

Braga, D.; Coelho, L.; Resende Jr., F.G.V., Dias, M. S. 2007. "Subjective and Objective Evaluation of Brazilian Portuguese TTS Voice Font Quality", in *Advances in Speech Technology*, 14th International Workshop, June 27-29 2007, Maribor, Slovenia.

Subjective and Objective Evaluation of Brazilian Portuguese TTS Voice Font Quality

Daniela Braga¹, Luís Coelho², Fernando Gil V. Resende³ and Miguel Sales Dias¹

¹ Microsoft Language Development Center, MLDC, Portugal
{i-dbraga, midias}@microsoft.com

² Polytechnic Institute of Oporto, Portugal
luiscoelho@eseig.ipp.pt

³ Federal University of Rio de Janeiro, Brazil
gil@lps.ufrj.br

Abstract. In this paper, it is our aim to define a set of objective acoustic criteria, based on subjective listeners assessment of talent voices, that can help to rate the voice font quality, bearing in mind the development of Text-to-Speech systems (hereafter TTS). The voice talent selection process is seldom discussed in scientific papers, since academic community is obviously more concerned with synthetic voice evaluation than with modal voice talent assessment. However, the choice of a voice that will be used in TTS systems is definitely of extreme importance. In this work, we describe the selection process of the Brazilian Portuguese voice talent, based on socio and linguistic criteria and supported by subjective and objective tests and measurements of modal voice quality. Relationship between subjective and objective parameters is presented and discussed. A subjective test comparing natural and resynthesized voices was also conducted and its results are discussed.

1 Introduction

In the last few years, TTS technology has been widely improved in several aspects. The speech is more fluent and continuous, the intelligibility is higher, even emotions can be now synthesized. These developments lead to a real naturalness and created the new challenge of starting embedding TTS systems with working objects and tools, such as computers, cellular phones, Personal Digital Assistants (PDAs) and other devices, in order to create a voice help that is always present in daily human tasks. This permanent contact with a voice should be reliable and pleasant. But how good is a voice? Which acoustic and perceptive parameters make people prefer one voice to another? Is it possible to measure voice pleasantness? Questions like these together with the development of a new voice font for Brazilian Portuguese (BP) TTS systems are in the basis of the present study. In specialized literature, voice quality is often associated with voice impairments or disorders. There is an extensive bibliography on this subject, in which voice quality means normal speech characteristics that are

somehow affected for pathological reasons, leading to hoarseness, roughness, raspiness, effort to talk, breathiness, vocal fry, uncomfortable or abnormal pitch, and other abnormal vocal symptoms [1], [2], [3], [4]. Still related with the clinical field and from a physiological point of view, voice quality is also used in phonetic and acoustic descriptions, such as the one carried out by Laver [5]. In this paper, the author describes the dynamics of six basic phonation types (modal voice – the so-called “normal voice”, falsetto, whisper, creak, harshness and breathiness) and possibilities of combination of each one’s settings, generating compound phonation types (such as whispery falsetto, whispery creak, harsh whispery voice and so on). This author encompasses both supralaryngeal and laryngeal phonatory settings. This labelling proposal has been widely applied to speech analysis and synthesis. Another context in which voice quality is often used is the evaluation of telephone speech output [6] and synthesized speech. Many papers address several overall subjective tests, like the Mean Opinion Score (MOS) test [7], or intelligibility tests such as Semantically Unpredictable Sentences (SUS) test, the Diagnostic Rhyme Test (DRT), amongst many other tests (for a review on speech quality tests in TTS, see [8]). Voice quality is also associated with prosody generation and synthesis of emotions. Some authors explore the role of voice quality (meaning in this context fundamental frequency (hereafter F0), intensity, tempo aspects and spectral parameters) in emotion, mood and attitude expression [9], [10]. Campbell et al., within the same framework, call it “the 4th prosodic dimension” [11]. However, the assessment of modal voice quality (in Laver’s expression, explained above [5]) is rarely discussed. An interesting study addressing this subject was conducted by Syrdal et al. [12] in order to check the suitability of a speaker’s voice to develop a TTS system, based on the assumption that the perceived quality of a natural voice does not necessarily mean synthesized voice quality. The authors explore the correlation between acoustic characteristics (RMS energy, breathiness, long-term spectra, f0, formants and bandwidths, speaking rate, concatenation and target costs) and the subjective attributes of synthetic speech quality (intelligibility, naturalness and pleasantness). Nevertheless, there is not much description on the speakers’ selection process, which is said to have been made “empirically”, although all the candidates (6 females and 9 males) were said to be professional speakers. The lack of studies in the area of subjective and objective criteria for voice talent selection opens the possibility for the research presented in this paper. This paper is structured as following: in section 2, the voice talent selection process that led to the identification of the TTS voice talent is described; two conducted subjective tests and corresponding results are presented and discussed in section 3; in section 4, objective measurements are shown; in section 5, tests and results comparing natural and resynthesized voice samples are discussed; in section 6, main conclusions are pointed out and future work is foreseen.

