

Provinsi wana wana di maka di manana manana di mari ata dataka ata dataka na manana manana matana -

Approved for public release; distribution unlimited

٠.

INTERACTIVE VOICE TECHNOLOGY: VARIATIONS IN THE VOCAL UTTERANCES OF SPEAKERS PERFORMING A STRESS-INDUCING TASK

James D. Mosko Naval Aerosyace Medical Research Laboratory

> Kenneth N. Stevens '/ Hancock Place Cambridge, Massachusetts 02139

> > and

Glenn R. Griffin Naval Aerospace Medical Research Laboratory

62758N MF58528 MF5852801B 0001

<u>Reviewed by</u> Ashton Graybiel, M.D. Chief Scientific Advisor

Approved and Released by Captain W. M. Houk, MC, USN Commanding Officer

16 August 1983

NAVAL AEROSPACE MEDICAL RESEARCH LABORATORY NAVAL AIR STATION PENSACOLA, FLORIDA 32508

SUMMARY PAGE

THE PROBLEM

If speech recognition systems are to be successfully implemented in advanced military aircraft and other military weapons platforms, they must be able to respond correctly to voice commands produced under wide-ranging operational conditions. Information is needed concerning the modifications of speech behavior to be expected under various operational conditions so that knowledge about such vocal variability can be incorporated into the design of appropriate recognition systems.

FINDINGS

Detailed acoustical analyses were conducted of the words produced by four speakers undergoing a motion disorientation-inducing performance task. The results indicated that the speakers differed markedly in the types and magnitudes of the changes that occurred in their speech. For at least some of the speakers, the stress-inducing experimental condition caused an increase in fundamental frequency, changes in the pattern of vocal fold vibration, shifts in vowel production, and changes in the relative amplitudes of sounds containing turbulence noise. All of the speakers showed greater variability in their production of the words in the experimental condition than in a more relaxed control situation. This variability was manifested in the acoustical characteristics of individual phonetic elements, particularly in speech sounds occurring in the vicinity of unstressed syllables. The kinds of changes and variability observed serve to emphasize the limitations of speech recognition systems based on template matching. There is need for a better understanding of these phonetic modifications and for developing ways of incorporating knowledge about these changes into speech recognition systems.

INTRODUCTION

ግም ርጭር ጭን

Of the many possible military applications of speech recognition devices, the high-performance aircraft weapons platform appears to be one environment for beneficial integration. In this often "hands-busy", "eyes-busy" environment, alternate methods of systems query and control could enhance aircrew performance (1). This platform also offers many challenges to implementing speech recognition technology because the physical, cognitive, and psychological demands on the pilot are changing constantly. Often the changes are extreme, transitory, and unpredictable. Under such conditions the speech behavior of pilots will exhibit wide variations which can degrade voice communication performance with speech recognition devices. As plans develop to incorporate interactive voice systems in advanced military aircraft, it is important to determine how the speech of a pilot is likely to be modified in stressful operational situations.

Two of the authors have reported elsewhere (2, 3) the results of two experiments wherein attempts were made to document some of the vocal changes occurring in stressful conditions encountered in aviation -- high-noise levels (2) and motion disorientation (3). The conditions were simulated in the laboratory and variations in voice fundamental frequency, word token duration, and voice amplitude were reported. Since such variations rarely occur in isolation, it is important to understand other acoustic manifestations of such changes. In this report, we present a more detailed acoustical analysis of the vocal utterances of four selected speakers from one of the earlier experiments (3). The results of this analysis will serve as a guide for future work in defining the variations in speech produced by aviators in strens-inducing environments.

METHOD

The vocal utterances of four young adult male volunteer subjects, recorded as the subjects performed on the Visual-Vestibular Interaction Test (VVIT) (4), were subjected to detailed acoustical analyses. None of the subjects exhibited any hearing, speech, or physical condition to preclude their participation.

Description of VVIT.

In the VVIT individual subjects are seated in a blacked-out rotatable device facing a front lighted 17x17cm matrix. The matrix contains coordinate letters and digits on the left and upper margins, respectively, and randomly arranged digits in the body of the matrix. The subject receives an aurally presented coordinate set cue and his task is to locate the intersection of the coordinates in the body of the matrix, verbally report the digit at the intersection and the next two digits immediately below it in the same column. A typical cue is "A-2", and a typical verbal response is "six...five...nine." Each subject performs the task twice: once in a stationary, STATIC, mode and once in an oscillating, DYNAMIC, mode (sinusoidal oscillation, 0.02 Hz, 30 rpm peak).

The motion stimulus in the DYNAMIC mode involves two aspects of motion stress: 1) vestibular degradation of visual performance and 2) motion sickness. The first is generated by the vestibulo-ocular reflex (VOR) which tends to drive the eyes relative to the head-fixed display, thereby degrading visual target acquisition. With the sinusoidal stimulus, this aspect is definitely cyclic; the VOR and visual performance degradation reach a maximum twice in each cycle (corresponding roughly to the peak angular velocities), and diminish approximately to zero twice each cycle, as zero angular velocity is approached. In some individuals, exposure to this form of visual-vestibular interaction induces motion sickness which can build to the point of emesis within a 5-minute exposure (4).

In the present experiment, subjects were permitted to terminate the DYNAMIC mode if they felt severe motion sickness would result. Two subjects (speakers 2 and 4) completed both modes of the VVIT and two speakers (speakers 1 and 3) terminated the DYNAMIC mode prior to a complete trial.

Target Words.

The subjects were presented a target word preceding the coordinate cue. The subjects were instructed to repeat the target word and then to report the digits corresponding to the coordinate cue. A typical cueing sequence was "TACAN...A...2" and a typical subject response was "TACAN...seven...nine... two". The ten target words are listed in Table I. Two randomized lists of the target words and correct digit sequences are listed in Table II. Fortythree items per condition were used to conform to the time constraints of the VVIT.

Table I

altitude	marker	
bogev	monitor	
contact	pattern	
heading	radar	
holding	tacan	

List of target words used in the experiment

ACOUSTIC ANALYSIS

The number of utterances available for acoustic analysis was approximately 400 under the CONTROL condition, 160 under the STATIC condition, and 120 under the DYNAMIC condition. Wideband spectrograms, 0 - 5000 Hz frequency range, were made of all of the words produced under the STATIC and DYNAMIC conditions, and for a selected subset of the words in the CONTROL condition. Various measurements and observations were made from all of these spectrograms,

