

MICROPHONEMICS - HIGH QUALITY SPEECH SYNTHESIS BY WAVEFORM CONCATENATION

Konrad Lukaszewicz

Inst. of Biocybernetics and
Biomed. Eng., PAN, Warsaw
Poland

Matti Karjalainen

Helsinki Univ. of Technology
Acoust. Lab., Otakaari 5 A, Espoo
Finland

ABSTRACT

Speech synthesis by waveform concatenation has been the subject of many attempts with fairly low quality results. We have reformulated the approach and found its potential to natural and personal-sounding speech by rule-based synthesis. Our study in Finnish and Polish shows that the method called *microphonemics* could be implemented by standard micro-processors and D/A-converters without any expensive signal processing hardware.

The main problems to be solved in the microphonemic method were the interpolation of pitch-sized phoneme and allophone units in wide formant transitions, the synthesis of fricatives and some other consonant classes, and the control of pitch and intonation. We found that the waveform interpolation works if the formant transitions are narrower than 2 Barks (critical bands), which implies the use of intermediate units in wide transitions. Fricatives are realized by time-randomized selection of 10 ms signal units from 50 ms unvoiced prototypes. Pitch and intonation problems can be solved by several windowing techniques in the formation and concatenation of pitch-sized units. The paper describes our experiments and proposes synthesis-by-rule strategies for implementation.

INTRODUCTION

The methods of speech signal generation in speech synthesis are often divided into two main classes: *model-based* source-filter models (formant and LPC-synthesis) and *waveform-based* time-domain synthesis methods. The advantage of model-based synthesis is the flexibility of generating an infinite number of signals according to parametric controls that can be computed by rules, tables etc. This has become the major method especially in speech synthesis by rule.

Time-domain synthesis can be based on a collection of varying sized speech signal units like waveform cycles, pitch periods, sound segments, phonemes, diphones, syllables etc., taken from real speech. Concatenation of speech samples is a simple method that has been used in synthesis experiments of low to moderate quality. In principle the sound quality could be very high if it were possible for enough samples of natural speech to be stored and carefully combined. This method takes more memory than the model-based synthesis but otherwise it is not as complex and arithmetically intensive.

A well known example of time domain speech synthesis is the *Mozier method* [1], where pitch-period-sized prototype units of real speech are manipulated to take as little memory as possible but are still able to be reconstructed in an intelligible form. This moderate quality, low bit rate method is used in some limited vocabulary synthesizers. Our experiments show that the tricks like zero-phasing the signal to lower the bit rate tend to remark-

ably reduce the quality and speaker identity. The phase properties are important to be retained for very high quality, natural sounding speech in a similar way as in multipulse LPC-coding [2].

The term "*microphonemic method*" that is used in our study was adopted from early experiments of similar principles in Poland. *Patryn* [3] synthesized with phonemic units without transitions and pitch changes. His work was continued by *Kielczewski* in his doctoral thesis (1979). This microphonemic method applied pitch changes for intonation and transitions by mixing parts of neighbouring phoneme prototypes. *Lukaszewicz* et al. have worked on the method at the Institute of Biocybernetics, Warsaw, since 1980. Their synthesizer has found applications in a talking typewriter and a talking calculator.

The quality of speech in all of these synthesizers has been low to moderate. The objective of our study was to find methods to overcome the inherent difficulties in concatenating speech waveforms. Our experiments show that it is feasible to develop simple and inexpensive synthesizers with natural and high-quality human-like characteristics. This concerns also speech synthesis by rule with unlimited vocabulary.

PROBLEMS TO BE SOLVED IN THE USE OF SPEECH WAVEFORM CONCATENATION

The microphonemic method is based on modelling the time domain signal by using a dictionary of prototypes. These are derived from natural speech utterances and their size can be of different lengths. It is possible to store whole words, syllables, phonemes (allophones) or shorter segments. Using a dictionary of microphonemes and several rules it is possible to generate synthetic speech by concatenating prototypes one after another. Waveform interpolation and concatenation are applied to realize the transitions between consecutive units. There are several problems that need to be solved in order to obtain high quality synthetic voice, e.g.:

- * realizing dynamic and static variations of the units, especially in the generation of smooth and natural transitions between consecutive segments and phonemes,
- * synthesizing consonants, like tremulants (Finnish /r/), etc.,
- * modifying parameters to control intonation, stress and rhythm,
- * determining the prototype set which is needed for a good representation of natural speech,
- * extracting these prototypes and their positions in the uttered speech examples,
- * formulating a good strategy when using waveform concatenation for synthesis by rule.