2 Voice Talent Selection

The voice talent selection process was developed in four stages. In the first stage, a national call for voice talents was launched. A few mandatory profile requirements

were defined, such as: being female, having BP as mother language, having studied in Brazil up to the university level, speaking standard Brazilian Portuguese (although other dialectal varieties were expected to be considered) and having preferably some radio or theatre vocal experience. From a total of 97 presented female candidates, 62 were admitted to the second stage of evaluation. At this stage, candidates were invited to send us samples of their voices, with no restrictions of duration, with the maximum quality they could produce (acquired by their own means, either in professional recording studios, from their portfolio, or from their own PC). A subjective test was conceived with 13 questions, based in the MOS scale (with 5 points ratings scale), and was conducted by 7 adult listeners (2 females and 5 males), familiarized with Speech Processing technology. The test was blind, that is to say each candidate and corresponding voice sample were given a number, in order to avoid any partial judgment and the test was individually conducted. The listeners were asked to hear the voice samples (with headphones) and to rate them according to a set of previously defined subjective attributes (section 3). After this test, the 10 best scored candidates were invited to record a small text selected by us (of 219 words, phonetically and prosodically rich, with emotion expressions) in a professional recording studio (third stage of evaluation). Our goal with this casting was to produce a final survey in which all the voices would be evaluated under the same conditions of sampling rate and text type. We produced a subjective test survey with 13 questions, assessing essentially the same attributes as in the previous subjective test, but which were differently rated. In other words, the first test had a 5 points rating scale, whereas the second test was an exclusive multiple choice questionnaire where only the best voice for each attribute could be selected. The survey was carried out by 124 listeners (35 females and 89 males) that were not familiarized with Speech Processing technology (third stage of evaluation). The voice ranking was obtained through the sum of votes each voice received along the survey. Then in the fourth stage of evaluation, an objective analysis was carried out to confirm evidence provided by the subjective tests (section 4).

3 Subjective Tests and Results

After stage 1, where the voice's selection was based on curriculum analysis, we proceeded with an acoustical based selection. In the 2nd stage of evaluation, Test A was conducted. 62 voice samples of up to 2 minutes recording time were evaluated according with the following subjective parameters: pleasantness (PLS), intelligibility (INT), articulation (ART), accent pronunciation (ACP), expressiveness (EXP), exceptionalness (EXC), sensuality (SNS), attitude (ATT). Three more questions were asked addressing the listeners' judgment on the suitability of those voices for the following applications: e-mail, news or instructions reading. This was a 5 points rating scale, which means that all voices were classified with marks from 1 (bad) to 5 (excellent) in every subjective attribute. From the 62 voices the 10 best overall scores have been selected, sorted alphabetically by owner's name and numbered 1 to 10. In Test A results show that: all the voices have been rated above 2.5 which is the middle of the considered scale; the most pleasant voices are 6 and 7 (both rated 4.0); the sexiest

voice was 8 (rated 4.0); voice 3 showed more attitude, exceptionalness and intelligibility (rated 5, 4.3 and 4.7) but oppositely it was the less sexiest voice; almost all the voices had very good articulation (rated 5.0 with exceptions to voices 1, 6 and 9); voices 3 and 7 seemed to be more special and hard to forget (rated 4.3); regarding accent pronunciation, most of the voices showed very small dialectal marks but few have achieved the best score; voices 1 and 3 were judged as the most expressive (rated 4.4). Concerning voice applications, voices 8 and 10 received the highest preferences for e-mail (both rated 3.9 which is not a very high score) and news reading (rated 4 and 4.3); for instructions voices 10 and 3 were preferred. In these last parameters, related with frequent interactions, voice 10 had always the best score. In Figure 1, the assessed subjective parameters are displayed.

