

VOICE QUALITY BEFORE AND AFTER RADIOTHERAPY: ACOUSTICAL, CLINICAL AND PERCEPTUAL PITCH MEASURES*

Irma M. Verdonck-de Leeuw and Florian J. Koopmans-van Beinum

Abstract

Speech samples (read aloud text and sustained /a/) of patients with early glottic carcinoma before and after radiotherapy, and of control speakers, were analysed for various pitch measures: acoustical pitch, EGG-pitch, perceptual pitch evaluations by trained and untrained raters, and by the speakers themselves and their partners. Results showed that the speakers before radiotherapy differed from the speakers 6 months after radiotherapy, according to the acoustical measured pitch. There were no differences between the speaker groups on any of the other pitch measures. The results of pitch ratings by trained listeners, EGG, and acoustical pitch correlated strongly; a factor analysis resulted in five factors: one factor for read aloud text (acoustical pitch, EGG, and pitch evaluations by trained raters), another factor for read-aloud text with the pitch evaluations by untrained raters, one factor for the sustained /a/ speech material (acoustical pitch, pitch evaluations by trained and by untrained raters), a factor for the pitch evaluations by the speakers themselves, and a factor for the pitch evaluations by their partners; relations between acoustical measured pitch and pitch ratings by trained and untrained raters seem to be dependent on the voice quality of the speakers, in terms of patients versus control speakers.

1 Introduction

Within the scope of a co-operative study between the Netherlands Cancer Institute (Antoni van Leeuwenhoek Hospital) and the Institute of Phonetic Sciences of the University of Amsterdam, research was carried out on the effect of radiotherapy on voice quality. The aim of this study was to obtain parameters that can describe voice quality of patients with early glottic cancer (before and after radiotherapy) and of normal speakers. Voice quality can be described by several perceptual, clinical, as well as acoustical methods. In this paper we focused on various pitch measures. Pitch is supposed to be one of the parameters that can be influenced by the presence of a tumour or by side-effects of radiotherapy on the vocal fold tissue, such as mucositis,

* Preliminary results were presented at the XIIIth International Congress of Phonetic Sciences in Stockholm (Verdonck-de Leeuw & Koopmans-van Beinum, 1995)

tissue oedema, etc.. Before radiotherapy, the actual tumour can cause changes of the vocal folds such as mass change, stiffness change, and asymmetry. Little is known about voice quality after radiotherapy. Some studies report voice improvement to a normal or near-normal level, 6-12 months after radiotherapy, for about 70% of the patients (Colton et al., 1978; Harrison et al., 1990; Hoyt et al., 1992; Miller et al., 1990). Other studies report abnormal post-radiation voices (Lehman et al., 1988; Stoicheff, 1983).

Furthermore, pitch measurements are important cues for other acoustical and perceptual measurements, such as spectral noise (Emanuel et al., 1974; Wolfe et al., 1991; Verdonck-de Leeuw & Boersma, 1995), and breathiness and tension (Hammarberg et al., 1980). In a later stage of our study the results will be compared with other perceptual parameters of voice quality (evaluations on semantical scales, such as breathiness, harshness, creakiness), with other clinical methods (phonetogram, phonation-quotient, evaluations of stroboscopic recordings of vocal fold vibration, quality of life questionnaires), and with other acoustical analyses (LTAS, SNR, perturbation).

The specific aims of the present paper were to investigate the pitch of voices before and after radiotherapy, and of control voices, and to investigate the relations between the various perceptual, acoustical and clinical pitch measures.

Acoustical and EGG pitch analyses were taken into account as 'objective' pitch measures; the results were compared with perceptual pitch evaluations by trained and untrained raters, and with perceptual pitch evaluations by the speakers themselves and by their partners. The trained raters were used to provide an analytic description of voice quality. The role of the untrained raters was to find out how 'ordinary' people evaluate voice quality. In order to investigate if changes of voice quality also influence quality of life evaluations, the evaluations of the speakers themselves and their partners were taken into account. The untrained and trained raters were asked to evaluate voice quality of both read-aloud text and sustained /a/ produced by the speakers. Analyses of a sustained /a/ are common practice in clinical settings (for instance voice range profile, phonation flow) and are therefore included as speech material in our project. In order to assess the practical relevance of voice changes in the patients' home environment, running speech fragments were used as well, because these are more representative for conversational speech. Fragments of read-aloud text were used, rather than spontaneous speech, in order to avoid variance between speakers caused by unequal texts.