Some of these problems were studied by us at the Helsinki University of Technology, Acoustics Laboratory, by using the following experimental techniques.

WIDE FORMANT TRANSITIONS

The first problem to be solved in high quality waveform concatenation is to realize formant transitions e.g. in diphthongs (like /ui/ in Finnish) and in glides. The original idea of the microphonemic method was to apply simple linear interpolation from one pitch prototype to another by amplitude mixing (see Fig. 1). In our experiments we found that this works satisfactorily only if the glide in formant frequencies is less than 2 Barks (critical bands). In wider transitions, the amplitude-based interpolation is not sufficient to introduce a perceptually acceptable formant movement effect. For highest quality speech even 1 Bark transitions may be needed.

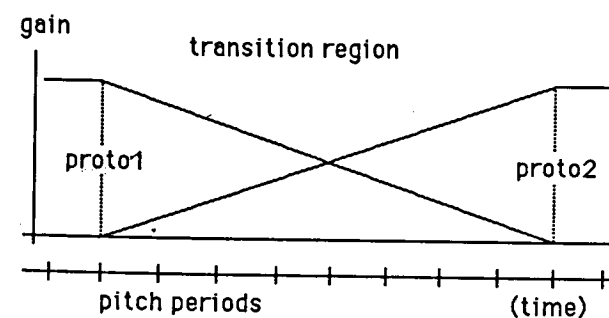


Fig. 1. Linear amplitude-based interpolation between two pitch-sized prototypes to simulate formant transitions.

If the formant distances between sound segments larger than 2 Barks are needed, some intermediate prototypes should be used to interpolate through (see Fig. 2.). It was possible for all transitions found in Finnish and Polish to be synthesized in this way.

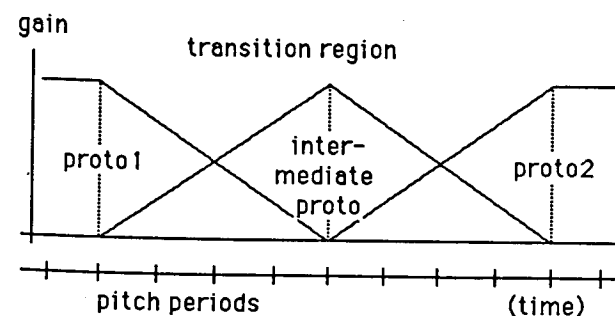


Fig. 2. Linear amplitude-based interpolation between two pitch-sized prototypes with an intermediate prototype.

SYNTHESIS OF CONSONANTS

Many consonants need special processing. Short non-repetitive units like bursts in stop consonants can be stored as direct waveform segments and as several variants in the context of different vowels or vowel groups. Sometimes the effect of neighbouring consonants must also be analyzed and the context stored for synthesis.

Fricatives need special treatment, too. Prototypes of about 50 ms in total length were found to be suitable and 10 ms units from them were randomly taken for concatenation. The same

interpolation rule as in vowels can be applied. Most voiced consonants behave in the same manner as vowels except that the variability according to the context is only higher.

PITCH AND INTONATION CONTROL

Prosodic features reveal some difficulties in concatenation. A simple and fairly successful method to control intonation is the use of minimum-pitch-period-sized prototypes and insertion of zero-signal segments to obtain the desired effective pitch for each moment (Fig. 3). A suitable windowing technique and the overlapping mixing of pitch periods could improve the results still further (Fig. 4). Timing is controlled by counting a proper number of pitch periods.