For age assessment, listeners were asked to guess voices' age in the following intervals: 1 – under 25 years old), 2 – between 25 and 30, 3 – between 30 and 35, 4 between 35 and 40, 5 – over 40. All 10 finalist voices were perceived to be between 2.5 and 3.4 points, which means that the preferred perceptual voice age rounded 27 and 31.5 years old. This evidence is interesting when analyzing age perception of the less scored candidates, whose age was perceived as being over 35 years old. Concerning speaking rate, all 10 best scored voices were considered normal, neither too fast nor too slow. Another interesting issue is that perceived age does not correspond to real age, since most of the 10 final speakers are over 31. Although we have used most of the MOS scale parameters, we have excluded naturalness, comprehension and listening effort. In fact, these parameters are crucial in synthetic voice quality, but are redundant when assessing normal and professional modal voices that are recorded in good quality conditions.

In the 3rd evaluation stage, Test B was carried out with the purpose of deciding the best Brazilian voice that will be integrated in our BP TTS system. Therefore, the test could not be too much time consuming, in order to be able to be carried out by as many listeners as possible. That is why we excluded some subjective parameters that were assessed in Test A, such as EXC and ACP. INT and ART were assessed in the same question. The test was designed as a survey. First questions were on listener's gender and mother language. Middle questions were on the PLS, SNS, ATT, INT, SPR and voice applications (e-mail, news and instructions reading). Last question was on age perception of the preferred voice. Test B was rated according with the listeners' best choice for each question.

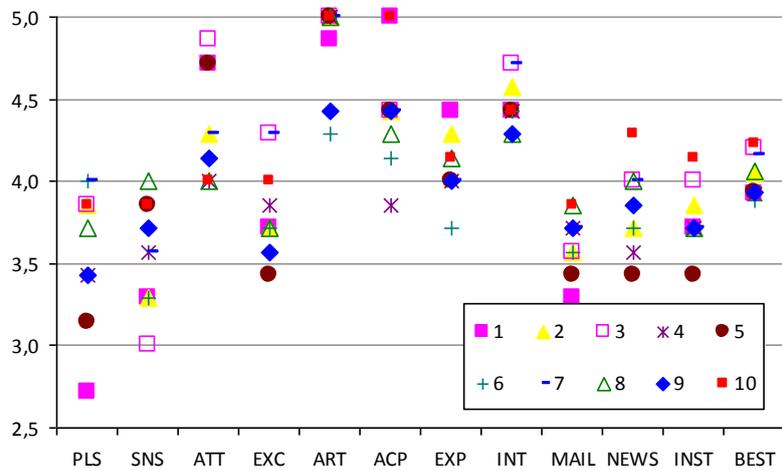


Fig. 1. Subjective test results for the 10 best scored voices on test A.

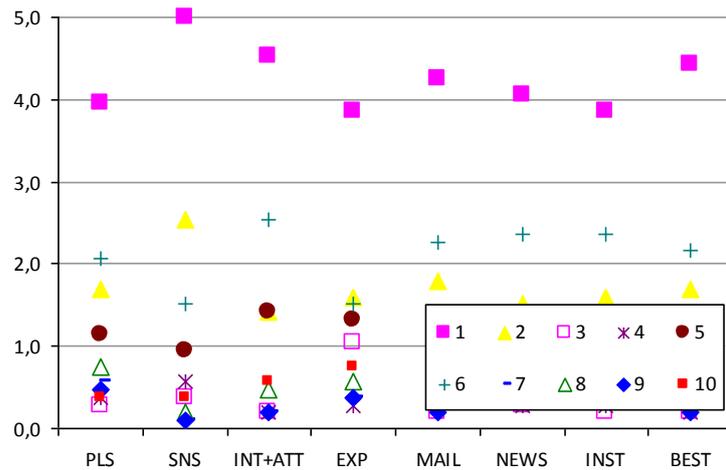


Fig. 2. Subjective test results for the 10 best scored voices on test B.

