 ms prior to the F_3 minimum is indicated by dashed lines. The movement of the tongue can thus be roughly characterized as a transition from the dashed tongue shape to the solid tongue shape. This movement corresponds to the articulation toward the “acoustic center” of /r/; i.e., the portion of the movement up to the point in time of the F_3 minimum.

Inspection of the lingual articulations for some subjects suggests that these subjects utilize different articulatory gestures, aimed at different vocal tract shapes, to produce /r/ in different phonetic contexts. For example, the backward movement of the tongue, with a slight downward movement of the tongue blade, used by subject 1 to produce the /r/ in /wadrav/ does not appear to be aimed at the same vocal tract

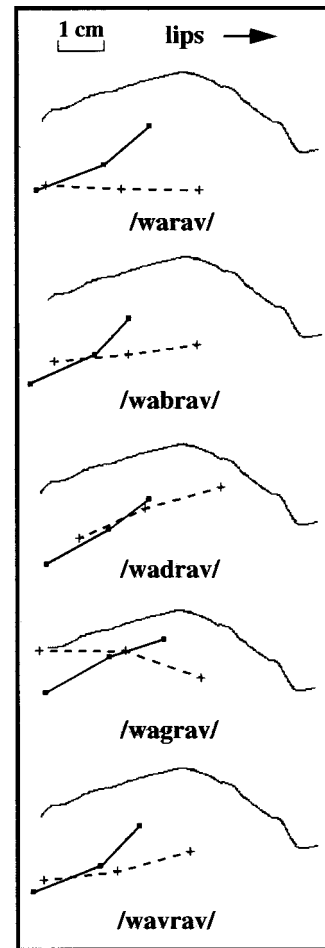


FIG. 4. Sample lingual articulations used by subject 4 to produce /r/ in the five phonetic contexts. For each context, two schematized tongue shapes and a palatal trace are shown. Each tongue shape schematic was formed by connecting the three tongue transducers with straight lines. The tongue shape at the F_3 minimum for /r/ is drawn with solid lines. The tongue shape 70 ms prior to the F_3 minimum is drawn with dashed lines.

shape for /r/ as the upward movements of the tongue blade used by the same subject to produce /r/ in the /warav/, /wabrav/, and /wavrav/ contexts (Fig. 1). Similarly, the downward movement of the tongue blade used by subject 2 to produce the /r/ in /wadrav/ does not appear to be aimed at the same vocal tract shape as the upward movements of the tongue blade used by the same subject to produce /r/ in /warav/, /wabrav/, or /wavrav/ (see Fig. 2). Additional examples of this phenomenon can be seen in Figs. 1–7. The possible relevance of these observations to theories of speech motor control will be addressed in Sec. III.

B. Tongue shapes at acoustic center of /r/

Figure 8 shows tongue configurations at the F_3 minimum of /r/ for each of the seven speakers. For each utterance, the three tongue transducer positions are connected by a straight line. The tongue configurations for all repetitions in all phonetic contexts are superimposed for each speaker. Thus the fact that different numbers of utterances were analyzed for different subjects and contexts is reflected in this figure. As previously reported elsewhere (e.g., Delattre and Freeman, 1968; Hagiwara, 1994, 1995; Ong and Stone,

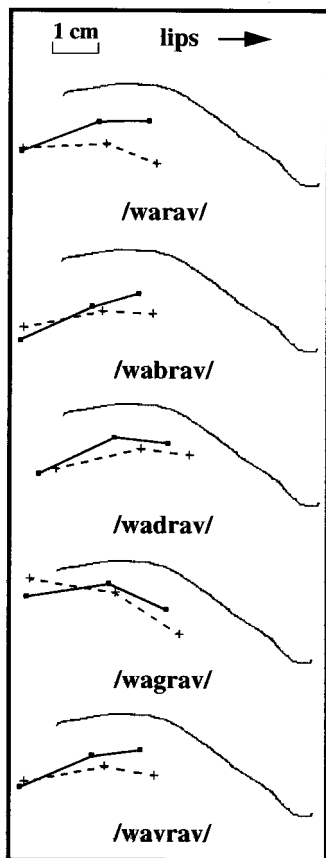


FIG. 5. Sample lingual articulations used by subject 5 to produce /r/ in the five phonetic contexts. For each context, two schematized tongue shapes and a palatal trace are shown. Each tongue shape schematic was formed by connecting the three tongue transducers with straight lines. The tongue shape at the F_3 minimum for /r/ is drawn with solid lines. The tongue shape 70 ms prior to the F_3 minimum is drawn with dashed lines.

1998; Westbury *et al.*, 1995), a wide range of tongue shapes is seen both within and across subjects. Also of note is the fact that, although most subjects seem to use an approximate continuum of tongue shapes (e.g., S2, S3, S6, S7), others show a more bimodal distribution of tongue shapes (e.g., S4, S5). Finally, the tongue shapes across subjects appear to form an approximate continuum between a bunched configuration (e.g., S6) and a retroflexed configuration (e.g., S4). A more detailed indication of the effects of the different phonetic contexts on the tongue shapes for /r/ can be gained from Fig. 9, which shows the average tongue shapes used by each subject in each phonetic context, coded by phonetic context. Figures 10–16 show the corresponding average F_3 traces, starting from the point of the F_3 minimum for /r/ and continuing for 70 ms, for each speaker coded by phonetic context. With the exception of the /wadrav/ productions of subject 2, which had considerably higher values of F_3 than the other utterances for that subject, the subjects showed minimum F_3 values well below 2000 Hz in all contexts, as expected from earlier studies of /r/ production.