							······
	STATIC	CONDITION			DYNAM	C CONDITION	N
	Cue	Carrier	Digits		Cue	Carrier	Digits
1.	A-Z	Pattern	342	1.	G - 7	Radar	476
2.	H-10	Bogey	914	2.	D-5	Altitude	176
3.	L-7	Heading	765	3.	в-4	Holding	824
4.	K-9	Monitor	939	4.	L-1	Pattern	496
5.	C91)	Marker	129	5.	J-3	Bogey	434
6.	H-12	Altitude	273	6.	G-2	Heading	763
7.	J-3	Heading	434	7.	A-5	Radar	427
8.	E-10	Pattern	965	8.	K-12	Monitor	793
9.	B-4	Holding	824	9.	G-4	Contact	143
10.	A-10	Radar	813	10.	I-8	Tacan	487
11.	G-1	Bogey	315	11.	A-3	Holding	181
12.	L-2	Altitude	654	12.	I-4	Bogey	948
13.	J-10	Heading	342	13.	B- 3	Altitude	443
14.	E-4	Radar	415	14.	L-2	Contact	654
15.	F-3	Marker	677	15.	C-1	Radar	831
16.	I-2	Bogey	836	16.	K-3	Tacan	766
17.	G-9	Heading	566	17.	L-5	Marker	549
18.	C-1	Tacan	831	18.	I-7	Bogey	283
19.	C-10	Monitor	251	19.	L-3	Altitude	5/37
20.	F-8	Marker	167	20.	E-12	Heading	685
21.	J4	Holding	674	21.	K-1	Monitor	479
22.	L-1	Tacan	496	22.	I-12	Altitude	652
23.	K-12	Contact	793	23.	D-10	Monitor	987
24.	E-8	Radar	159	24.	D-2	Bogey	717
25.	G5	Holding	631	25.	E-7	Marker	657
26.	I-9	Radar	521	26.	G-10	Heading	547
27.	B-1	Contact	369	27.	1-5	Contact	368
28.	D-3	Tacan	298	28.	B-12	Tacan	697
29.	H-4	Holding	142	29.	H-2	Radar	471
30.	D-9	Pattern	192	30,	F-9	Monitor	523
31.	K-1	Tacan	479	31.	H-10	Holding	914
32.	A-7	Monitor	1 3 4	32.	C-4	Contact	438
33.	G-10	Bogey	547	33.	L-4	Tacan	185
34.	E7	Contact	657	34.	A-9	Marker	641
35.	H-12	Monitor	273	35.	L-12	Pattern	968
36.	I-4	Altitude	948	36.	J-10	Heading	342
37.	D-2	Marker	717	37.	E-4	Fattern	415
38.	A-2	Altitude	342	38.	H-7	Holding	147
39.	J-3	Contact	434	39.	I-9	Marker	521
40、	B-1.0	Pattern	431	40.	F-1	Pattern	985
41.	J-7	Bogey	4 2 8	41.	A-7	Radar	134
42.	D-8	Radar	629	42.	G5	Bogey	631
43.	F-3	Monitor	677	43.	K-7	Monitor	686

Гab	1e	ΪĪ
_	_	

Sequence of Words Presented to Speakers in STATIC and DYNAMIC Conditions

as described in detail later. For selected utterances, additional analysis and observation was carried out using a computer. In order to do this additional analysis, speech waveforms were displayed in some cases, discrete Fourier transforms at selected points in the utterances were calculated and displayed, measurements of fundamental frequency were made, formant frequencies were determined using a linear predictive coding (LPC) procedure (5), and measurements of durations of certain speech events were made. All the computer-based data were obtained with the waveform low-pass filtered (4.8 kHz), sampled at 10 kHz, and first-differenced. Measurements of every property or parameter for each utterance were not attempted. A more limited goal was to make a sufficient number of measurements and observations to allow a cataloguing of the kinds of changes that occurred in the utterances under the various experimental conditions.

In order to guide the data analysis, we recognized that two broad categories of change can occur in an utterance when it is repeated by a speaker under different conditions. One of these is associated with overall changes in the "posture" or "state" of the speaker's speech production system, and these changes are reflected in certain average characteristics of the utterances. Thus, for example, a speaker may speak with greater effort by using a higher subglottal pressure. This higher pressure could lead to a greater overall amplitude of the speech sounds, a modification of the average spectrum of the glottal output, and a possible modification of sounds produced with abrupt release or with turbulence noise. Or, a greater overall tension of the vocal folds or a change in the characteristics of the vocal-fold surfaces could lead to a modification of the frequency, spectrum, and regularity of the glottal output.

A second type of change that can occur when a speaker repeats an utterance under different conditions is a modification of particular phonetic elements within the utterance. In English (as in any language) a number of rules can be applied optionally to specify the way in which particular phonetic elements may be produced in specified phonetic environments. These rules generally apply to so-called redundant features (i.e., phonetic features which, by themselves, are not utilized for signaling a phonetic distinction in the language, but which can provide cues a listener might use in decoding an utterance), in addition to the cues associated with features that are distinctive. These types of optional rules are illustrated in the following examples:

- A word-initial /b/, /d/, or /g/ in English can be optionally prevoiced (e.g., the initial /b/ in the word bogey);
- An unstressed vowel between two voiceless consonants can be optionally produced as a voiceless vowel (e.g., the unstressed vowel in altitude);
- An utterance-final stop consonant may or may not be released (e.g., the final /d/ in <u>altitude</u>);

 A /t/ preceding an unstressed vowel may be optionally produced as a flap (e.g., in the word monitor).

As a preliminary to a detailed discussion of the data, and to indicate the acoustic structure of the words used in the experiment, we show in Figure 1 a spectrogram of each word produced by one of the speakers. Eight of the words are bisyllabic, all with stress on the first syllable, and two of the words have three syllables. The words contain a variety of stop and nasal consonants in initial, intervocalic, and final position, and both back and front vowels are represented in the corpus.

RESULTS AND DISCUSSION

AVERAGE PROPERTIES OF UTTERANCES

Fundamental frequency

In the utterances produced under the CONTROL condition, the fundamental frequency (F_0) was usually considerably lower than in the two experimental conditions. This difference presumably arose for several reasons: 1) the CONTROL utterances were generated in a phrase following the carrier word say; 2) these utterances were generated in rapid succession in the form of a list; and 3) in the experimental condition the words were followed by a number sequence. In view of these differences, we shall not compare F_0 patterns for GONTROL and experimental conditions, but only within the different experimental conditions.

The values of F_0 to be reported here were obtained by measuring the glottal period for three successive periods at a selected point in the word (usually 50 msec following the onset of the vowel selected for study) and averaging the reciprocals of these three numbers. There was some variation in the average F_0 from word to word throughout an experimental run, presumably because of variations in intrinsic F_0 from vowel to vowel, and because of influences of voicing characteristics of adjacent consonants. Consequently, in examining the effect of experimental conditions on F_0 , we will either observe the F_0 for particular words throughout an experimental run, or we shall average over groups of vords.

There were considerable differences among speakers in the way F_0 varied throughout the STATIC and DYNAMIC experimental runs. Selected data for the four speakers are presented in Figure 2. For Speaker 2, there was essentially no change in F_0 as the experiment proceeded through the STATIC and DYNAMIC conditions. Data for this speaker are shown for both stressed and unstressed vowels. For the other three speakers, there was some rise in F_0 throughout the experiment, particularly in the DYNAMIC condition. In the case of Speaker 1, for whom F_0 measurements were made in both stressed and unstressed vowels, the increase in F_0 for the DYNAMIC condition seemed to be more marked for the stressed vowels. The speaker showing the greatest F_0 rise appeared to be Speaker 3. Figure 3 displays the F_0 values for the stressed vowels of this speaker, averaged over successive groups of five words, throughout the experiment. There appears to be a gradual rise in F_0 as the



Figure 1 (Part a.). See caption Figure 1 (Part b).

K K



Figure 1 (Part b). Spectrograms of one example of each of the ten test words produced by Speaker 4. The number identifying each word is listed in Table 2, where S = STATIC and D = DYNAMIC.