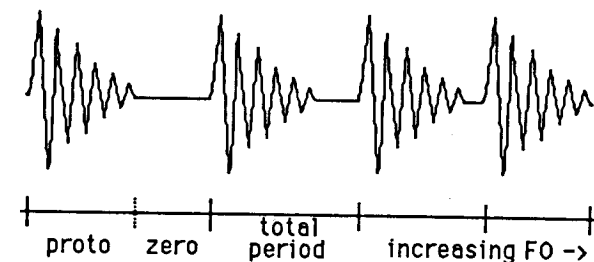


Fig. 3. Zero signal insertion as a method of controlling pitch in concatenation.

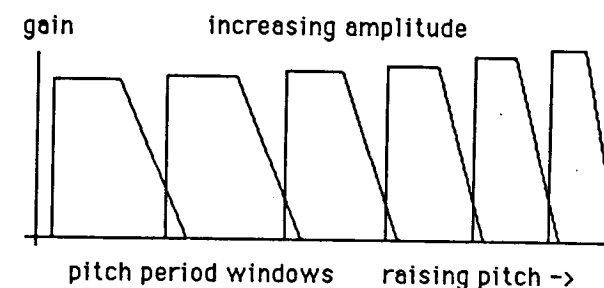


Fig. 4. Overlapping window summation in pitch control.

EXPERIMENTAL STUDY

About 70 Finnish and Polish phoneme pairs, concentrating on the synthesis of diphone-like transition segments, were studied experimentally by the microphonemic method. Some other larger units (syllables, words) were also modeled.

A microprocessor-based signal editor (SPS-02) was used to extract the prototype units from real speech. The same editor system was further applied to scale the amplitude, adjust the pitch period and to mix the prototypes for concatenation and synthesis experiments. Another analysis system, ISA [4], with auditory spectrum and spectrogram display was used to pick up the best positions of the prototypes and to compare the original against the synthetic speech examples. The principle of the auditory model for this analysis is presented in [5].

Prototypes from the original speech were used to model the phoneme pair transitions with two different principles of prototype selection. The first one was for producing intelligible, moderate quality speech with a minimum number of prototypes which were located in the middle of the quasiperiodic steady-state phonemes and one prototype in the middle of the transition.

The other method was to produce higher quality speech with a larger number of prototypes. This was accomplished by choosing the prototypes at each point where the formant frequencies started to change. If the change was larger than 2 Barks an extra prototype between the starting and ending points was taken. The maximum difference in any formant transition to be interpolated was always less than 2 Barks. A prototype was selected also at the points where the formants changed their direction of movement. For a full synthesis system some of the intermediate prototypes may be selected so that they can be used in several contexts.

As an example, the number of prototypes in the Finnish diphthong /ia/ was three for intelligible and seven for high quality speech. The maximum number of prototypes was never larger than nine for any diphthong-like unit. The size of a prototype was usually equal to one pitch period. However, in the case of stop consonants the length of a prototype was two to five times longer and for fricatives five times longer.

Fig. 5. shows the auditory spectrogram of the original diphthong utterance /ia/ with the related loudness function. Vertical lines with the capital letters A to C mark the positions of the prototypes in the lower-quality experiment. The auditory spectrogram of the synthesized version is shown in Fig. 6. Lines related to digits 1 through 7 in Fig. 5 indicate the places of the prototypes in the case of the highest quality reconstruction. The corresponding auditory spectrogram is in Fig. 7.

The pitch-sized prototypes from real speech naturally inherit some speaker-specific features and personality of the voice. The time-domain signal carries the tone quality features related to the detailed amplitude and phase spectrum. Our experiments show that the phase, especially rapid phase transitions can be very important to the naturalness of some allophones (nasals, liquids etc.) and their combinations. The prototypes may also include inherent pitch and amplitude data of the allophones that will be modified according to the context during the resynthesis.

MICROPHONEMIC SYNTHESIS BY RULE

Lukaszewicz et al. have implemented a low-to-moderate quality rule-based microphonic synthesizer in Polish with some practical applications. The objective of the present study was to find the feasibility of the microphonic method in high-quality synthesis by rule. Because of the relative high storage required and the tedious prototype preparation the method is not as well suited to limited vocabulary synthesis.