Figure 17 shows the midsagittal palatal outline (thick solid line) and mean tongue shapes at the time of the F_3 minimum for /r/ for each of the seven subjects. For each subject, mean configurations from two phonetic contexts (solid and dashed lines) are shown to illustrate the range of

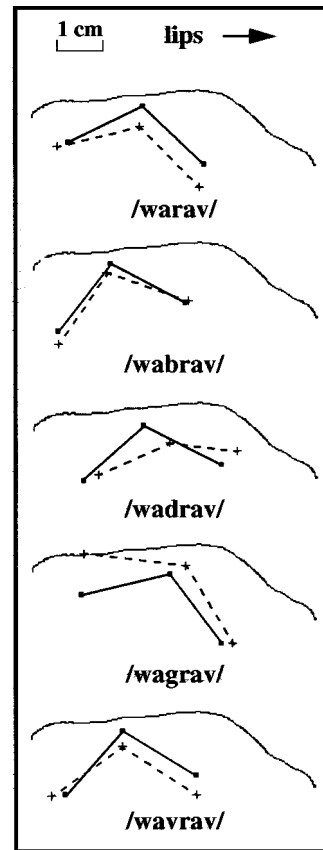


FIG. 6. Sample lingual articulations used by subject 6 to produce /r/ in the five phonetic contexts. For each context, two schematized tongue shapes and a palatal trace are shown. Each tongue shape schematic was formed by connecting the three tongue transducers with straight lines. The tongue shape at the F_3 minimum for /r/ is drawn with solid lines. The tongue shape 70 ms prior to the F_3 minimum is drawn with dashed lines.

tongue shapes used by that subject. Tongue outlines were created by connecting the average positions of the three tongue transducers for a given utterance with a smooth curve to roughly approximate tongue shape.⁵ A line was then extended downward from the tongue front transducer position, then forward to the lower incisor transducer position, to provide a rough estimate of the relative size of the front cavity across contexts.⁶ Also shown in the upper left corner of this figure are two superimposed, highly schematic vocal tract outlines that illustrate trading relations for maintaining a relatively stable F_3 . The effect on F_3 of the longer front cavity of the dashed outline, which can be roughly characterized as a retroflexed /r/, is counteracted by the effects of the longer and slightly narrower constriction of the solid outline, which can be roughly characterized as a bunched /r/. Similarly, the vocal tract outlines for all subjects indicate that shorter front cavity lengths are accompanied by a compensating increase in constriction length and/or decrease in the constriction area. Furthermore, the tongue shapes during /wagrav/ (solid lines) are generally much closer in shape to tongue shapes for /g/ than are the /r/ shapes for /wabrav/ or /warav/ (dashed lines), suggesting that subjects utilize /r/ configurations that are reached relatively easily in the current phonetic context.

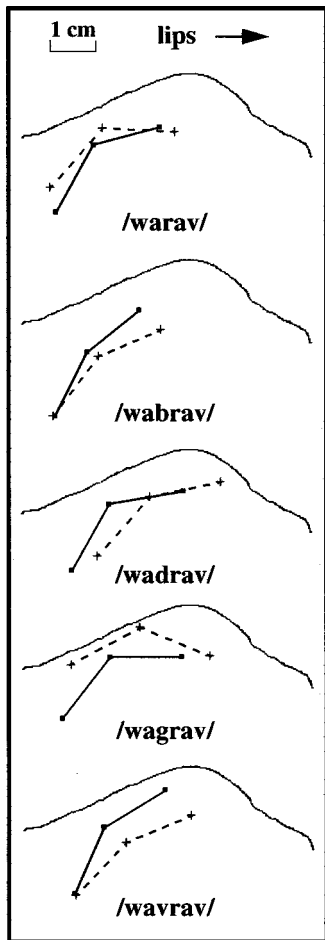


FIG. 7. Sample lingual articulations used by subject 7 to produce /r/ in the five phonetic contexts. For each context, two schematized tongue shapes and a palatal trace are shown. Each tongue shape schematic was formed by connecting the three tongue transducer positions with straight lines. The tongue shape at the F_3 minimum for /r/ is drawn with solid lines. The tongue shape 70 ms prior to the F_3 minimum is drawn with dashed lines.

C. Articulatory trading relations

For each subject, Pearson correlation coefficients corresponding to the predicted covariances described in Sec. ID were estimated across utterances at the point of F_3 minimum and are listed in Table I. All subjects showed a significant positive correlation between tongue back height (TBY in Table I) and tongue front horizontal position (TFX), indicative of a trading relation between constriction length and front cavity length. Six of seven subjects also showed a second strong trading relation: five subjects showed a trading relation between constriction length and constriction area as evidenced by a negative correlation between TBY and tongue front height (TFY), and one subject showed a trading relation between front cavity length and constriction area as evidenced by a positive correlation between TFX and TFY. One subject (subject 7) showed only very weak correlations other than the strong trading relation between tongue back height and tongue front horizontal position.