2 Method

2.1 Speakers/recordings

Patients with early glottic cancer (T1N0M0) were treated with radiotherapy (60 Gy in 30 fractions, or 66 Gy in 33 fractions). Voice samples of the same 10 patients were recorded before radiation, as well as 6 months, and 2 years after radiation (longitudinal group, see table 1). Recordings were also made of 3 other groups of 10 patients each (separate groups, see table 1), before radiation, 6 months, and 2 years after radiation, and of 20 patients longer than 3.5 years after radiation. Finally, recordings were made of 20 control speakers (table 1). The matching between patients and control speakers took into account sex (all male), age (47-81), as well as smoking habits. The speakers read out a text for about 5 minutes and produced a sustained /a/. Fragments (ca 45 s.) of all texts and the onset + 2 seconds of the sustained /a/ speech material were digitised by means of the Sound editor of an Iris Indigo R4000 with a sample frequency of 48 kHz, and with 16 bit resolution.

Table 1. Composition of the subject sample: longitudinal group before, 6 months after, and 2 years after radiation; separate groups before, 6 months after, 2 years after, and longer than 3.5 years after radiation, and control group; totals and mean ages are given in the last two rows.

	before	6 months	2 years	>3.5 years	control
Longitudinal group	10 -->	10 -->	10		
Separate groups	10	10	10	20	20
Totals	20	20	20	20	20
Mean age (years)	64	64	66	70	65

2.2 Raters/rating procedure

An adapted version of the Vocal Profile by Laver (1981) was used for the evaluations by the trained raters as their role was to provide an analytic description of voice quality. The untrained raters, the speakers themselves and their partners evaluated the voice quality of the speakers by means of the scaling instrument developed by Fagel et al. (1983); this instrument was developed to obtain ratings of untrained Dutch raters. All raters evaluated the read-aloud text in a first session; the sustained /a/ speech material was judged one week later. First the raters heard 10 examples of various voices, ranging from extreme pathological to normal, in order to get a reference frame. After the examples 110 fragments of read aloud text were presented (10 training fragments and 100 fragments of the speakers as indicated in table 1).

The trained raters were three female phonetic researchers; two had followed a training course on the Vocal Profile by Laver (1981), the third was trained by one of the others. All three rated the voices independently from each other on various voice quality scales adapted from the Vocal Profile. On average, the rating of each scale took about 30 seconds. Results of the scales *low-high* and *sonorous-shrill* (13-point) are presented in the present paper.

The untrained raters in this experiment were 20 students (6 male, 14 female), without any experience for this listening task. They were paid for their participation. The raters received written instructions. The raters listened to the tapes in a quiet room, individually. On the average, the whole rating procedure (instructions + rating) took about 1 1/2 hours for the read-aloud text and 1 1/2 hours for the sustained /a/. The raters judged the speech samples on various voice quality scales; in the present paper, only results of the scales *low-high* and *deep-shrill* (7-point) are presented.

The speakers and their partners received score forms with a written instruction. They were asked to evaluate the voice of the speaker at home by filling out the form independently from each other on various voice quality scales; again, only the scales *low-high* and *deep-shrill* (7-point) are presented in this paper. The speakers and the partners were treated as separate rater groups. Notice the differences between at the one hand these two rater groups (judged one voice every time) and at the other hand the trained and untrained raters (judged all voices at one time).

2.3 Acoustical Pitch

The acoustical pitch was determined by means of the program "Praat" developed by Boersma (1995). The pitch period of a sound was determined by the position of the maximum of the autocorrelation function of the sound. The complete 9-parameter algorithm, as is implemented into the speech processing program Praat, is extensively described by Boersma (1993). The VoicingThreshold and the Silence-Threshold were

adapted to exclude speech pauses, while all other parameters were kept default, so that the Hamming/Hanning-equivalent window length was 60 ms. For the sustained /a/, the most representative pitch value was calculated as the median of all measured frames. For the read text fragments, the average value over all voiced frames was taken instead, in order to be able to compare the results with the next to be described clinical fundamental frequency.