In Figure 2, the assessed subjective parameters in test B are displayed using a normalized 5 points scale for easy comparison with the results in test A. The “best” parameter corresponds to the question: “In your opinion, what is the best voice in general?”. Voices 1, 6 and 1 have the best scores and are highly separated from the remaining. Voice 1 occupies the first place in all of the subjective parameters. Voice 6 has the 2nd best score, except in SNS, in which Voice 2 is best scored. We believed that the listeners could be influenced by the first voice they hear during the survey/test and hence that the obtained results for voice 1 could be biased. In half of the surveys the voice listening order was changed. The score of the voices that were listened first

started to increase as suspected however voice 1 continued to receive a number of votes above the average.

It is interesting to compare the obtained results for test A and B, since the voice ranking has dramatically changed. In test A, the best scored voices were voices 10 and 3, but in Test B, the best scored voices were 1 (rated 4.8), followed by voices 6 (rated 2.2), 2 (rated 1.7) and 5 (rated 1.2) some of the worst scored voices in Test A. A possible explanation for these results is the content of voice samples used in test A. Sample for voice 8, for instance, was a popular TV spot advertising a prestigious car brand with a sophisticated music behind. Many samples were commercial spots or parts of documentaries. Another fact that may have influenced the first group of listeners was the presence of image data, as some voice samples were video clips. To prevent these variation factors in Test B, the voice content was a previously provided text to all candidates recorded with no background music or effects, under the same studio conditions.

The results were totally opposite as firstly expected. Voice 1 is positioned in the best place in most of the considered subjective attributes, immediately followed by Voice 6, which had almost the same overall score in Test A. These results seem to confirm the importance of this methodology in voice talent selection process.

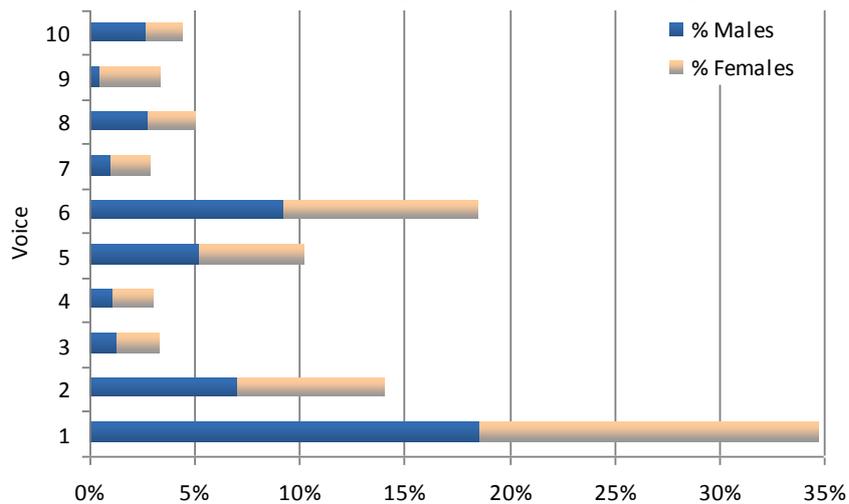


Fig. 3. Test B overall results for the 10 best scored voices according with listeners' gender. The percentage is normalized according to the total number of survey answers for each gender (89 males and 35 females)

The victory of Voice 1, in Figure 3, shows that both male and female listeners agreed when judging the voice quality. The same happened with the other voices with ratios around 50% for each gender. In a more detailed analysis we concluded that women seem to prefer dynamic and professional voices, rather than sexy voices.

In Table 1, a subjective parameters correlation matrix is drawn, using the following equation:

Table 2. F0 analysis.