日本のないようというないのないである

Figure 2 (Part b). Examples of fundamental frequency changes in selected words throughout an experimental session for each of the four speakers. Each point represents the mean F_0 for a single token over three successive glottal periods, sampled at a point 50 msec from the onset of the vowel. The numbers on the abscissa represent token numbers (from Table 2) for the STATIC (S) and DYNAMIC (D) conditions. In the case of Speakers 1 and 2, data are shown for both the stressed vowel (solid dots) and the final vowel (open circles). These data show that the stress-inducing conditions of the experiment have different effects on different speakers, and that the effect on the F_0 of the stressed vowels tends to be greater than the unstressed vowels.



Figure 3. Fundamental frequency for stressed vowels in words produced by Speaker 3 in the course of an experimental session. F_0 values for each vowel were obtained by procedures indicated in Figure 2 and in the text. Each point usually represents the average of these F_0 values over the vowels in five successive words in the experiment, as indicated on the abscissa.

speaker progresses through the STATIC condition, and F_0 remains at the higher value in the DYNAMIC condition until the experiment is terminated at item 19D. The local peaks in the curve could represent brief intervals in which a transient change occurred in the speaker, possibly as a result of increased stress.

These increases in F_0 for some speakers as the experiment progressed presumably reflect either a rise in vocal-fold tension or a rise in subglottal pressure, or an overall increase in tension that occurs throughout the articulatory, laryngeal, and lower respiratory systems.

State and configuration of vocal folds

The principal acoustic consequence of a change in the tension of the vocal folds is a change in the frequency of vocal-fold vibration. However, another consequence of a modification in the state of the vocal-fold surfaces or in the configuration of the vocal folds is an alteration in the waveform of the volume velocity source at the glottis. This kind of change is more difficult to quantify, since the glottal waveform is filtered by the vocal tract, and the sound wave provides only indirect evidence for this waveform.

For one speaker (Speaker 4) informal listening to the recording for the STATIC and DYNAMIC conditions indicated that a change in voice quality occurred as the utterances for the DYNAMIC condition were produced. An attempt was made to establish some kind of acoustic manifestation of this change, through observation of the waveforms or spectra of the vowels. Figure 4 shows the waveform and the spectrum sampled 50 msec from onset of voicing for the vowel $/\partial C/$ in the word <u>pattern</u> produced by Speaker 4. The two utterances of the word were at the beginning of the STATIC condition, and near the end of the DYNAMIC condition. Also shown in the figure are spectrograms of these two words.

Probably the most striking difference between the two utterances is in the waveform. For the DYNAMIC condition, there is a more rapid decay in the first-formant oscillation in the early part of the glottal cycle, and the oscillation is almost extinguished in the later part of the cycle (presumably during the most open phase of the glottal cycle). In the spectra, the main difference is in the relative amplitude of the lower harmonics, particularly the first. Thus for utterance 40D, the amplitude of the first harmonic is about 12 dB below the peak in the first formant (F1), whereas in item 1S this difference is about 18 dB. In the spectrogram there is evidence for a "filling in" of spectral energy at low frequencies, below F1, in item 40D.

These acoustic differences can be ascribed to changes in the glottal configuration for the two utterances. In the case of 40D, it is probably that the vocal folds are more abducted, so there is never a complete closure of the glottis during the cycle. This configuration would lead to greater acoustic losses in the region of F1, as indicated by the more rapid decay of the waveform.



The present study did not examine this kind of modification in detail for every speaker. For Speakers 1, 2, and 3, informal observations of waveforms and spectra of the type shown in Figure 4, together with informal auditory evaluations, failed to show consistent changes in the glottal waveform or voice quality that were as marked as those illustrated for Speaker l_i in Figure 4. The figure illustrates, however, the nature of the changes that can occur in the acoustic characteristics of a vowel when there are modifications in the manner of vibration of the vocal folds.

Formant frequencies

The frequencies of the first three or four formants were measured at selected points in a number of the utterances produced by each of the speakers in the CONTROL, STATIC, and DYNAMIC conditions. An example of the kind of data that emerged from these measurements for Speaker 4 is shown in Figure 5a. The formant frequencies in each utterance of the vowel /2c/ in pattern were measured (using an LPC algorithm) at two points located approximately 40 and 60 msec from the onset of voicing. Each point in Figure 5a represents an average of these two values for F1 and formant 2 (F2) for one utterance. This figure shows that the formant frequencies for utterances of this word in the CONTROL condition are tightly clustered, indicating that the successive repetitions of the word are very similar. The scattering of the points suggests that there is greater var' bility in the utterances for the STATIC condition, and still greater variability for the DYNAMIC condition. Although it is difficult to draw firm conclusions on the basis of so few utterances, there is a tendency for F2 (and, to some extent, F1) to be higher for the experimental conditions than for the CONTROL condition. A possible interpretation is that the larynx is positioned slightly lower for the more relaxed CONTROL condition, leading to a longer vocal tract and hence lower formant frequencies. Similar tendencies can be observed for the first vowel in the word tacan for Speaker 4, shown in Figure 5b, although in this case the points for the CONTROL condition are not quite as tightly clustered.

Examples of the formant frequency data for the other speakers are given in the various panels of Figure 6. The formant values were obtained in the same way as described above, i.e., each point represents the average of two measurements spaced 20 msec apart. In the stressed vowel in the word <u>tacan</u>, Speaker 1 shows essentially the same amount of variability in F2 for the CONTROL and the experimental conditions. For the /C/ in <u>heading</u>, on the other hand, the points for the CONTROL condition are more tightly clustered. For both vowels, Fl seems to be lower for the CONTROL condition than for the experimental conditions. In the two vowels of <u>altitude</u> for Speaker 1, there is again substantial variability in F2 for the experimental conditions. For all the utterances of Speaker 1 there is no obvious trend in the data when STATIC and DYNAMIC conditions are compared.

The limited amount of data collected for Speaker 2 suggest that he was less variable in his vowel productions during the experimental conditions than were the other speakers. Speaker 3 shows more variability, particularly in the stressed vowel in <u>altitude</u>. Apparently the influence of the following /1/ was different from one repetition to the next.



Date on formant frequencies Fl and F2 for the stressed vowel in two words produced by Speaker variability in vowel production for the CONTROL condition than for the experimental conditions, as weil measurements spaced about 20 msec apart, centered 50 msec from vowel onset. There appears to be less Formant frequencies represent averages of two as a tendency for F1 and F2 to increase in the stress-inducing situation. 4 in the CONTROL, STATIC, and DYNAMIC conditions. Figure 5.



Figure 6. Examples of formant frequency data from several words produced by Speakers 1, 2, and 3 in the CONTROL condition (squares), the STATIC condition (triangles), and the DYNAMIC (circles). Formant measurements are obtained by the method indicated in Figure 5.

ł

The principal conclusions to be drawn from the formant-frequency data are: 1) there is more variability in vowel formant frequencies in successive productions of words in the STATIC and DYNAMIC conditions than in the CONTROL condition, with Fl showing a range up to 100 to 200 Hz, and F2 a range up to 200 to 300 Hz; 2) for some speakers and some vowels there are systematic shifts in formant frequencies from the CONTROL condition to the experimental conditions, suggesting that there is a shift in "posture", possibly a change in positioning of the laryngeal structures.