D. Analysis of acoustic and articulatory variabilities

The results in Sec. IIC indicate that most subjects exhibited two of three hypothesized articulatory trading rela-

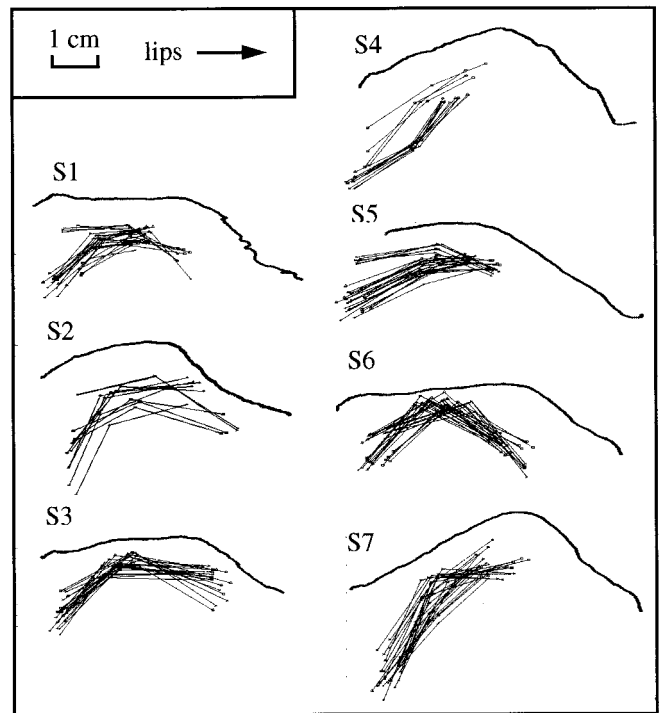


FIG. 8. Tongue configurations at the F_3 minimum of /r/ for each of the seven speakers. For each utterance, the three tongue transducer positions are connected by straight lines. The tongue configurations for all repetitions in all phonetic contexts are superimposed for each speaker.

tionships that act to reduce acoustic variability. Furthermore, as described in Sec. ID, it is unlikely or impossible for a subject to utilize all three trading relations because they counteract one another. However, it is still possible that the

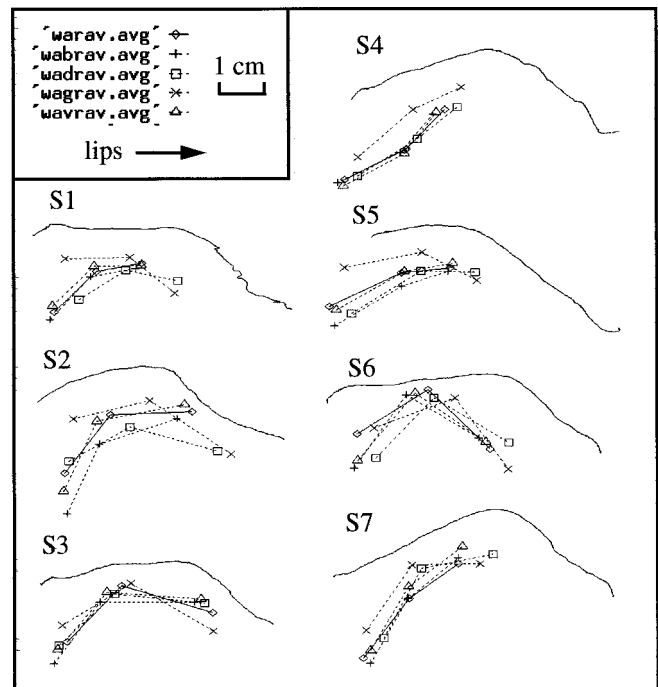


FIG. 9. Averaged tongue configurations at the F_3 minimum of /r/ for each of the seven speakers. The averaged positions of the three tongue transducer positions for each of the five phonetic contexts are connected by straight lines.

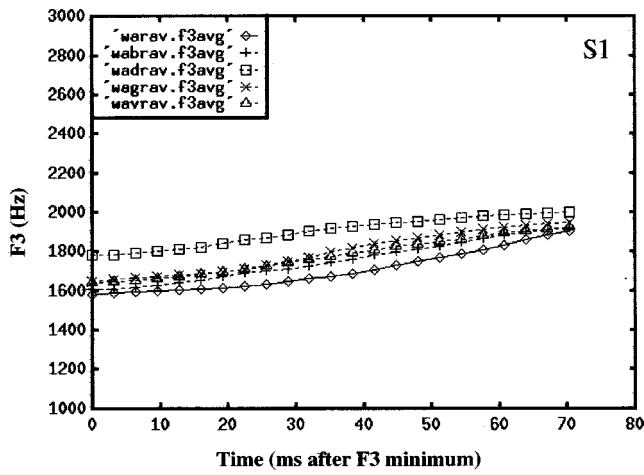


FIG. 10. Averaged F_3 values for /r/ in the five phonetic contexts for subject 1. For each context, the averaged F_3 is plotted as a function of time starting from the F_3 minimum.

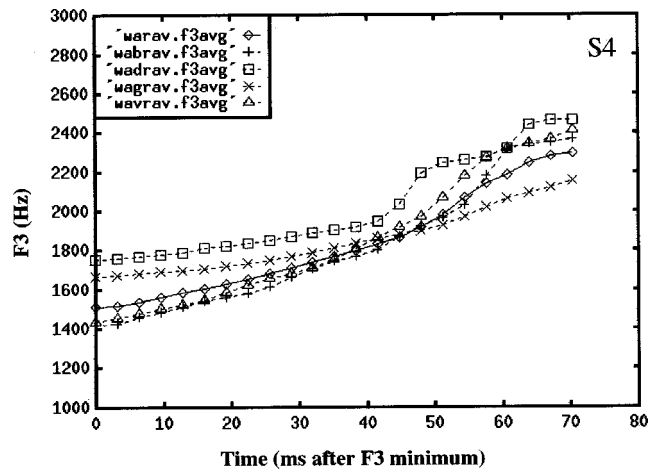


FIG. 13. Averaged F_3 values for /r/ in the five phonetic contexts for subject 4. For each context, the averaged F_3 is plotted as a function of time starting from the F_3 minimum.