#	Avg. (Hz)	St. Dev. (Hz)	Min. (Hz)	Max. (Hz)	Dif. (Hz)
1	226,6	39,0	105,7	253,7	148,0
2	223,5	47,1	112,0	303,4	191,5
3	252,6	37,6	143,2	304,3	161,1
4	231,3	35,8	112,0	283,8	171,7
5	242,1	35,8	104,3	286,3	182,0
6	233,2	37,0	159,7	299,1	139,4
7	258,2	37,2	158,6	307,7	149,1
8	261,5	30,4	179,3	305,4	126,2
9	235,5	40,4	158,1	303,1	145,0
10	243,5	34,4	162,0	284,9	122,9
Min.:	223,5	30,4	104,3	253,7	122,9
Avg.:	240,8	37,5	139,5	293,2	153,7
Max.:	261,5	47,1	179,3	307,7	191,5

Table 3. Energy analysis.

#	Total (dB)	St. Dev.	Avg. Int.
1	62,80	0,09	0,07
2	62,60	0,08	0,06
3	69,50	0,14	0,12
4	68,40	0,12	0,10
5	58,30	0,06	0,05
6	63,20	0,08	0,07
7	64,60	0,11	0,09
8	67,00	0,10	0,08
9	69,80	0,14	0,11
10	66,60	0,09	0,08
Min.:	58,30	0,06	0,05
Avg.:	65,28	0,10	0,08
Max.:	69,80	0,14	0,12

Table 4. Speaking rate analysis.

#	Speak. Rate (words/sec)	Speak. Rate (%)	Pause Rate (%)
1	3,07	54,00	46,00
2	2,01	69,50	30,50
3	2,90	64,70	35,30
4	2,62	62,70	37,30
5	2,60	78,30	21,70
6	2,48	66,90	33,10
7	3,11	62,90	37,10
8	2,40	68,50	31,50
9	2,49	70,40	29,60
10	2,68	56,10	43,90
Avg.:	2,64	65,40	34,60
St. Dev.:	0,33	7,08	7,08

Table 5. Score correlations.

Feature	Correlation
F0 Average (Hz)	-0,57
F0 Min. (Hz)	-0,46
F0 Max. (Hz)	-0,71
F0 Std. Dev. (Hz)	0,27
F0 Range (Hz)	0,05
SPR (words/sec)	0,13
SPR (%)	-0,33
Pause Rate (%)	0,33
Total Energy (dB)	-0,55
Energy St. Dev.	-0,52
Avg. Intensity	-0,58

Table 6. Score correlation with formant frequencies and related bandwidths calculated for some of the most representative vowels.

A		6		i		o		u	
F1	F2								
-0,36	-0,30	-0,57	-0,52	-0,32	-0,28	0,04	0,02	-0,45	0,21
B1	B2								
-0,18	-0,19	0,36	-0,27	-0,19	-0,11	0,29	0,17	-0,26	-0,21

4 Objective Tests and Results

The following acoustic parameters were assessed in the fourth stage of evaluation: F0 (average, maximum, minimum, range and standard deviation), energy (mean and standard deviation), intensity, speaking rate (SPR in words per minute excluding pauses) and pausing rate (PAR) (total duration of voice sample without pauses).

In Tables 2 to 5, F0 analysis, energy, speaking rate and pausing are displayed for each voice. Tables' inspection show that listeners prefer: 1) female voices with low mean F0 ranging between 223 and 240 Hz; 2) a discourse with well defined pauses between utterances and with a high speaking rate. Fundamental frequency and energy seems to be the most important parameters to evaluate when judging voice quality. Data displayed in Table 5 show that high F0 peaks contribute to a low scores. The importance of a voice with well-defined pausing is confirmed and shows that listeners seem to prefer voices with low minimum F0s. The correlation equation (1) was also used in this analysis.

Three phonetically rich sentences for each voice (around 15 seconds of speech) were manually segmented and labeled. The voice score was correlated with the vowel's formant frequencies and bandwidth. Table 6 shows that the obtained results for the fundamental frequency can also be extended to specific phonemes although /o/ seems to have very little importance on voice scoring. It is also interesting to note that the correlation values for each phoneme are very similar but they show discrepancies between them.

5 Resynthesis tests

Based on some authors who state that the perceived quality of a natural voice may not predict its synthetic quality [12], a different subjective test with the resynthesized voices was also conducted. Our goal with this test was assessing the 10 natural voices used in Test 2 after resynthesis. The chosen manipulated parameters in resynthesis were F0 and durations which were both changed with the PSOLA algorithm. F0 was set constant to the mean value and energy was normalized. Durations were changed from 80% in the beginning to 120% in the end of the sentence. Two sets of voice samples were presented to 8 different listeners who were not familiarized with TTS technology. The most interesting conclusion drawn from this test is the consensual opinion of Voice 10's good quality, which was secondly scored in natural and resynthesized voice quality rankings.