The synthesis process in microphonemics consists of the concatenation of precompiled prototype units with some context dependent modification rules applied to the prototypes. When compared to traditional model-based parametric synthesis this means more like operating with discrete symbol-like units instead of continuous-time control parameters. Some arithmetic and numeric computation can be avoided in this way.

The higher levels of text-to-speech synthesis transform the input text to a string of phonemic symbols. This process is very language dependent. In Finnish it is almost a one-to-one mapping from grapheme string to phoneme string with some prosody control analysis /6/, while e.g. in English it is a much more complicated task. The assembly of microphonemes into speech signals by phonemic level control information follows the same guidelines in all languages. A set of rules defines how the prototypes are to be modified and concatenated and how the prosodic control information is taken into account.

In our semimanual experiment we used a special notation to describe the assembly of microphonic units. The following forms were used:

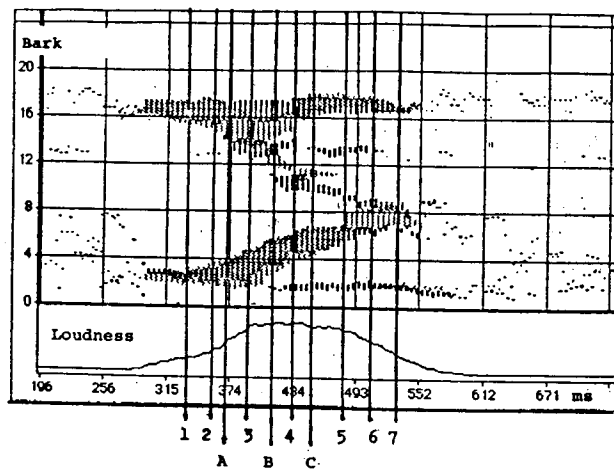


Fig. 5. Auditory spectrogram of the original speech, (Finnish /ia/)

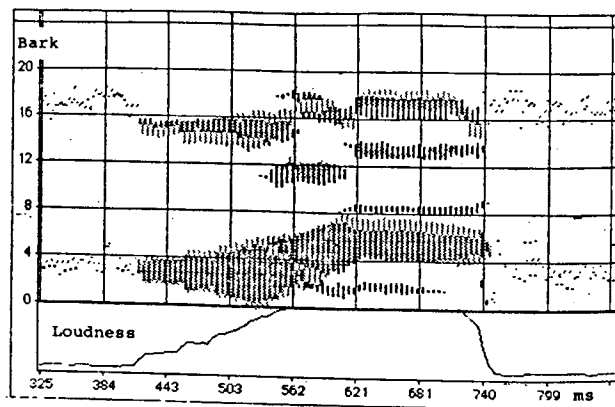


Fig. 6. Auditory spectrogram of the lower-quality reconstruction by the microphonic method with three prototypes (A, B, C in Fig. 5.)

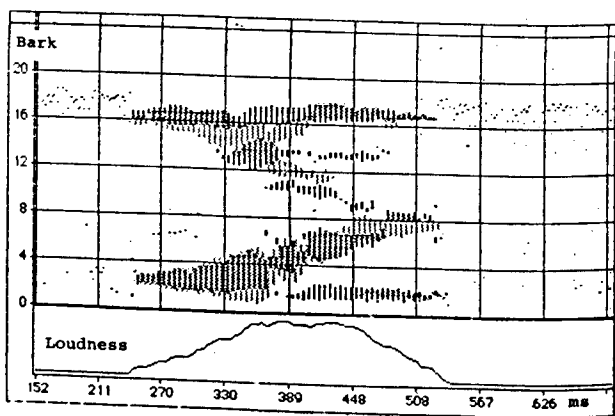


Fig. 7. Auditory spectrogram of the higher-quality reconstruction by the microphonic method with 7 prototypes (1 to 7 in Fig. 5.)

- x means a phoneme
- . is a transition region
- x_1x_2 is one prototype (one pitch period) which was taken from the beginning of a transition between x_1 and x_2 .
- $x_1 \cdot x_2$ one prototype taken from the middle of the transition between x_1 and x_2
- $\cdot x_1x_2$ one prototype from the end of the transition between x_1 and x_2
- n(...) integer to show the number of repetitions of some units, e.g. $5(x_1x_2)$
- n(-) n (integer) periods of linear interpolation of two neighbouring prototypes, e.g. $4(-)$, - equals to $1(-)$.