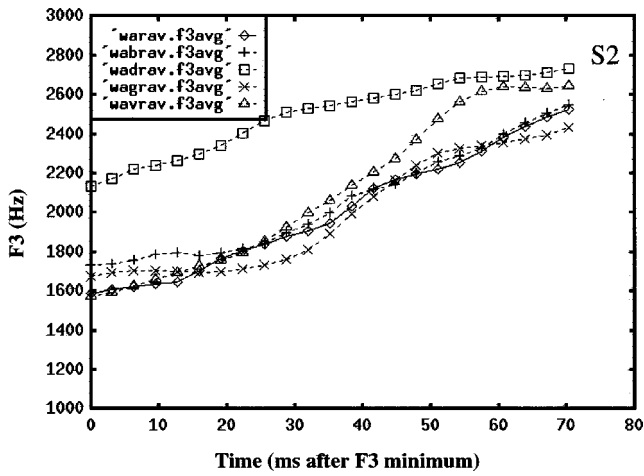


FIG. 11. Averaged F_3 values for /r/ in the five phonetic contexts for subject 2. For each context, the averaged F_3 is plotted as a function of time starting from the F_3 minimum.

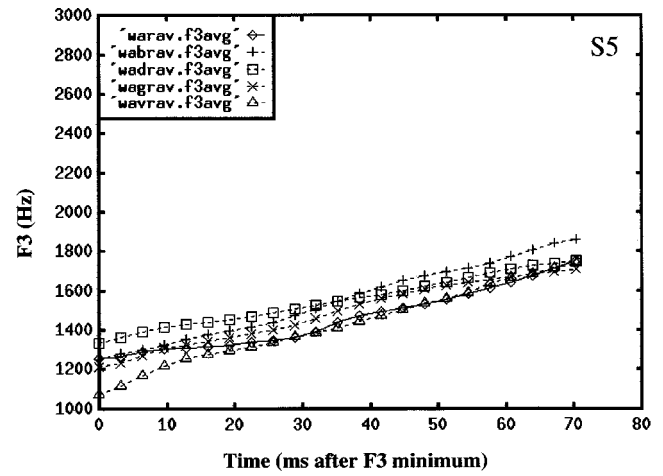


FIG. 14. Averaged F_3 values for /r/ in the five phonetic contexts for subject 5. For each context, the averaged F_3 is plotted as a function of time starting from the F_3 minimum.

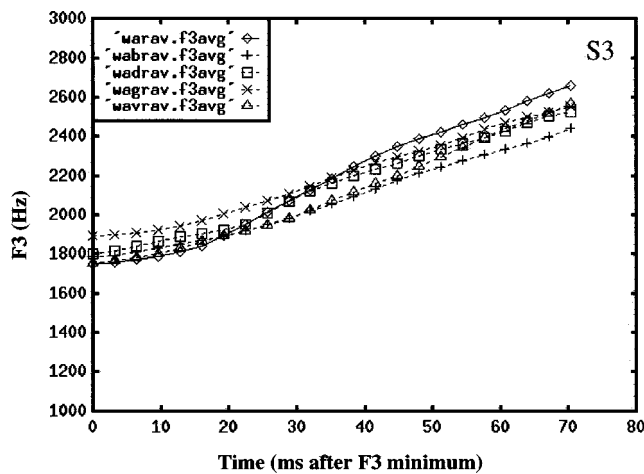


FIG. 12. Averaged F_3 values for /r/ in the five phonetic contexts for subject 3. For each context, the averaged F_3 is plotted as a function of time starting from the F_3 minimum.

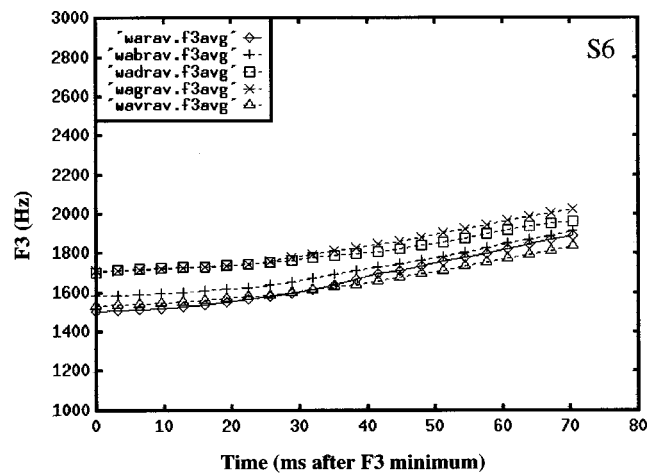


FIG. 15. Averaged F_3 values for /r/ in the five phonetic contexts for subject 6. For each context, the averaged F_3 is plotted as a function of time starting from the F_3 minimum.

