

# EFFICIENT SUBSET SELECTION FROM PHONETICALLY TRANSCRIBED TEXT CORPORA FOR CONCATENATION-BASED EMBEDDED TEXT-TO-SPEECH SYNTHESIS

Aleš Mihelič<sup>1</sup>, Jerneja Žganec Gros<sup>1</sup>, Nikola Pavešić<sup>2</sup>, Mario Žganec<sup>1</sup>

<sup>1</sup>Alpineon, Ljubljana, Slovenia

<sup>2</sup> Faculty of Electrical Engineering, University of Ljubljana, Slovenia

**Key words:** embedded text-to-speech synthesis, speech corpus design, sentence subset selection

**Abstract:** In this paper we describe the design concept of a corpus-based concatenation text-to-speech (TTS) system for Slovenian, suitable for implementation in embedded applications. Because memory and processing power requirements are important factors when designing TTS systems for embedded devices, lexica and speech corpora need to be reduced. We describe a simple and efficient implementation of a greedy subset selection algorithm that extracts a compact subset of high coverage text sentences out of a larger set of sentences. The experiment on the Slovenian text corpus showed that the subset selection algorithm produced a compact sentence subset with a small redundancy. We conclude that the proposed sentence selection algorithm is capable of selecting a rather modest subset of sentences out of the reference text corpus, and the resulting sentence subset covers the most frequent collocations, words, quadphones and triphones in a given language.

## Postopek za izbor podmnožice stavkov iz besedilnega korpusa za sintezo govora v vgrajenih sistemih

**Ključne besede:** vgrajena sinteza govora, govorne podatkovne zbirke, izbira reprezentativne podmnožice stavkov

**Izvleček:** V članku opisujemo postopek zasnove in tvorjenja govorne zbirke za korpusno sintezo slovenskega govora, ki je primerna za implementacijo v vgrajenih sistemih. Postopek obsega izbiro besedila, snemanje in segmentacijo ter označevanje govornega gradiva. Sprva smo izvedli frekvenčno analizo pogostosti pojavljanja glasovnih sklopov za slovenski jezik nad obsežnim vhodnim besedilom, ki smo ga predhodno pretvorili v fonetični prepis. Nadalje opisujemo postopek, ki iz množice besedil v pisni obliki izbere kompaktno podmnožico stavkov, ki vsebujejo vsa želena pogosta zaporedja glasov v danem jeziku. Sledi snemanje ter posamodajno označevanje posnetega govornega gradiva. Članek sklenemo z rezultati preskusa naravnosti in razumljivosti sintetizatorja govora.

### 1. Introduction

A vital part of speech technology applications in modern human-machine user interfaces is a text-to-speech (TTS) engine. Text-to-speech synthesis enables automatic conversion into spoken form of any available textual information. Despite its considerable promise, text-to-speech synthesis is still not being used on a wide-scale basis in public service contexts. Wider acceptance of TTS devices will depend on three factors: better quality, smaller footprints, and more attractive pricing.

With the evolution of small portable devices, porting of high quality text-to-speech engines to embedded platforms has been made possible /1/, /2/. Many applications in mobile telephony and portable computing require high-quality speech synthesis systems with a very modest computational and memory footprint.

TTS systems, which are using a corpus-based concatenative approach, yield close-to-natural sounding speech /3/, /4/, /5/. However, the linguistic resources required to build embedded TTS modules need to be scaled down to meet the hardware specifications of the embedded devices. The major memory and processing power consuming linguistic resources that need to be reduced are lexica and speech

corpora /6/, /7/. The application of these reductions to Slovenian is demonstrated in the paper by utilizing efficient exception lexicon and speech corpus reductions. They were performed on the baseline full-size AlpSynth TTS system /8/, which uses a 95 Mb read speech corpus. A new, compressed speech corpus with a reduced set of read sentences was designed, which covers the most frequent allophone sequences in the language. The sentence subset has been selected from a large phonetically transcribed text corpus. The goal was to extract a sentence subset with high phonetic coverage and small size.

First an overview of the AlpSynth TTS components is given. We continue to describe the small-footprint speech corpus design process, with an emphasis on the sentence subset selection method. The quality of the synthesized speech was assessed in a listening experiment in terms of intelligibility and naturalness of pronunciation, whereby the small-footprint TTS system was compared against the baseline full-size AlpSynth TTS system. We conclude the paper by discussing the evaluation results and outlining plans for future work and concept implementation (Figure 1).

## 2. Concatenation-based TTS

In the AlpSynth TTS system, input text is transformed into its spoken equivalent by a series of modules: a grapheme-to-phoneme module produces strings of phonetic symbols based on information in the written text; a prosodic generator assigns pitch and duration values to individual phones; final speech synthesis is based on speech unit concatenation, where the elemental units are selected from a pre-recorded and annotated speech corpus and later concatenated using a pitch-synchronous overlap-and-add technique. The linguistic front-end speech synthesis phases used in the system are described in the following two subsections.