By using this notation we can express phoneme strings in the way of the following (Finnish) examples:

- /ai/ -> $12(@a) 5(-) a.i 5(-) 10(i.@)$
- /anna/ -> $17(@a) 15(.an) 9(na.) - n.a - 15(a.@)$
- /olli/ -> $@.o 5(-) ol. 5(-) 6(.ol) 5(-) 6(li.) 3(-) .li - 8(i.@) 3(-) 10(i.@)$

where symbol @ denotes space (pause).

This notation could be developed towards a formal rule language to be used in the implementation of the rule-based synthesis. It should also be possible to express the prosody-related control information, durations of the concatenated units (instead of counting periods), relative pitch and amplitude, special effects etc. To do this, the basic unit in the language could be an event object that contains fields or slots for different properties and relations.

The automatic generation of speech from phonemic code could proceed as follows. A rule-based match of the phonemic code to a set of templates is carried out to give the best candidate string of allophonic units and corresponding microphonic prototypes. Slot values related to prosodic features are filled based on context-dependent prosody rules. An experimental study of this kind is under development.

IMPLEMENTATION ASPECTS

An estimate of the memory capacity that is needed for prototypes in a moderate-quality synthesizer (Finnish) is: some 30 "phonemes", in average 8 variants (vowel contexts), and the same amount of intermediate prototypes. This results in a total number of about 500 units, each of 12ms in duration times 14 samples/sec (8 bits), which amounts to less than 100 kilobytes. At the level of present ROM-memory technology it is feasible to use up to 256 kbytes of memory for the prototype storage and synthesis rules, thus achieving high-quality synthetic speech with personal-sounding voice.

A single microprocessor like the Motorola 68000 is capable of doing this synthesis in real time. Serial and/or parallel ports are needed for input and a single D/A-converter (8 to 12 bits) with a reconstruction filter may be used to form the analog output. Another possibility is to design with multiplying D/A-converters so as to avoid software multiplications for amplitude scaling in the interpolation. The microphonic method is also well suited to software-based speech synthesis in microcomputers with special D/A-hardware to support fast analog output. The software for the microphonic synthesis by rule can be based on the manipulation of prototypes along the guidelines stated earlier.

The selection of the prototypes during the development of the system is a laborious and critical task that is difficult to be automated. A semiautomatic segmentation algorithm and pitch period detector could help if the voice of several speakers must be modeled. We are working to create two different development

systems to continue the studies on the microphonic method. One will be based on a personal computer, another in an artificial intelligence programming environment.

CONCLUSIONS

Our experiments showed clearly that the microphonic method by waveform interpolation and concatenation has potential for high-quality speech synthesis by rule. Its main technical advantage is that no computationally intensive signal processing is required. To achieve the highest-quality results optimal extraction of prototype segments from real speech and a good strategy for rule-based concatenation is needed. Auditory spectra and spectrograms were found important in the extraction process to find the best segments that meet the requirements of human auditory perception.

REFERENCES

- /1/ Costello B.C., Mozer F.S., Time-Domain Synthesis Gives Good-Quality Speech at Very Low Data Rates. Speech Technology Sept/Oct. 1982, p. 62-68.
- /2/ Atal B.S., Remde J.R., A New Model of LPC Excitation for Producing Natural-Sounding Speech at Low Bit Rates. Proc. of ICASSP-82, Paris, p. 614 - 617.
- /3/ Patryn R., Transitionless Synthesis of Speech. Acoustica Vol. 48 (1981) no. 4, p. 275-276.
- /4/ ISA, Intelligent Speech Analyser, Instruction Manual, Vocal Systems, Finland, 1987.
- /5/ Karjalainen M., A New Auditory Model for the Evaluation of Sound Quality of Audio Systems. Proc. of ICASSP-85, Tampa 1985, p. 608-611.
- /6/ Karjalainen M., An Approach to Hierarchical Information Processes with an Application to Speech Synthesis by Rule. Ma 29, Acta Polytechnica Scandinavica, Helsinki 1978.