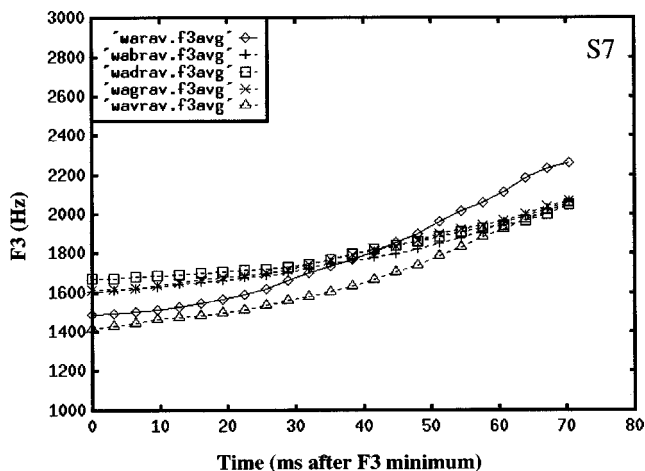


FIG. 16. Averaged F_3 values for /r/ in the five phonetic contexts for subject 7. For each context, the averaged F_3 is plotted as a function of time starting from the F_3 minimum.

significant correlations that violate the trading relations could effectively “override” the beneficial articulatory tradeoffs, potentially nullifying or even reversing the effect of the utilized trading relations on acoustic variability. It is therefore necessary to estimate the net effect of all three articulatory covariances, as outlined in Sec. I E.

F_3 variance estimates with and without covariance

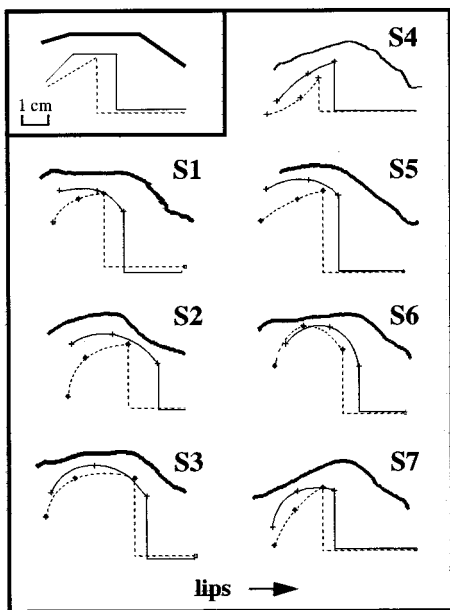


FIG. 17. Trading relations during /r/ production. The upper left corner shows two superimposed, highly schematized vocal tract outlines (dashed and solid lines) illustrating trading relations between front cavity length and palatal constriction length and area. Also shown are vocal tract outlines that illustrate the range of tongue shapes used by each of the seven subjects to produce /r/ in different phonetic contexts. Thin solid lines correspond to the tongue shapes for /r/ in /wagrav/ (averaged across repetitions), and dashed lines correspond to the /r/ in /wabrav/ or /warav/, depending on the subject. Thick solid lines indicate palatal outlines. Each outline is formed by connecting the three tongue transducer positions with a smooth curve, then projecting downward and forward from the frontmost tongue transducer to the lower incisor transducer. All seven subjects show tradeoffs between the front cavity length and the constriction length and/or area when producing /r/ in the two different contexts.

TABLE I. Articulator correlation coefficients. Significant correlations that are consistent with hypothesized trading relations are shown in boldface. TBY=tongue back height; TFX=tongue front horizontal position; TTY=tongue front height.

Subject	TBY-TFX	TBY-TFY	TFX-TFY
1	0.74^a	-0.75^a	-0.75 ^a
2	0.92^a	-0.69^a	-0.88 ^a
3	0.77^a	-0.74^a	-0.46
4	0.91^a	0.89 ^a	0.82^a
5	0.64^a	-0.49^a	-0.57 ^a
6	0.55^a	-0.81^a	-0.60 ^a
7	0.84^a	0.05	0.06

^aStatistically significant ($p < 0.01$).

terms [Eqs. (2) and (3), respectively] were calculated using the tongue back height, tongue front horizontal position, and tongue front height transducer coordinates. The corresponding F_3 standard deviations were then averaged across subjects. The A_i values for each speaker were estimated using multiple linear regression across utterances and time bins and are provided in Table II; the value of E for a particular time bin was simply the residual of the regression in that time bin. R^2 values for the F_3 fit (without the residual term) ranged from 0.75 to 0.87 for the different subjects, with an average R^2 of 0.79. If covariances are high and the actual effect of an articulator’s position on F_3 is very low, regression analysis can possibly result in estimates of transducer contributions that have the wrong sign, which could in turn cause some articulatory covariances to decrease estimated F_3 variability when in reality they increase or have no significant effect on F_3 variability. The fact that none of the transducer contribution estimates produced by the regression were of the opposite sign as expected from acoustic theory considerations and the MRI-based area function analysis indicates that this potential problem did not affect our results.

F_3 standard deviation estimates with and without covariance terms are shown in Fig. 18 as a function of time starting at the F_3 minimum for /r/, averaged across subjects. (Standard deviations were plotted in place of variances to produce values whose units are Hz.) Also plotted is the standard deviation obtained from measured values of F_3 . When articulatory covariances are included, the F_3 standard deviation estimate is equal to the measured F_3 standard deviation; this is as expected because of the inclusion of the residual term in the variance estimate calculations. The solid line in the figure thus represents both the measured F_3 standard deviation and the estimated F_3 standard deviation including

TABLE II. Regression coefficients indicating the relationship between transducer coordinates and F_3 . Units are Hz/mm.