2.4 Clinical fundamental frequency

By means of an electroglottograph (Stopler Teltec GFA06) the average fundamental frequency was measured for the read aloud text. The same text was used as was recorded for perceptual and acoustical analyses described above, but the recording itself was another realization. The speakers read aloud the text for about 5 minutes while up to 1000 voiced samples were analysed and averaged. We experienced difficulties in obtaining EGG data for some of the speakers; these difficulties may be due to incomplete vocal fold closure, resulting in weak and noisy signals or disruptions of the signal; also fat necks seem to raise problems. Cases we judged as unreliable were left out of consideration, resulting in 18 speakers before radiotherapy, 11 speakers 6 months after, 16 speakers 2 years after, and 14 speakers longer than 3.5 years after radiotherapy, and 20 control speakers. Each group originally contained 20 speakers.

3 Results

3.1 Intrarater reliability

Since during the rating procedure ten voice samples, selected from the available material itself and ranging from extremely deviant to normal, were presented twice: first as training samples (the first 10 voice samples that had to be judged), the second time as part of the 100 test samples, we could determine intrarater reliability. A matched sample t-test compared the first and second ratings of each voice over all 3 trained and 20 untrained raters, and indicated that the trained and untrained raters were highly reliable: none of the t-tests was significant ($p > 0.05$); they never differed more than 1/2 scale value. This counted for the read-aloud text as well as for the sustained /a/. This test could not be done for the judgments by the speakers and their partners because no repeated judgments were available here.

3.2 Interrater reliability

An interrater reliability coefficient was calculated for all scales: Cronbach's alpha. This is a measure of the reliability of the means of the ratings given by a panel of raters. Alpha is defined as $(MS_{betw} - MS_{res}) / MS_{betw}$ in which MS_{res} = Mean Square residual and MS_{betw} = Mean Square between people. A low reliability can be caused either by a high MS_{res} (the raters disagree), a low MS_{betw} (there is little variation between the speakers, i.e. the true variance is low), or by both (Rietveld, 1993).

All ratings were reliable ($\alpha > .80$); on the read aloud text, the results for the three trained raters were $\alpha = .83$ for *low pitched-high pitched* and $\alpha = .82$ for *sonorous-shrill*; for the 20 untrained raters $\alpha = .91$ for *low-high* and $\alpha = .89$ for *deep-shrill*. On the sustained /a/ speech material, the results for the trained raters

were $\alpha = .90$ for *low pitched-high pitched* and $\alpha = .87$ for *sonorous-shrill*, for the 20 untrained raters $\alpha = .93$ for *low-high* and $\alpha = .90$ for *deep-shrill*. The differences between trained and untrained raters were due to low MSbetw (=true variance) by the trained raters. This test could not be done for the ratings of the speakers themselves and their partners because they rated only one voice at the time.

3.3 Differentiation among speaker groups

Two variance analyses for the acoustical, perceptual and clinical pitch parameters were carried out: trend analyses on the three longitudinal speaker groups (randomized block design, with repeated measures) and variance analyses on the five separate speaker groups (one-way factorial design, without repeated measures). No significant differences ($p < 0.05$) were found between the separate speaker groups (before radiotherapy, 6 months after, 2 years after, longer than 3.5 years after radiotherapy, and the control speakers), nor for the longitudinal speaker group (before radiotherapy, 6 months after, and 2 years after radiotherapy) either. The results for the acoustical pitch on read-aloud text revealed differences between the speaker groups at a lower level ($p < 0.10$). Posthoc tests after the variance analysis ($F = 2.34$, $p = 0.06$) for the separate speaker groups revealed that the difference between speakers before radiotherapy and the speakers after radiotherapy was significant ($p < 0.05$). The same results were found for the longitudinal group ($F = 3.18$, $p = 0.07$).

Obviously, pitch is not influenced by side-effects of radiotherapy. Speakers before radiotherapy seem to have high-pitched voices according to acoustical measured pitch on read-aloud text, but this conclusion was not perceptually confirmed.

3.4 Relations between the various pitch measures

In order to investigate interrelations between acoustical, perceptual and clinical pitch measures, Pearson correlations were calculated for the acoustical and perceptual pitch measures on read-aloud text and sustained /a/ and EGG (table 2). A Principal Component Analysis was used to decompose the correlation matrix into (varimax rotated) factors (PCA). With the criterion 'eigenvalue greater than one', the PCA produced 5 factors, together explaining 78 % of the total variance (table 3). On the basis of the factor loadings ($> .50$) the factors were mainly determined by:

- acoustical pitch, EGG, and *low pitched-high pitched* and *sonorous-shrill* by trained raters, on read-aloud text (factor 1),
- acoustical pitch, *low pitched-high pitched* and *sonorous-shrill* by trained raters, and *low-high* and *deep-shrill* by untrained raters, on sustained /a/ (factor 2),
- *low-high* and *deep-shrill* by untrained raters, on read-aloud text (factor 3),
- *low-high* and *deep-shrill* by the speakers themselves (factor 4),
- *low-high* and *deep-shrill* by the partners (factor 5).