The best voice will be used in TTS systems for desktop computers but also in mobile devices which in comparison have a reduced hardware performance and perform audio processing in smaller sample rates. In order to evaluate the voice quality in these situations the resynthesized acoustic utterances were downsampled to 8KHz presented to the same listeners. This test showed that the global rating for each voice changes due to frequency downsampling but the ranking position for each voice remained unaltered. However, as expected due to frequency restriction, the differences

between each voice were severely reduced with more similar results for voices 1 and 6 (1st and 2nd in the overall ranking).

6 Conclusions

In this paper, a methodology in four stages to assess TTS voice font quality is proposed. Subjective tests and objective analysis and correlations were described and results were discussed. The following main conclusions can be drawn from this work: 1) a good voice is subjectively pleasant, assertive, expressive and preferably between 26 and 30 years old; 2) the corresponding objective characteristics are F0 ranging between 220 and 240 Hz, low minimum F0s and high speaking rates combined with well-defined pausing. These results can be considered in synthetic speech generation. As future work, we intend to extract additional parameters and evaluate their correlation with the obtained subjective results. As an extension of this research, a similar study for European languages voice quality assessment is ongoing. Based on our present experience, we propose the following changes in order to improve our methodology: 1) not allow video clip files as voice samples; 2) randomly change the order of voice samples every time a listener goes through the test or survey; 3) final selected voices should be no more than 8, in order to prevent judgment dispersion. Though our study is mainly focused on voice analysis we can add that other psychosomatic parameters should be considered. The voice owner should be easy working and must have good stamina since the recording process can be long; If EGG recording is used (typical procedure for high quality databases) a neck with little fatness can lead to better signals and reduced discomfort.

References

1. R. J. Bake, *Clinical measurement of speech and voice*. College Hill, Boston, 1987, pp. 197-240.
2. J. Kreiman, B.R. Gerratt, G.B. Kempster, A. Erman, "Perceptual Evaluation of Voice Quality. Review, tutorial and a framework for future research", *J. Speech and Hearing Research*, vol. 36, 1993, pp. 21-40.
3. L. Eskenazi, D. G. Childers, D. M. Hicks, "Acoustic Correlates of Vocal Quality", *J. of Speech and Hearing Research*, Vol.33, 1990, pp. 298-306.
4. J. Kreiman, B. R. Gerratt, "Sources of listener disagreement in voice quality assessment", *J. Acoust. Soc. Amer.*, 108 (4), 2000, pp. 1867-1876.
5. J. Laver, *The Phonetic Description of Voice Quality*, Cambridge University Press, Cambridge, 1980.
6. ITU-T Recommendation P.85, Telephone transmission quality subjective opinion tests. A method for subjective performance assessment of the quality of speech voice output devices.
7. M. Viswanathan, M. Viswanathan, "Measuring speech quality for text-to-speech systems: development and assessment of a modified mean opinion score (MOS) scale", *Computer Speech and Language*, vol. 19, January 2005, pp. 55-83.
8. S. Lemmety, *Review on Speech Synthesis Technology*, Master's Thesis, Helsinki University of Technology, 1999.

9. C. Gobl, A. Chasaide, "Testing affective correlates of voice quality through analysis and resynthesis", Proc. of the ISCA Workshop on Speech and Emotion, 2000, pp. 178-183.
10. O. Turk, M. Schröder, B. Bozkurt, L. M. Arslan, "Voice Quality Interpolation for Emotional Text-to-Speech Synthesis", Proc. of Interspeech 2005, Lisbon, Portugal.
11. N. Campbell, P. Mokhtari, "Voice Quality: the 4th Prosodic Dimension", 15th Intern. Congress of Phonetic Sciences, 2003, pp. 2417-2420.
12. A. Syrdal, A. Conkie, Y. Stylianou, "Exploration of Acoustic Correlates in Speaker Selection for Concatenative Synthesis", Proc. of ICSLP 98, 1998.