Frication and aspiration noise; stop bursts

ι.

A number of speech sounds are produced by creating turbulence at a constriction in the vocal tract, thereby causing random noise to be generated in the vicinity of the constriction. The noise is usually called frication noise if it is generated primarily at a narrow constriction in the oral cavity, and is called aspiration noise if the vocal tract is relatively unconstricted, the glottis is somewhat open, and noise is produced primarily in the vicinity of the glottis. When a stop consonant is released, a burst of frication noise is usually generated in the vicinity of the constriction, and this may be accompanied by a transient acoustic excitation of the vocal tract as the intra-oral pressure is abruptly released.

These various kinds of sources of excitation of the vocal tract could potentially undergo some modification as the speaker is exposed to a stressful situation. For example, a change in the respiratory force, giving rise to a modification in the subglottal pressure, could influence the amplitude and spectrum of the turbulence noise. The detailed configuration of the constriction and the state of the surfaces forming the constriction could also influence the turbulence noise. The properties of the burst for a stop consonant could be affected by the abruptness of the stop release, which in turn may be determined by the state of the surfaces of the structures forming the constriction.

To investigate these possible effects, several measurements were made on noise portions of selected utterances in the corpus. One set of measurements was made on the noise bursts at the /t/ and /k/ releases in the word tacan. The method for obtaining quantitative measures of the burst is illustrated in In the case of the /t/ burst, the spectrum was sampled with a 25.6 Figure 7. Hamming window centered 5 msec following the release, and another spectrum was obtained with a similar window centered at the onset of the third glottal period following the beginning of normal voicing. Both spectra were smoothed using a 14-pole LPC procedure. The amplitude of the spectral peak nearest the fourth formant region in the smoothed spectrum was determined for each spectrum, and the differences between these amplitudes for the burst and for the vowel onset were obtained. A similar procedure was followed for the /k/burst and the vowel immediately following this burst, except that in this case the amplitudes that were measured were the spectral peaks in the third formant region.



Figure 7. Illustration of the measurement of burst amplitude for the two stop consonants in the word <u>tacan</u>. The upper left panel shows LPC spectra sampled 1) immediately following (5 msec) the onset of the /t/ burst, and 2) at the third pitch period of the following vowel. The lower left panel shows similar spectra for the /k/ burst in the same word. The sampling points are indicated in the spectrogram at the right. The relative burst amplitude for the /t/ burst is defined as the amplitude of the F4 peak of the burst relative to that in the spectrum of the vowel (i.e., the F4 peaks shown in the upper left panel). In the case of the /k/ burst, the measure used is based on the amplitudes of the F3 peaks, as indicated.

Data for Speaker 2 from the CONTROL and experimental situations are summarized in Table 3. Two trends are evident in these data when one examines the amplitude of the burst relative to the adjacent vowel (the differences in the table). One is that the amplitudes of the bursts in relation to the vowels are, on the average, greater in the experimental conditions than in the CONTROL conditions. This difference is especially evident for the /t/ bursts. The other is that there is more variability in this relative burst amplitude in the experimental utterances than in the CONTROL atterances. А conclusion to be drawn from these data is that this speaker, in the experimental condition, is using a more abrupt release for the stop consonants (possibly with a higher respiratory effort) and he is more variable in the way he implements these articulatory and respiratory gestures.

Table III

Amplitudes of Spectral Peaks (in dB) in Burst and in Adjacent Vowel for the Stop Consonants in Several Repetitions of the Word Tacan by Speaker 2

Control Utterances	F4 peak in /t/	F4 peak in first /2C/	Difference	F3 peak in /k/	F3 peak in second /XC/	Differ- ence
1	30	26	4	26	40	-14
2	29 29	33	4	32 21	37 41	-5
4	29	30	-1	27	39	-12
5	29	34	-5	31	37	-6
6	27	33	-6	30	37	-7
7	30	37	-7	27	37	-10
8	3 C/	35	5	22	35	13
mean			3.2			-10.9
s.d.			3.3			4.6
Experiment Utterances	al					
18S	33	28	5	22	44	-22
225	33	30	3	36	42	-6
285	30	29	1	34	40	-6
315	35	27	8	26	39	-13
10D	34	27	7	41	41	0
16D	33	26	7	31	41	-10
28D 33D	38 20	27	9	42	41 54	1
μει	27	22		20	44	-0
mean			4.5			-7.8
s.d.			4.1			6.9

There were substantial individual differences between speakers, both in the relative burst amplitudes for the two stops in <u>tacan</u> and in the amount of variability from one repetition to the next in the experimental conditions. For example, Table 4 shows that Speaker 4 had a relatively weak /t/ burst (in the 74 region) in relation to the F4 amplitude of the vowel, but was quite

Table IV

Man and Standard Deviation of Relative Burst Amplitude (in dB) of /t/ and /k/ in tacan for Repetitions in Experimental Conditions, for Three Different Speakers. Method of Measuring Relative Amplitude is Shown in Figure 7.

	Speaker 2	Speaker 3	Speaker 4
Relative amplitude of /t/ burst s.d.	4.5	-5.8	-9.1
	4.1	5.8	3.3
Relative amplitude of /k/ burst	-7.8	9.3	3.9
s.d.	6.9	6.3	5.2

consistent in producing this burst amplitude (s.d. of 3.3 dB). On the other hand, this speaker (as well as Speaker 3) had a much higher amplitude /k/ burst than did Speaker 2.

In order to illustrate some other aspects of turbulence noise generation, measurements were made of aspiration noise in the initial /h/ in heading, and of the /p/ burst in pattern. Measurements of the /p/ burst were made in a manner similar to the t/and k/bursts in tacan (Figure 7). The amplitudes of the spectral peaks in burst and in vowel in the F3 region were used as the measure. In the case of the /h/, the noise spectrum was sampled at a point 15 msec prior to the onset of voicing in $/\xi/$. The spectrum was smoothed (as before) and the amplitude at the peak in the F3 region was measured. The following vowel spectrum was sampled at the third glottal period, and again the F3 peak was measured. The data for three speakers are summarized in Table 5. Again we see large individual differences in the amplitude of the noise in relation to the vowel and in the amount of variability shown by the different speakers. Speaker 4 has the greatest relative amplitude of the noise, but also shows considerable variability in the noise generation mechanism.

Table V

	Speaker 2	Spf.aker 3	Speaker 4
Amplitude of F3 peak in /h/	20	16	2.6
Amplitude of F3 peak in <i> E</i> /	38	29	33
Difference s.d. of difference	-18(8) 3.6	-13(3) 0.5	-7(8) 5.9
Amplitude of F3 peak in /p/	19	15	24
Amplitude of F3 peak in /æ/	34	35	30
Difference s.d. of difference	-15(8) 4.3	-20(5) 3.9	-6(8) 6.4

Amplitudes (in dB) of Spectral peaks in Noise and in Following Vowel for /h/ in <u>heading</u> and /p/ in <u>pattern</u>, for Experimental Conditions. The Numbers in Parentheses represent the Number of tokens Measured.