### 2.1. Grapheme-to-allophone conversion

Input to the TTS system is unrestricted Slovenian text. It is translated into a series of allophones in two consecutive steps. First, input text tokenization and token-to-word conversions are performed. Abbreviations are expanded to form

equivalent full words using a special list of lexical entries. The text normalizer converts further special formats, such as numbers or dates, into standard grapheme strings. The rest of the text is segmented into individual words and basic punctuation marks.

Next, phonetization or grapheme-to-phoneme conversion is performed. Word pronunciations are derived based on a user-extensible pronunciation dictionary and letter-to-sound rules. We constructed a dictionary that covers over 1,400,000 Slovenian inflected word forms. When dictionary derivation fails, words are transcribed using automatic lexical stress assignment and letter-to-sound rules. The use of rules enables the TTS system to generate a first attempt at pronunciations of neologisms and named entities.

To further reduce the memory footprint of the grapheme-to-allophone conversion module, we compiled an exception dictionary that contains only the differences from the phonetic transcriptions obtained by applying the letter-to-sound rule set. Similar to /7/, a compression factor of ten was

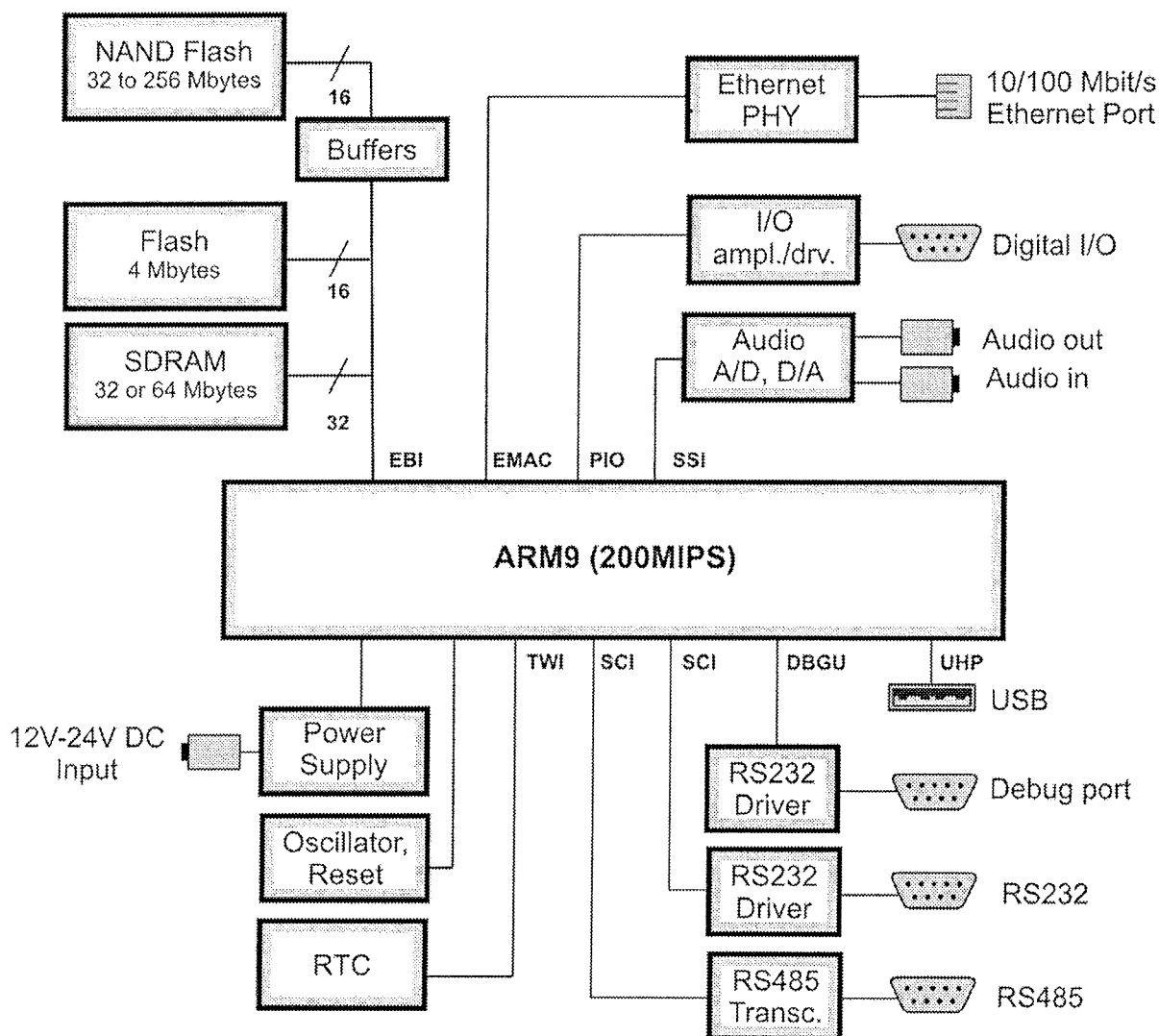


Fig. 1: Outline of hardware implementation of the reduced-footprint TTS system.