Subject	A_1 (TBY)	A_2 (TFX)	A_3 (TTY)
1	-28.25	15.20	-35.61
2	-81.13	92.77	-35.25
3	-12.28	25.93	-51.50
4	-36.83	77.34	-36.55
5	-24.04	21.44	-30.49
6	-21.68	10.99	-30.95
7	-46.29	33.71	-31.87

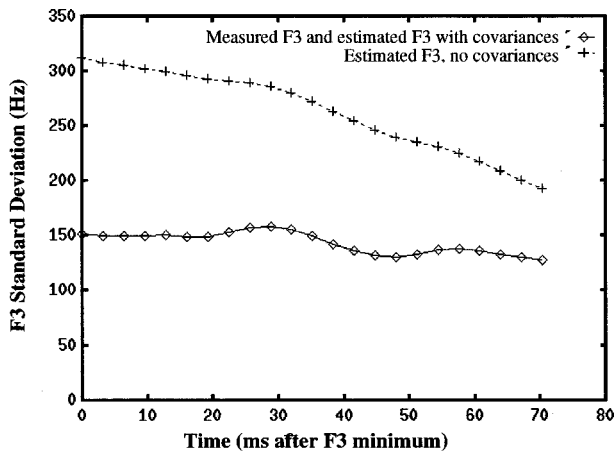


FIG. 18. Comparison of the measured $F3$ standard deviation with estimates of $F3$ standard deviation derived from articulator transducer positions, averaged across subjects. Removal of the articulatory covariance terms results in a much higher estimate of $F3$ standard deviation, indicating that articulatory covariances greatly reduce acoustic variation.

the covariance terms. When articulatory covariances are removed from the estimates, the estimated $F3$ standard deviation increases substantially. The dashed line in Fig. 18 represents estimated $F3$ standard deviation without covariances using the three tongue transducer coordinates. According to this estimate, then, $F3$ standard deviation would be 105% higher at the acoustic center of /r/ if the articulatory tradeoffs had not been present.

The increase in the $F3$ variance estimate without covariances is seen at the $F3$ minimum for all subjects. This observation suggests that ambiguous results from previous studies may have been at least partly due to analyzing only one articulatory tradeoff at a time, since in our study no subject exhibited all three hypothesized trading relations, but all subjects showed a net decrease in acoustic variability due to the *combined* effects of the articulatory covariances. Assume, for example, that the data listed in each column of Table I were the result of an independent research study. Researchers investigating the trading relation in column 2 (in which five of seven subjects used the trading relation) would sharply disagree with researchers investigating the trading relation in column 3 (in which only one of seven subjects used the trading relation) as to whether or not trading relations are reliably exhibited, and both sets of researchers would report ambiguous results since in neither case do all subjects behave in the same way. A much clearer picture has emerged from the current study due to the analysis of the combined effects of the articulatory covariances.

Also evident in Fig. 18 is a steady decrease with time of the effects of the covariance terms on $F3$ as the /r/ transitions into the following /a/. This decrease, evident in six of the seven subjects, is suggestive of a decrease in the use of trading relations as vocal tract shape differences across utterances (due to the different phonemes preceding /r/ in different utterances) diminish.

III. DISCUSSION

The results of this study indicate that the widely varying tongue shapes used by each subject to produce /r/ in different phonetic contexts exhibit systematic articulatory tradeoffs

that act to reduce acoustic variability across contexts despite the large variations in vocal tract shape. The first piece of evidence arises from inspection of the tongue shape extremes for each subject (i.e., the averaged tongue shapes from the two phonetic contexts whose averaged tongue shapes differed the most), which show clear tradeoffs between the length of the front cavity and the length and/or area of the palatal constriction (Fig. 17). An acoustic theory analysis of the effects of these aspects of vocal tract shape on $F3$ indicates that the systematic tradeoffs would be expected to reduce acoustic variability across contexts. Analysis of articulatory covariances indicates that most speakers exhibit two of three articulatory trading relationships that were hypothesized based on acoustic properties of the vocal tract (Sec. II C). Furthermore, the use of all three trading relationships is very unlikely if not impossible given that they counteract each other. Finally, analysis of the combined effects of these articulatory covariances indicates that they strongly influence $F3$ variability across contexts, effectively cutting $F3$ standard deviation in half compared to what it would have been without the articulatory covariances (Sec. II D; Fig. 18).

Unlike earlier trading relations studies that reported mixed results across subjects (e.g., de Jong, 1997; Perkell *et al.*, 1993, 1994; Savariaux *et al.*, 1995a), the reduction of $F3$ variability due to articulatory covariances was seen at the acoustic center of /r/ for all seven subjects in the current study. We believe that the following factors contributed to this difference. First, the current study investigated a phoneme known to exhibit a large amount of articulatory variability across contexts. Such a sound would be expected to exhibit stronger trading relations due to the larger overall articulatory variability. Second, the current study investigated the combined effects of multiple articulatory covariances. Although the combined effect of articulatory covariances was a reduction of $F3$ variability in all seven subjects according to the analysis of Sec. II D, different subjects used different combinations of the individual articulatory trading relations (Table I). It is therefore not surprising that in earlier studies, which investigated articulatory covariances individually, some subjects did not use an hypothesized trading relationship. The results of the current report suggest that these subjects may well have used other, unanalyzed trading relations that reduced acoustic variability.

The issue of articulatory trading relations in speech production is relevant to current theories concerning the control of speech movements. Roughly speaking, computational models of speech motor control can be classified according to the type of phonemic “targets” that they use to command movements of the speech articulators. One type of computational model, exemplified by the task-dynamic model of Saltzman and Munhall (1989), utilizes a target for each phoneme that specifies important aspects of the shape of the vocal tract for that phoneme. This “vocal tract shape target” view is closely related to theories of speech perception and production in which the articulatory gesture serves as the basic unit of speech. These include the motor theory of speech perception (Liberman and Mattingly, 1985; Liberman *et al.*, 1967), the direct realist theory of speech perception (Fowler, 1986, 1996), and the linguistic-gestural theory of

phonology (Browman and Goldstein, 1990a,b). A second type of computational model, exemplified by the DIVA model of speech acquisition and production (Guenther, 1995; Guenther *et al.*, 1998), utilizes only an acoustic or auditory target for each phoneme, with no explicit vocal tract shape target. These models may use different shapes of the vocal tract to produce the same acoustic signal for a phoneme depending on things like phonetic context. Theories related to this “auditory target” view have been posited by various researchers (for some recent examples, see Johnson *et al.*, 1993; Perkell *et al.*, 1993, 1995, 1997; Savariaux *et al.*, 1995a,b) and many of the roots for this line of thinking can be traced to the influential work of Jakobson *et al.* (1951).