In summary, then, these samples of data from bursts and aspiration noise show rather large fluctuations in the relative amplitude of the noise during the experimental conditions in which the speakers are placed under stress. Data from one speaker cuggest that the variability is less when the test words are being repeated under more relaxed conditions, and that the mean values of the noise amplitude are different under these conditions. The data also indicate rather large individual differences in relative spectral amplitude of the noise for a given speech sound from one speaker to another.

Timing and segment durations

A number of measurements were made of the time between various acoustic boundaries in several of the words. The aim was to determine whether there were systematic changes in the temporal characteristics of the utterances between the CONTROL condition and the two experimental conditions. The measurements included the durations of stressed vowels in syllable-initial position, durations of vowels with secondary stress in syllable-final

position, voice-onset time for voiceless stop consonants, and the duration of the stop gap for medial voiceless stop consonants.

Some of the results are summarized in Table 6. In general, there were no systematic differences in durations for the STATIC and DYNAMIC conditions, and consequently the data for these two experimental conditions were averaged. The main effects emerging from these data are: 1) final vowels in the words are consistently longer in the experimental condicions than in the CONTROL condition; and 2) there is greater variability in vowel durations in the experimental conditions. These two results are presumably a consequence of the fact that the CONTROL words were read rather rapidly as a list (each word preceded by say), with the speaker in a relatively stable physiological state. The data also show greater variability in final vowel durations than in initial vowel durations, both in the experimental and CONTROL conditions. This result is not unexpected, since the amount of syllable-final lengthening occurring in this situation is presumably not well controlled. Apparently, as long as some final lengthening is produced by a speaker to mark the end of an utterance, the actual amount of lengthening is not crucial.

Further detailed examination of the data in Table 6 indicates trends that characterize the timing of some speakers but not others. For example, there is a tendency for Speaker 3 to lengthen voice-onset time (VOT) in the experimental conditions, whereas Speaker 1 tends to shorten it. There are also substantial differences in the way the consonant /k/ preceding the second vowel in <u>tacan</u> is produced. For two speakers (3 and 4), the VOT is long relative to that for the initial /t/, whereas for the other two speakers, the reverse is true.

In summary, these data on timing provide additional evidence for the observation that speakers show more variability when they are placed under stress in the experimental conditions than in the more relaxed CONTROL condition. There is little evidence (beyond the syllable-final lengthening), however, for significant shifts in timing strategies between the various CONTROL and experimental situations. Some shifts in durations are exhibited by some speakers when they are in the experimental situation, but these changes in timing are not very large.

MODIFICATIONS OF INDIVIDUAL PHONETIC ELEMENTS

A number of modifications were observed in individual phonetic elements in repetitions of words in the STATIC and DYNAMIC conditions. These modifications are usually the result of application of optional low-level phonetic rules in English. Apparently when a speaker is producing utterances under various amounts of stress, fluctuations in the application of these rules give rise to phonetic modifications. Most of the phonetic modifications described below were observed on the spectrograms made of the utterances.

Medial unstressed vowel in "altitude"

The medial unstressed vowel in the word altitude lies between two

Table VI

Average Durations of Various Speech Events for Words Produced in CONTROL and Experimental Situations. Values are in Milliseconds. Numbers in Parentheses are Number of Tokens Measured. For Experimental Conditions, Data from STATIC and DYNAMIC Situations are Averaged.

hÌ

			CONTRO	<u>)L</u>	EXPERIMEN	TAL
			Duration	<u>s.d.</u>	Duration	s.d.
Stressed v	owels	3				
Speaker 1	/20/ /20/ /01/	tacan pattern holding	126 135 148	6 (4) 7 (6) 9 (9)	131 157 168	9 (5) 6 (4) 16 (4)
Speaker 2	/ 2 ///	tacan holding	115 172	6 (5) 8 (5)	115 158	9 (7) 19 (8)
Speaker 3 Speaker 4	/28/ /28/ /01/	tacan tacan holding	146 134 201	20 (5) 6 (4) 7 (4)	142 157 209	17 (5) 13 (8) 17 (8)
		mean	147	 9	155	13
Final vowe	ls.					
Speaker 1	/æ/ /Iŋ/	contact holding	173 134	9 (4) 11 (8)	204 239	30 (7) 48 (6)
Speaker 2	/æ/ /IŊ/	tacan holding	239 106	13 (5) 28 (5)	307 234	26 (7) 23 (7)
Speaker 3 Speaker 4	/æ/ /æ/ /IN/	tacan tacan holding	260 267	28 (5) 6 (4) 11 (4)	236 284 235	27 (6) 29 (8) 30 (8)
		mean	184	15	248	30
Stop gap						
Speaker 1 Speaker 2 Speaker 3 Speaker 4	/k/ /k/ /k/ /k/	tacan tacan tacan tacan mean	71 65 74 <u>88</u> 74	0 (4) 8 (5) 6 (5) <u>8 (</u> 4) 6	85 68 88 <u>81</u> 81	9 (6) 6 (7) 4 (5) <u>7</u> (8) 7
Voice-onse	et ti	ne				
Speaker 1	/t/ /k/ /p/ /k/	tacan tacan pattern contact	67 29 41 60	5 (4) 3 (4) 10 (6) 12 (4)	55 24 46 49	14 (5) 7 (5) 13 (4) 6 (4)
Speaker 2	/t/ /k/	tacan tacan	57 39	8 (5) 16 (5)	50 45	7 (7) 5 (7)
Speaker 3	/t/ /k/	tacan tacan	58 74	6 (5) 6 (5)	62 88	20 (5) 4 (5)
Speaker 4	/t/ /k/	tacan tacan mean	54 _ <u>88</u> _ 57	$ \begin{array}{c} 12 (4) \\ -\frac{8}{8} (4) \end{array} $	46 <u>-81</u> - <u>55</u> -	10 (8) $-\frac{7}{9}$ (8)

voiceless consonants and, in English, such a vowel is subject to an optional rule: <u>devoicing</u>. Examples of two utterances of this word by the same speaker - one with the vowel voiced and the other devoiced - are shown in Figure 8. Of the 28 utterances of this word by the four speakers in the STATIC and DYNAMIC conditions, the medial vowel was voiceless in seven. There was no significant difference in the number of devoiced vowels in the two experimental conditions. Only one speaker (Speaker 4) was consistent in voicing all of these unstressed vowels. The least consistent speaker was Speaker 3, who devoiced the vowel in one-half of the utterances. In cases where the vowel was voiced, there was often considerable variability in the duration of voicing. For example, in the five utterances by Speaker 1 in which the vowel was voiced, the number of glottal vibrations occurring during the open phase of the vowel varied from two to five.

Prevoicing of initial voiced stop in "bogey"

An initial voiced stop in English may be prevoiced. In the list of words used in this study there was just one initial voiced stop - the /b/ in <u>bogey</u>. Of the 34 utterances examined for the two experimental conditions, seven showed prevoicing. Only one speaker (Speaker 3) failed to prevoice any of the utterances. There was a slight tendency for more of the /b/'s to be prevoiced in the DYNAMIC condition (five prevoicings versus two in the STATIC condition). Of the utterances that were prevoiced, there was some variation in the amplitude and duration of the prevoicing.