The relations between the acoustical pitch analysis, the EGG data, and the pitch ratings by the trained raters on the read-aloud text are clear; they loaded highly on the same factor and correlations between these parameters were high ($.60 < r < .75$), as can be seen in table 2 and table 3, respectively. For the sustained /a/ speech material, the same can be concluded and it can be added that the untrained raters were also able to evaluate pitch according to acoustically measured pitch, although correlations were somewhat lower ($.40 < r < .56$) compared to the read-aloud text.

ficant over all speakers. We decided to compare the results of Pearson correlations between the acoustical and perceptual pitch measures of speakers before radiotherapy with the results of the control speakers, since previous research revealed that speakers before radiotherapy have the most deviant voices (i.e. whisper, harshness etc.) compared to the 'normal' control speakers (De Leeuw, 1991). Results are given in table 4 and reveal that the correlations between perceptual and acoustical measured pitch are always lower for the speakers before radiotherapy than for the control speakers, indicating that pitch ratings are influenced by the voice quality of speakers.

Table 4. Significant ($p < 0.01$) Pearson correlations over 20 patients before radiotherapy and over 20 control speakers for acoustical pitch and perceptual pitch ratings by 3 trained and 20 untrained raters on read-aloud text, and for acoustical pitch and perceptual pitch ratings by 3 trained and 20 untrained raters on a sustained /a/.

	read-aloud text acoustical pitch		sustained /a/ acoustical pitch	
	patients	control	patients	control
<i>low-high</i> untrained	.61	.79	.71	.82
<i>deep-shrill</i> untrained	.65	.83	.60	.82
<i>low pitched-high pitched</i> trained	.69	.69	.59	.83
<i>sonorous-shrill</i> trained	.67	.81	.53	.83

4. Conclusion

It can be concluded that there was no statistically significant effect found on the various pitch measures by radiotherapy. The tendency for acoustical measured pitch, that patients before radiation had high pitched voices compared with patients six months after radiotherapy, may be due to mechanical effects of the tumour on the vocal folds. Another explanation may be an increased tension of the vocal folds by the patient in order to compensate for his voice loss. Also, little is known about the effect of microlaryngeal surgery that most of the patients have undergone before radiation.

Although the other results in this experiment did not differentiate significantly between the speaker groups, strong relations were found between the acoustical pitch analysis, and the EGG data, and the pitch evaluations by the trained raters on the read-aloud text. The expectation that what one can hear should also be measurable, becomes true in this experiment, at least for the trained raters. The same counts for acoustical pitch, pitch evaluations by the trained raters and by the untrained raters on the sustained /a/ speech material. The speakers themselves and their partners judged their voices differently from the trained and untrained raters; two separate factors were formed due to nonsignificant correlations with the other pitch measures. This may be due to the different way they were asked to evaluate the voices: one voice at one time instead of all voices at one time from tape recordings, as the raters did.

In the present study an attempt was made to describe the relations between acoustical and perceptual pitch measures on the basis of the classification of speakers in 'speakers before radiation' and 'control speakers'. Relations between acoustical and perceptual pitch measures seem to depend on the voice quality of the speakers: relations appeared to be lower for speakers before radiotherapy than for control speakers and earlier research showed that speakers before radiation have deviant voices compared to control speakers. An explanation for this conclusion might be that extreme

whisper or harshness for instance, colours the pitch evaluations. Further research is needed on relations between acoustical measured pitch and specific voice quality parameters like whisper, harshness and so on.

The untrained and trained raters were asked to evaluate voice quality of read-aloud text and of a sustained /a/ produced by the speakers. Earlier research has shown that read-aloud text and sustained /a/ are equally reliable with respect to perceptual evaluations of (laryngeal) voice quality (De Krom, 1995; De Leeuw, 1991). Results in the present study revealed that read-aloud text and sustained /a/ are two different types of speech material at least in as far as pitch analyses are concerned: separate factors were found for read-aloud text and sustained /a/. Therefore, both types should be taken into account in voice quality analyses. More voice quality research is necessary to investigate the voice production process in running speech and sustained /a/.