A major difference between the auditory target and vocal tract shape target computational model classes is that the former explicitly predict the existence of articulatory trading relations when producing the same phoneme in different contexts, whereas the latter do not. Because the current results show the existence of trading relations in all seven subjects, they appear to favor acoustic target models over vocal tract shape target models. A potential reason for the use of articulatory tradeoffs is that they can reduce the amount of effort required to move the articulators through a set of acoustic targets. For example, the tongue shapes for /r/ in /wagrav/ were generally closer to the tongue shapes for /g/ than the tongue shapes for /r/ in other contexts, suggesting that, to a first approximation, subjects moved to the closest vocal tract shape that could be used to produce the appropriate sound in the prevailing phonetic context. The acoustic target control scheme used by the DIVA model has this property.

It is possible, however, that the trading relations found in the current study are artifactual, rather than the result of a movement control process that is making use of them to minimize effort while maintaining an acceptable acoustic signal. For example, one might imagine that the tongue shape variations simply reflect a controller that undershoots the vocal tract shape target for /r/ (thus maintaining some aspects of the vocal tract shape from the preceding phoneme at the point of *F3* minimum), rather than one that purposefully utilizes a different vocal tract shape for /r/ in different contexts because it “knows” how to use different shapes to achieve the same acoustic/auditory result in different contexts. However, the fact that some subjects appear to use articulatory gestures toward different vocal tract shapes to produce /r/ in different contexts argues against an “artifact” interpretation. For example, the backward movement of the tongue, with a slight downward movement of the tongue blade, used by subject 1 to produce the /r/ in /wadrav/ does not appear to be a truncated movement aimed at the same vocal tract shape for /r/ as the upward movements of the tongue blade used by the same subject to produce /r/ in the /warav/, /wabrav/, and /wavrav/ contexts (see Fig. 1). Instead, it appears that the /r/ gesture in /wadrav/ is toward a different vocal tract shape that nonetheless suffices for /r/. (See Sec. II and Figs. 1–7 for further examples.) Further research into this issue will be carried out in a subsequent modeling study that will address other possible explanations for the trading relations found in the current study.

ACKNOWLEDGMENTS

We thank Ken Stevens, Harlan Lane, Dan Bullock, and John Ingram for their help and comments, Shinji Maeda for the VTCALCS program, Abeer Alwan and colleagues for the MRI data, and John Westbury and Anders Löfqvist for constructive reviews of an earlier draft of this paper. Supported by NIDCD Grants Nos. 1R29-DC02852 to Frank Guenther, 5R01-DC01925-04 to Joseph Perkell, and 1R03-C2576-01 to Suzanne Boyce, and NSF Grant No. IRI-9310518 to Carol Espy-Wilson. Frank Guenther is also supported by the Alfred P. Sloan Foundation.

¹In the International Phonetic Alphabet, this phoneme is indicated by the symbol /r/. We use the symbol /r/ here because it is more widely recognized and is unambiguous with regard to American English.

²Although the direction of the effect on *F3* for all of these vocal tract shape manipulations should be the same, the magnitude of the effect on *F3* will be different for each manipulation. We account for this by including the magnitudes of the effects when analyzing the combined acoustic effect of these movements; this is done through the A_i terms in Eqs. (1)–(3) below.

³The vocal tract area functions were provided by Abeer Alwan and colleagues from the Electrical Engineering Department at the University of California, Los Angeles.

⁴The palatal traces for subjects 2 and 4 were slightly misaligned relative to the tongue transducer data. To correct for this, the palatal trace for these subjects have been raised approximately 3 mm relative to the tongue transducer positions in all figures.

⁵Because any form of curve through the three tongue transducer points would represent only a rough approximation to the actual tongue shape, the curves in Fig. 17 were simply hand-drawn using computer drawing software.

⁶The lower incisor location was not available for subject 2, so the vocal tract outlines for this subject in Fig. 17 are based on a lower incisor position estimated from the lower lip position.

Alwan, A., Narayanan, S., and Haker, K. (1997). “Toward articulatory-acoustic models for liquid approximants based on MRI and EPG data. Part II. The rhotics,” *J. Acoust. Soc. Am.* **101**, 1078–1089.

Boyce, S., and Espy-Wilson, C. Y. (1997). “Coarticulatory stability in American English /r/,” *J. Acoust. Soc. Am.* **101**, 3741–3753.

Browman, C., and Goldstein, L. (1990a). “Tiers in articulatory phonology, with some implications for casual speech,” in *Papers in laboratory phonology. I: Between the grammar and physics of speech*, edited by J. Kingston and M. E. Beckman (Cambridge University Press, Cambridge, UK), pp. 341–376.

Browman, C., and Goldstein, L. (1990b). “Gestural specification using dynamically-defined articulatory structures,” *J. Phonetics* **18**, 299–320.

De Jong, K. J. (1997). “Labiovelar compensation in back vowels,” *J. Acoust. Soc. Am.* **101**, 2221–2233.