Release of final stops

Two of the test words ended in stop consonants (altitude and contact), and these words provided an opportunity to examine the speakers' habits with regard to final stop release. Of the 57 versions of these two words that were examined, the final stop was released in 39. Most of the unreleased stops were the voiced final stop in altitude. Some speakers consistently released the stop in contact but no speaker was consistent in releasing the stop in altitude. There was no significant difference in the incidence of final stop releases for the STATIC and DYNAMIC conditions. When the final stop was released, there was considerable variability in the amplitude of the burst at release. Figure ' shows spectrograms of three versions of contact produced by Speaker 1. In 34S, there is no stop release, in 23S there is a weak release, and in 14D a strong release. Furthermore, ir 23S the release is clearly an alveolar release, whereas in 14D, the burst shows that the release is from the velar position. For this word with a final stop cluster, then, we observe this additional variability in the place of consonant release. In none of the utterances of this word were there two successive releases corresponding to the two different places of articulation for the consonant sequence.

Velar stop before an unstressed vowel

In the word <u>bogey</u> the velar stop /g/ occurs before an unstressed vowel. Observation of spectrograms of this word indicated, for some utterances of



Figure 5. Spectrograms of the word <u>altitude</u> produced by Speaker 2 during the experimental session. These spectrograms illustrate a version of the word in which the unstressed vowel is devoiced (lower panel) and one in which there is a brief interval of voicing in the vowel (upper panel).



Figure 9. Spectrograms of three versitons of the word contact produced by Speaker 1 during the experimental session. These spectrograms illustrate three kinds of release of the final consonants: On alveolar release (top panel), no release (middle panel), and a velar telease (bottom panel). Also evident on these spectrograms are varying amounts of nasalization in the vowel preceding the nasal consonant: substantial nasalitation (bottom panel), and a small amount of nasalization (top panel). the word, the /g/ was produced without a complete closure of the vocal tract, so that a velar continuant resulted. Intermediate cases between a stop and a continuant were also observed. This "weakening" of a velar stop is not unexpected, since there is no opposition in English between a velar stop and a velar fricative, and consequently there is not a strong motivation to represent this distinction in the sound wave.

Two stamples of the word <u>bogey</u> produced by Speaker 2 are shown in Figure 10. In a of these utterances (5D), there appears to be a complete velar closure, whereas in the other (12D), the continuation of the second formant is evidence that the closure was not complete.

Of the 33 versions of <u>bogey</u> examined in the two experimental conditions, observation of the spectrograms indicated that about one-half of the velars were produced with a complete closure and a release burst, and about one-half showed no evidence of a complete closure. All speakers produced utterances of both types, and there was not a significant difference in the incidence of the two types in the STATIC and DYNAMIC conditions.

Nasalization of vowels preceding nasal consonants

In English and many other languages when a nasal consonant follows a vowel, some anticipatory nasalization is produced in the vowel. The amount and duration of this nasalization depends to some extent on the phonetic context, but some variation is to be expected within the utterances of one speaker and from one speaker to another.

The list of test words contains a number of examples of vowels followed by nasal consonants: contact, heading, holding, monitor, pattern, and tacan. The acoustic correlates of nasalization are not understood sufficiently to make it possible to establish with any precision from the sound wave the time at which velopharyngeal opening occurs, or the degree of velopharyngeal opening. The general acoustic correlates of nasalization for vowels are 1) a weakening of the first formant, 2) a broadening of the first formant, 3) introduction of an additional "nasal" formant in the vicinity of the first formant, and 4) a shift in the frequency of the first formant. From observation of spectrograms, we have attempted to determine the time at which nasalization becomes evident in each of the words listed above. This procedure is clearly subject to some uncertainty, and the results must be interpreted with this in mind.

Examples of words in which a final vowel precedes a nasal consonant are shown in Figure 11. In the case of 7S, there is evidence of nasalization in the vowel immediately following the release, whereas in 6D acoustic evidence for nasalization does not appear until about 50 msec following the /d/ release. Evidence for variability in nasalization of a vowel preceding a nasal in intervocalic position can be observed in the spectrograms of the word <u>contact</u>, shown earlier in Figure 9. In 23S there is little nasalization in the vowel, and the onset of the nasal consonant is rather abrupt. For 14D, on the other hand, nasalization can be seen over about one-half of the vowel



Figure 10. Spectrograms of two versions of the word <u>bogey</u> produced by Speaker 2 during the experimental session. The version at the left illustrates a token in which the preunstressed /g/ is produced with essentially complete vocaltract closure, whereas for the right-hand utterance complete closure is not achieved.



Figure 11. Spectrograms of two versions of the word <u>heading</u> produced by Speaker 4 during the experimental sessions. These spectrograms illustrate different degrees of nasalization for the final vowel, with the token at the left showing the greater amount of nasalization. Other examples of vowel nasalization can be seen in Figure 9.

duration (as evidenced by the appearance of an additional resonance at low frequencies) but for 34S there is an intermediate degree of nasalization.

There are large individual differences in the average time from consonant release to onset of nasalization in syllable-final vowels (such as <u>heading</u>, <u>tacan</u>, etc.). These average times are about 38, 54, 17, and 22 msec, respectively, for Speakers 1, 2, 3, and 4. Within a given speaker there are also substantial differences in onset of nasalization (as in the examples in Figures 9 and 11), but there are no consistent differences for the DYNAMIC compared with the STATIC conditions. The kind of variability that exists for a particular speaker is illustrated in Figure 12. For all of the utterances with a final nasal consonant produced by Speaker 4 (32 in all), the percent of the final vowel duration that was nasalized (as determined by spectrographic observation) was measured and displayed as a distribution in the figure. The distribution of this percentage across all 32 utterances is very broad, indicating that this speaker shows a large amount of variability in timing the onset of nasalization in these vowels when he is operating in the stress-producing experimental conditions.

Weakening of alveolar stops and nasals in pre-unstressed position

Several of the phonetic modifications discussed above are related to the "weakening" that occurs in the components of a syllable when it is unstressed. When an alveolar stop or nasal precedes an unstressed vowel, this weakening can take the form of flapping of the consonant. Within the list of utterances that were used in this study, there are six places where this optional weakening of an alveolar consonant can occur. These are listed as follows: heading, holding, monitor, pattern, and radar.

Examination of the spectrograms showed that there were at least three different ways of producing the stop consonants in these words: 1) as a full-fledged stop with a well-defined burst; 2) as a flap with a brief closure interval (less than about 20 msec); 3) as a continuant with no evidence of closure. Examples of two of these manifestations of pre-unstressed stop consonants are given in Figure 13, and the third can be seen in Figure 11, item 6D. The stop in monitor was almost always produced with a full-fledged stop closure by the four speakers (except on one occasion by Speaker 2). All the other alveolar stops in pre-unstressed position were produced in various ways by the different speakers, and no speaker was completely consistent in producing one of these words in the same way each time during the experimental session.

Likewise the /n/ in <u>monitor</u> was produced in several ways: 1) a continuant with no evidence of nasal closure, 2) a nasal flap, 3) a non-nasal flap, i.e., with no velopharyngeal opening, and 4) a full-fledged nasal consonant. Examples of some of these manifestations of this pre-unstressed nasal consonant are shown in Figure 14. Most speakers produced more than one of these versions of /n/ in the course of the experimental session.



Figure 12. Degree of nasalization exhibited by Speaker 4 in vowels occurring in final syllables with a following nasal consonant. For each of 32 such vowels (in the words <u>heading</u>, <u>holding</u>, <u>pattern</u>, and <u>tacan</u>), the percent of the vowel duration that was nasalized was estimated from the spectrograms. The figure represents a distribution of these 32 measurements.