5. Acknowledgments

We are indebted to the subjects who provided the voice material for this study, and judged their own voice, just as their partners. We would like to thank J.M.A. De Jong of the 'Radiotherapeutisch Instituut Limburg' for referring part of the patients, A. Greven of the Phoniatic Department of the Academic Hospital of the Free University of Amsterdam, and G. Vreeburg of the Department of Otolaryngology and J. Uytendijk-Winkel of the Department of Voice, Speech and Language Disorders of the Academic Hospital of Maastricht for their co-operation during recording sessions. Furthermore, we would like to thank the untrained raters, and the trained raters J. van Rie and R. van Bezooijen of the University of Nijmegen, who gave the perceptual pitch descriptions. The software for processing the acoustical pitch was kindly provided by P. Boersma of the Institute of Phonetic Sciences of the University of Amsterdam. G. Baris, H. Bartelink and R. Keus of the Department of Radiotherapy of the Netherlands Cancer Institute are gratefully acknowledged for referring the patients and for making useful comments on earlier versions of this paper, as did F.J.M. Hilgers (Department of Otolaryngology-Head&Neck Surgery of The Netherlands Cancer Institute), and L.C.W. Pols (Institute of Phonetic Sciences, University of Amsterdam).

6. References

- Boersma, P. (1993). "Accurate short-term analysis of the fundamental frequency and the harmonics-to-noise ratio of a sampled sound", *Proceedings of the Institute of Phonetic Sciences Amsterdam* 17, 97-110.
- Colton R. H., Sagerman, R., Chung, C.T., Yu, Y.W. & Reed, G.F. (1978). "Voice change after radiotherapy", *Radiology* 127, 821-824.
- De Leeuw, I.M. (1990). "The relation between perceptual and clinical parameters of voice quality of patients with early glottic cancer before and after radiotherapy and of normal speakers", *Proceedings of the Institute of Phonetic Sciences Amsterdam* 14, 27-38.
- Emanuel, F.W. & Smith, W.F. (1974). "Pitch effects on vowel roughness and spectral noise", *Journal of Phonetics* 2, 247-253.
- Fagel, W., Van Herpt, L. & Boves, L. (1983). "Analysis of the perceptual qualities of Dutch speakers' voice and pronunciation", *Speech Communication* 2, 315-326.
- Hammarberg, B., Fritzell, B., Gauffin, J., Sundberg, J. & Wendin, L. (1980). "Perceptual and acoustic correlates of abnormal voice qualities", *Acta Otolaryngologica* 90, 442-451.
- Harrison, L.B., Solomon, B., Miller, S., Fass, D.E., Armstrong, J. & Sessions, R. (1990). "Prospective computer-assisted voice analysis for patients with early stage glottic cancer: a preliminary report of the functional result of laryngeal irradiation", *Int. J. Radiation Oncology Biol. Phys.* 19, 123-127.
- Hoyt, D.J., Lettinga, J.W., Leopold, K.A. & Fisher, A. (1992). "The effect of head and neck radiation therapy on voice quality", *Laryngoscope* 102, 477-480.

- Laver, J., Wirz, S., Mackenzie, J. & Hiller, S.M. (1981). "A perceptual protocol for the analysis of vocal profiles", *Edinburgh University Department of Linguistics Work in Progress* **14**, 139-155.
- Lehman, J.J., Bless, D.M. & Brandenburg, J.H. (1988). "An objective assessment of voice production after radiation therapy for stage 1 squamous cell carcinoma of the glottis", *Otolaryngol. Head Neck Surg.* **98**, 121-129.
- Miller, S., Harrison, L.B., Solomon, B. & Sessions, R. (1990). "Vocal changes in patients undergoing radiation therapy for glottic carcinoma", *Laryngoscope* **100**, 603-606.
- Verdonck-De Leeuw, I.M. & Boersma, P. (1995). "The effect of radiotherapy measured by means of Harmonicity" (to appear in: Powell, T.W. (Ed.), *Pathologies of speech and language: contributions of clinical phonetics and linguistics.*)
- Verdonck-De Leeuw, I.M. & Koopmans-van Beinum F.J. (1995). "The effect of radiotherapy on various acoustical, clinical and perceptual pitch measures", *Proceedings of the International Congress of Phonetic Sciences 95, Stockholm* **4**, 610-613.
- Wolfe, V., Cornell, R. & Palmer, C. (1991). "Acoustic correlates of pathologic voice types", *Journal of Speech and Hearing Research* **34**, 509-516.