achieved, compared to the baseline full-lexicon representation, without sacrificing transcription accuracy. Further lexicon size reductions were achieved by modeling the exception lexicon in form of a decision tree, as proposed by /9/. We intend to test more approaches for efficient lexicon representation, among them finite-state transducers, which have already been applied for coding of language resources /10/, including Slovenian /11/.

## 2.2. Prosody Modeling

Corpus-based prosody modeling yields high-quality and close-to-natural sounding prosody parameter prediction; however, it requires a large amount of linguistic information upon which to rely. We used a compact rule-based prediction method to determine the target prosodic parameters in four phases: intrinsic duration assignment, extrinsic duration assignment, modeling of the intra word F0 contour, and assignment of a global intonation contour /8/.

Regardless of whether the duration units are words, syllables, or phonetic segments, contextual effects on duration are complex and involve multiple factors. A two-level duration model first determines the words' intrinsic duration, taking into account factors relating to the phone segmental duration, such as: segmental identity, phone context, syllabic stress, and syllable type: open or closed syllable /12/. Further, the extrinsic duration of a word is predicted, according to higher-level rhythmic and structural constraints of a phrase, operating at the syllable level and above. Here the following factors are considered: the chosen speaking rate, the number of syllables within a word and the word's position within a phrase, which may be isolated, phrase-initial, phrase-final or nested within the phrase.

Finally, intrinsic segment duration is modified, so that the entire word acquires its predetermined extrinsic duration. It is to be noted that stretching and squeezing does not apply to all segments equally. Stop consonants, for example, are much less subject to temporal modification than other types of segments, such as vowels or fricatives. Therefore, a method for segment duration prediction was used, which adapts a word with an intrinsic duration  $t_i$  to the determined extrinsic duration  $t_e$ , taking into account how stretching and squeezing apply to the duration of individual segments /12/.

Slovenian is a language with pitch accent, therefore special attention was paid to the prediction of tonemic accents for individual words. First, initial vowel fundamental frequencies were determined according to the parameters obtained from prior prosody measurements, creating the F0 backbone. Each stressed word was assigned one of the two tonemic accents characteristic for Slovenian. The acute accent is mostly realized by a rise on the post-tonic syllable, whereas with the circumflex the tonal peak usually occurs within the tonic.

## 3. Speech Corpus Design

For unit-selection and other types of concatenation-based text-to-speech synthesis, a speech corpus of recorded and

annotated elemental speech units is required /4/. The quality of the output synthetic speech depends crucially on the quality of the speech corpus. The longer elemental speech units are used, the better and more natural-sounding synthetic speech the TTS system can yield.

However, with longer elemental speech units the corpus size increases dramatically, as do the recording and annotation costs. Therefore, a compromise between the size of the speech corpus and the quality of the resulting speech has to be taken /13/ that is even more pronounced for embedded TTS.

If the corpus selection method is unbalanced or random, the recorded data may lack critical phone transitions and may be full of redundancies. Various corpus reduction methods have been reported, from those optimizing and reducing the contents of the prerecorded and annotated speech corpora to those that try to compress the initial text corpus to be recorded /14/, /15/, /16/, /17/, /18/. Often sentence pair exchanges are calculated using diphone and triphone entropies. In /14/, the unit coverage is maximized using prosody information. In /15/, a modified greedy algorithm is applied that maximizes the hit-rate and covering-rate for sentence selection criteria. A two-stage sentence recording script design presented in /17/ takes into account the balance of acoustic speech parts to provide variations in short-time speech features, and the linguistic parts provide long-time speech features, such as words or frequent word sequences.

We wanted the most frequent allophone sequences in a given language to be represented in the final sentence set, and therefore we implemented a greedy algorithm, similar to the one described in /15/, to reduce the initial text sentence set to a compact and efficient subset. The process of designing a speech corpus for concatenation-based TTS was divided into three phases:

- Representative sentence set selection,
- Recording of selected texts, and
- Segmentation and annotation of the recorded speech material.

### 3.1. Sentence subset selection algorithm

Initially, we collected a large corpus of texts covering various text styles, ranging from newspaper articles to fiction. All sentences shorter than five words or longer than 25 words were discarded from further analysis. The remaining reference text corpus contained 500,000 different sentences, corresponding to 50 Mb of text in ASCII format.