Figure 13. Spectrograms of two versions of the word <u>heading</u> produced by Speaker 1 during the experimental session. For the token at the right, the /d/ is flapped, whereas in the left-hand token there is little evidence for an alveolar closure. A more complete /d/ closure is illustrated in the right-hand spectrogram in Figure 11.



Figure 14. Spectrograms of three versions of the word <u>monitor</u>, illustrating three different acoustic manifestations of the pre-unstressed /n/. For the upper right token (produced by Speaker 4), there is no /n/ closure, and the only evidence for the nasal consonant is some nasalization during the vowel (slight weakening of F1). The other two tokens (produced vv Speaker 1) show examples of a complete /n/ closure (lower left) and of a flapped /n/ (upper left).

CONCLUSIONS

Two general observations can be made from the data presented: 1) some speakers will show consistent shifts in the properties of their speech when a task requires them to perform under stress; and 2) most speakers will exhibit more variability in the properties of their speech sounds when performing under stress than they will when performing under more relaxed and controlled conditions.

The first observation points to some consistent acoustic changes for which data should be extracted in future experiments:

- shifts in fundamental frequency in an upward direction, particulary for stressed vowels;
- changes in the waveform of the glottal pulses, leading to modifications in vowel spectra, particularly in the lowest one or two harmonics;
- shifts in the first two formants of vowels, primarily in the direction of increased F1 and F2 for the samples observed in this analyses; and
- 4) changes in the amplitude of turbulence noise in relation to the adjacent vowel for certain stop consonants; this change appears to be in the direction of increased relative amplitude of the noise for stressful situations.

The second observation appears to be more germane to the development of interactive voice systems, and that is the significant variability of words produced by a speaker in a stress-inducing situation. Apparently, because a number of production options are available for particular speech sounds in varying phonetic contexts, a speaker is less consistent in the choice of options in conditions of stress. For stressed syllables, options arise because certain acoustic characteristics are not required to make phonetic distinctions in English, and are free to vary. Similarly, "weakening" in the production of certain phonetic elements occurs because the production of the elements are made with less force or with a less precise realization of the ideal target states for the articulators.

We observed examples of such options in the production of <u>stress</u> or relatively strong syllables: 1) pre-voicing or lack of pre-voicing for word-medial voiced consonants; 2) the degree of nasalization that precedes a vowel; and 3) the amount by which a final vowel is lengthened.

For the task involved in this study, "weakening" or a decrease in the force or preciseness of articulation probably accounts for most of the variability observed. Significant variability in four types of weakening were observed: 1) variability in the amount of devoicing of an unstressed vowel between two voiceless consonants; 2) variability in the release of stop consonants; 3) variability in the degree to which a pre-unstressed velar stop becomes a continuant; and 4) variability in the degree to which pre-unstressed alveolar nasals or stops become flapped or become even weaker than flaps.

Three implications can be drawn from these results for the use of speech recognition systems in aircraft. First, and most obvious, is the importance of "training" recognition devices while the operator is not in a relaxed, stress-free environment. A common practice among users of recognizers is to train the device in the work space and during task performance. Updating the reference patterns at time intervals after initial training is also a common procedure. These practices increase the robustness of the device for changes in the acoustic structure of the speakers' utterances and changes in the ambient environment. Using similar procedures for aircraft environments will be difficult to realize since the variability of the environment and task requirements exceeds what is encountered in most common work spaces. Obtaining the number of training tokens of each utterance required to account for each speaker's variability in each different stressful situation would be unwieldly and tedious.

A second implication suggests training the users of interactive voice systems to speak in a uniform manner under all conditions, even conditions of stress. Users would be made aware of the phonetic modifications likely to occur under different conditions and would be trained to minimize them. Again, such a procedure will be hard to realize in aircraft because of the rapidly changing physical, cognitive, communicative, and psychological demands on aviators.

A third implication suggests that the most effective approach to the problem of changes in the speech signal due to stress is to design a speech recognition device capable of dealing with speech variations more directly. Since any recognizer designed for use in aircraft will use a finite vocabulary, grammar, and syntax, an integral part of such a device would be a set of rules specifying the possible acoustic modifications given the phonetic description of the word. A set of rules would indicate, for example, that prevoicing of an initial voiced stop consonant is an optional property providing evidence for the voicing feature, but is not a necessary acoustic correlate of this class of consonants in English. Within the context of the aviation environment, a large number of such rules would be incorporated into the recognition system to indicate alternative realizations of a word. Some rules would be specific, but many would be general in form, and would describe phonetic tendencies applicable to a large number of lexical items.

In summary, the kinds of changes and variability observed in the vocal utterances of speakers in a stress-producing situation serve to emphasize the limitations of speech recognition systems based on template matching. There is a need for a better understanding of these phonetic modifications and for developing ways of incorporating knowledge about these changes into speech recognition systems.

REFERENCES

1. Nye, J. M., Human factors analysis of speech recognition systems. <u>Speech</u> Technology, Vol. 1, No. 2, 50-57, 1982.

14 MAR 10 MAR

- 2. Mosko, J. D., Remington, R. W., and Griffin, G. R., Phonological variants in medial stop consonants under simulated operational environments: implications for voice activated controls in aircraft. In: Aural Communication in Aircraft. Advisory Group for Aerospace Research and Development (AGARD) Conference Proceedings 311. London: Technical Editing and Reproduction Ltd., Pp 3-1 - 3-4, 1981.
- 3. Mosko, J.D., and Griffin, G. R., Clear speech: A strategem for improving radio communications and automatic speech recognition in noise. In: Aural Communication in Aircraft. Advisory Group for Aerospace Research and Development (AGARD) Conference Proceedings 311. London: Technical Editing and Reproduction Ltd., Pp 4-1 - 4-4, 1981.
- Lentz, J. M., and Guedry, Jr., F. E., Motion sickness susceptibility: A retrospective comparison of laboratory tests. <u>Aviat. Space and Environ</u>. Med., 49:1281-1288, 1978.
- 5. Markel, J. D. and Gray, A. H., <u>Linear Prediction of Speech</u>. New York: Springer-Verlag, 1976.