Delattre, P., Freeman, D. C. (1968). “A dialect study of American r’s by x-ray motion picture,” *Linguistics* **44**, 29–68.

Espy-Wilson, C. Y., and Boyce, S. E. (1994). “Acoustic differences between “bunched” and “retroflex” variants of American English /r/,” *J. Acoust. Soc. Am.* **95**, 2823.

Espy-Wilson, C. Y., Narayanan, S., Boyce, S. E., and Alwan, A. (1997). “Acoustic modeling of American English /r/,” *Proceedings of Eurospeech*, September 1997, Rhodes, Greece (ESCA, Grenoble, France), Vol. 1, pp. 393–396.

Fowler, C. A. (1986). “An event approach to the study of speech perception from a direct-realist perspective,” *J. Phonetics* **14**, 3–28.

Fowler, C. A. (1996). “Listeners do hear sounds, not tongues,” *J. Acoust. Soc. Am.* **99**, 1730–1740.

Guenther, F. H. (1995). “A modeling framework for speech motor development and kinematic articulator control,” *Proceedings of the XIIIth International Conference of Phonetic Sciences* (KTH and Stockholm University, Stockholm, Sweden), Vol. 2, pp. 92–99.

Guenther, F. H., Hampson, M., and Johnson, D. (1998). “A theoretical investigation of reference frames for the planning of speech movements,” *Psychol. Rev.* **105**, 611–633.

- Hagiwara, R. (1994). "Three types of American /r/," UCLA Working Papers in Phonetics **88**, 63–90.
- Hagiwara, R. (1995). "Acoustic realizations of American /r/ as produced by women and men," UCLA Working Papers in Phonetics **90**, 1–187.
- Jakobson, R., Fant, G., and Halle, M. (1951). *Preliminaries to speech analysis: The distinctive features and their correlates* (MIT Press, Cambridge, MA).
- Johnson, K., Ladefoged, P., and Lindau, M. (1993). "Individual differences in vowel production," J. Acoust. Soc. Am. **94**, 701–714.
- Lieberman, A. M., Cooper, F. S., Shankweiler, D. P., and Studdert-Kennedy, M. (1967). "Perception of the speech code," Psychol. Rev. **74**, 431–461.
- Lieberman, A. M., and Mattingly, I. G. (1985). "The motor theory of speech revisited," Cognition **21**, 1–36.
- Maeda, S. (1990). "Compensatory articulation during speech: Evidence from the analysis and synthesis of vocal tract shapes using an articulatory model," in *Speech Production and Speech Modeling*, edited by W. J. Hardcastle and A. Marchal (Kluwer Academic, Boston), pp. 131–149.
- Ong, D., and Stone, M. (1998). "Three-dimensional vocal tract shapes in /r/ and /l/: A study of MRI, ultrasound, electropalatography, and acoustics," Phonoscope **1**, 1–13.
- Perkell, J., Cohen, M., Svirsky, M., Matthies, M., Garabieta, I., and Jackson, M. (1992). "Electromagnetic midsagittal articulometer (EMMA) systems for transducing speech articulatory movements," J. Acoust. Soc. Am. **92**, 3078–3096.
- Perkell, J. S., Matthies, M. L., Lane, H., Guenther, F. H., Wilhelms-Tricarico, R., Wozniak, J., and Guiod, P. (1997). "Speech motor control: Acoustic segmental goals, saturation effects, auditory feedback and internal models," Speech Commun. **22**, 227–250.
- Perkell, J. S., Matthies, M. L., and Svirsky, M. A. (1994). "Articulatory evidence for acoustic goals for consonants," J. Acoust. Soc. Am. **96**, 3326.
- Perkell, J. S., Matthies, M. L., Svirsky, M. A., and Jordan, M. I. (1993). "Trading relations between tongue-body raising and lip rounding in production of the vowel [u]: A pilot "motor equivalence" study," J. Acoust. Soc. Am. **93**, 2948–2961.
- Perkell, J. S., Matthies, M. L., Svirsky, M. A., and Jordan, M. I. (1995). "Goal-based speech motor control: A theoretical framework and some preliminary data," J. Phonetics **23**, 23–35.
- Saltzman, E. L., and Munhall, K. G. (1989). "A dynamical approach to gestural patterning in speech production," Ecological Psychol. **1**, 333–382.
- Savariaux, C., Perrier, P., and Orliaguet, J. P. (1995a). "Compensation strategies for the perturbation of the rounded vowel [u] using a lip tube: A study of the control space in speech production," J. Acoust. Soc. Am. **98**, 2428–2442.
- Savariaux, C., Perrier, P., and Schwartz, J. L. (1995b). "Perceptual analysis of compensatory strategies in the production of the French rounded vowel [u] perturbed by a lip tube," in *Proceedings of the XIIIth International Congress of Phonetic Sciences* (KTH and Stockholm University, Stockholm, Sweden), Vol. 3, pp. 584–587.
- Stevens, K. (1998). *Acoustic Phonetics* (MIT Press, Cambridge, MA).
- Westbury, J. R., Hashi, M., and Lindstrom, M. J. (1995). "Differences among speakers in articulation of American English /r/: An x-ray microbeam study," *Proceedings of the XIIIth International Conference of Phonetic Sciences* (KTH and Stockholm University, Stockholm, Sweden), Vol. 4, pp. 50–57.
- Westbury, J. R., Hashi, M., and Lindstrom, M. J. (1998). "Differences among speakers in lingual articulation of American English /r/," Speech Commun. **26**, 203–226.