REPORT DOCUMENT	ATION PAGE	READ INSTRUCTIONS
REPORT NUMBER	2. GOVT ACCESSION NO	. 3. RECIPIENT'S CATALOG NUMBER
NAMRL - 1300	An-A135 9	13.2
TITLE (and Subtitle)		5. TYPE OF REPORT & PERIOD COVERED
Interactive Voice Technolog	y: Variations in the	
Vocal Utterances of Speaker	s Performing a Stress	
Inducing Task	· .	6. PERFORMING ORG. REPORT NUMBER
AUTHOR(+)		8. CONTRACT OR GRANT NUMBER(*)
James D. Mosko, Kenneth N. S Glenn R. Griffin	tevens, and	
PERFORMING ORGANIZATION NAME AND	ADDRESS	10. PROGRAM ELEMENT, PROJECT, TASK
Naval Aerospace Medical Kese Naval Air Station	arch Laboratory	62758N MF58523 MF5852801B
Pensacola, Florida 32508		0001
CONTROLLING OFFICE NAME AND ADD	Fee	
Ways1 Modics1 Research and D	evelopment Command	16 August 1983
Naval Medical Command. Natio	nal Capital Region	13. NUMBER OF PAGES
Bethesda, Maryland 20814		36
4. MONITORING AGENCY NAME & ADDRESS	(it different from Controlling Office)	15. SECURITY CLASS. (of this report)
		UNCLASSIFIED
		154. DECLASSIFICATION DOWNGRADING
	·	
DISTRIBUTION STATEMENT (of this Repor	ri)	
Approved for public release	; distribution unlimi	ted
Approved for public release DISTRIBUTION STATEMENT (of the ebetre S. SUPPLEML.ITARY NOTES James D. Mosko, Ph.D. is cur Griffis Air Force Base, New D. KEY WORDS (Continue on reverse elde If no Supple Recognition	c; distribution unlimi ct entered in Block 20, if different for rrently at the Rome Ai York cesseary and identify by block numbe Motion stress	vom Report) nr Development Center,
Approved for public release - DISTRIBUTION STATEMENT (of the ebetre - SUPPLEM ITARY NOTES James D. Mosko, Ph.D. is cur Griffis Air Force Base, New - KEY WORDS (Continue on reverse elde if ne Speech Recognition Interactive voice technology	ct entered in Block 20, if different is rrently at the Rome Ai York Motion stress Vocal variabi	ted Tom Report) r Development Center, r) 1ity
Approved for public release DISTRIBUTION STATEMENT (of the abetra James D. Mosko, Ph.D. is cur Griffis Air Force Base, New KEY WORDS (Continue on reverse side if no Speech Recognition Interactive voice technology Speech communication	rrently at the Rome Ai York Motion stress Vocal variabi Acoustical an	ted om Report) r Development Center, '' lity alysis of speech
Approved for public release Approved for public release DISTRIBUTION STATEMENT (of the ebetre Supplementary notes James D. Mosko, Ph.D. is cur Griffis Air Force Base, New KEY WORDS (Continue on reverse elde if ne Speech Recognition Interactive voice technology Speech communication Stress effects on speech	c; distribution unlimi ct entered in Block 20, if different is crently at the Rome Ai York Cesseary and identify by block numbe Motion stress Vocal variabi Acoustical an	ted om Report) r Development Center, '' lity alysis of speech
Approved for public release Approved for public release DISTRIBUTION STATEMENT (of the ebetre James D. Mosko, Ph.D. is cur Griffis Air Force Base, New KEY WORDS (Continue on reverse elde If ne Speech Recognition Interactive voice technology Speech communication Stress effects on speech ABSTRACT (Continue on reverse elde If neo As plans develop to incor ilitary aircraft, it is impo- ikely to be modified in stre evices must be designed to r ariety of conditions. Acoustical analyses were	ct entered in Block 20, 11 different is ct entered in Block 20, 11 different is creently at the Rome Ai York Motion stress Vocal variabi Acoustical an Exercise and identify by block number porate interactive vo ortant to determine ho essful operational sit cespond correctly to c conducted of words pr	<pre>ted om Report) r Development Center, // lity alysis of speech // ice systems in advanced w the speech of a pilot is uations. Speech recognition ommands produced under a oduced by four speakers in</pre>
Approved for public release Approved for public release DISTRIBUTION STATEMENT (of the ebetre Supplementation statement (of the ebetre James D. Mosko, Ph.D. is cur Griffis Air Force Base, New KEY WORDS (Continue on reverse elde If ne Speech Recognition Interactive voice technology Speech communication Stress effects on speech ABSTRACT (Continue on reverse elde If neo As plans develop to incor ilitary aircraft, it is impo- ikely to be modified in stre evices must be designed to r ariety of conditions. Acoustical analyses were	e; distribution unlimi ct entered in Block 20, 11 different for crently at the Rome Air York Motion stress Vocal variabi Acoustical an essent to determine hor essful operational sit conducted of words pr	<pre>ted om Report) r Development Center, r/ lity alysis of speech) ice systems in advanced w the speech of a pilot is uations. Speech recognition ommands produced under a oduced by four speakers in</pre>
Approved for public release Approved for public release Supplementation statement (of the ebetre Supplementation statement (of the ebetre James D. Mosko, Ph.D. is cur Griffis Air Force Base, New KEY WORDS (Continue on reverse elde if ne Speech Recognition Interactive voice technology Speech communication Stress effects on speech ABSTRACT (Continue on reverse elde if nec As plans develop to incor ilitary aircraft, it is impo- ikely to be modified in stre evices must be designed to r ariety of conditions. Acoustical analyses were D I JAN 73 1473 EDITICN OF I NOV 68	ct entered in Block 20, 11 different is rrently at the Rome Ai York ceeseary and identify by block number Motion stress Vocal variabi Acoustical an resonate interactive vo ortant to determine ho essful operational sit cespond correctly to c conducted of words pr	<pre>ted om Report) r Development Center, r) lity alysis of speech) ice systems in advanced w the speech of a pilot is uations. Speech recognition ommands produced under a oduced by four speakers in LASSIFIED</pre>
Approved for public release Approved for public release DISTRIBUTION STATEMENT (of the ebetre SupplemITARY NOTES James D. Mosko, Ph.D. is cur Griffis Air Force Base, New KEY WORDS (Continue on reverse elde II ne Speech Recognition Interactive voice technology Speech communication Stress effects on speech ABSTRACT (Continue on reverse elde II nec As plans develop to incor ilitary aircraft, it is impo- ikely to be modified in stre evices must be designed to r ariety of conditions. Acoustical analyses were D FORM 1473 EDITICN OF 1 NOV 68 S/N 0102-014-6601	ct entered in Block 20, 11 different in ct entered in Block 20, 11 different in crently at the Rome Ai York Conserv and identify by block number Motion stress Vocal variabi Acoustical an conducted interactive vo conducted of words pr conducted of words pr SECURITY CL	<pre>ted om Report) r Development Center, // lity alysis of speech // ice systems in advanced w the speech of a pilot is uations. Speech recognition ommands produced under a oduced by four speakers in LASSIFIED AssiFICATION OF THIS PAGE (When Data Entered </pre>

and a first section of the section of the

a contrate and the contrate of the second of the Kill

0,5

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Dete Entered)

A motion stress-inducing situation. The aim of the analyses was to document the kinds of changes that occur in the vocal utterances of speakers who are exposed to motion stress and to comment on the implications of these results for the design and development of voice interactive systems.

The speakers differed markedly in the types and magnitudes of the changes that occurred in their speech. For some speakers, the stress-inducing experimental condition caused an increase in fundamental frequency, changes in the pattern of vocal fold vibration, shifts in vowel production and changes in the relative amplitudes of sounds containing turbulence noise. All speakers showed greater variability in the experimental condition than in a more relaxed control situation. The variability was manifested in the acoustical characteristics of individual phonetic elements, particularly in speech sounds occurring in the vicinity of unstressed syllables. The kinds of changes and variability observed serve to emphasize the limitations of speech recognition systems based on template matching of patterns that are stored in the system during a training phase. There is need for a better understanding of these phonetic modifications and for developing ways of incorporating knowledge about these changes within a speech recognition system.

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE(When Date Entered)