The text corpus was processed by a grapheme-to-allophone converter from the TTS system in order to obtain an allophone transcription of the text corpus. A statistical analysis of frequent phone sequences of allophones, diphones, triphones and quadphones was performed on this corpus. It provided us with information about how frequently certain phone combinations occur in spoken Slovenian. In addi-

tion, the analysis has shown that only a few triphones have frequent occurrences.

Therefore, it makes sense to select only the most frequent triphones to be represented in the final speech corpus. We opted for the first 1,000 triphones: these represent 1% of the complete triphone set but cover almost 50% of all triphones in the transcribed reference text corpus. In a similar way, the 500 most frequent quadphones were selected.

To synthesize high-quality speech, the speech corpus was required to contain a wide variety of speech parts: from collocations and words to diphones and sub-phoneme parts. With the most frequent triphones and quadphones selected, we wanted to select an optimal compact subset of corpus sentences that cover all the chosen allophone sequences, including most frequent collocations and words in a given language.

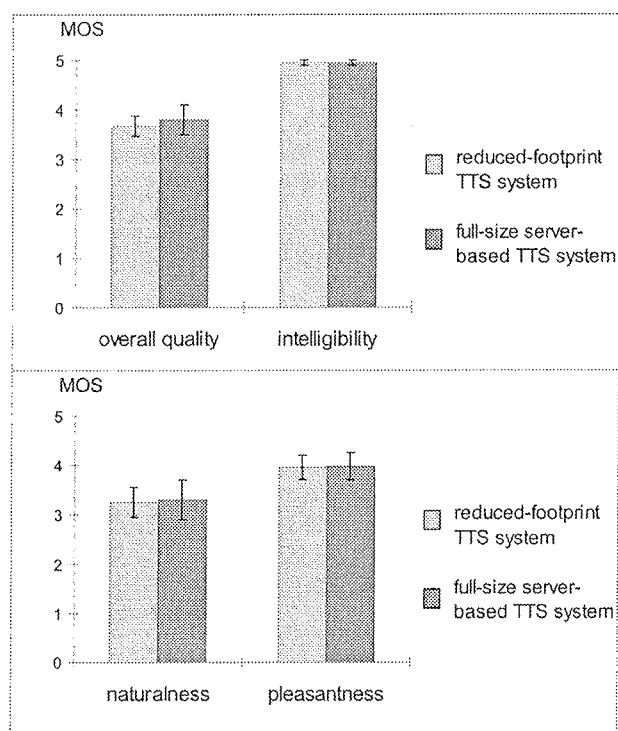


Fig. 2: Subjective evaluation results of the listening tests for both tested TTS systems. The results are given as MOS (Mean Opinion Score) ratings for the following categories: overall quality, intelligibility, naturalness, and voice pleasantness.

A greedy sentence selection algorithm was implemented for this purpose. Each sentence in the reference text corpus was equipped with a cost attribute based on the amount of the preselected frequent allophone sequences they contained. The highest cost value was attributed to a rare preselected quadphone or collocation, and the lowest to a frequent preselected triphone. In order to avoid the selection of long sentences, which contain more allophone sequences than shorter sentences, the cost value was normalized by the total number of allophones within the sentence.

The sentence with the highest score was selected for the final text corpus. The preselected allophone sequences covered by this sentence were eliminated from the list. Then the cost derivation and sentence selection process was performed for this new set of preselected allophone sequences and a new sentence was chosen for the final text corpus. The same process was repeated in a loop until all of the initial preselected allophone sequences were covered in the resulting corpus of selected sentences.

The sentence-selection algorithm was capable of selecting a rather modest subset of sentences out of the reference text corpus that cover the most frequent collocations, words, quadphones and triphones in the given language. A total of 299 sentences were selected out of the initial 500,000 sentences from the reference text corpus. The phonetic transcription of the selected sentence set covered all preselected most-frequent triphones and quadphones.

### 3.2. Recording and Segmentation

The selected sentence subset was recorded along with logatoms containing all phonetically possible diphone combinations for spoken Slovenian. The speaker was instructed to read the phonetically transcribed sentences and logatoms in supervised recording sessions.

The recorded speech material was segmented and annotated. A semi-automatic procedure was used for segmentation of elemental speech units. A dynamic time-warping acoustic alignment procedure between the synthesized voice and the recordings /19/ was used to obtain preliminary phone boundaries because it performed better in detecting consonant segment boundaries than the HMM approach /20/. The performance of the acoustical clustering plus dynamic time-warping method was upgraded along with boundary specific corrections by means of a decision tree, as recently proposed by /21/. Manual corrections were needed on consonant boundaries within some consonant clusters.

The final speech corpus contains read sentences with 1,993 words. In addition, 1,635 logatoms were recorded. For use in embedded devices, the speech corpus was compressed. Several compression techniques, outlined in /22/, were examined. Finally, a 14.1:1 corpus size reduction was chosen by using Ogg Vorbis encoding, without noticeably degrading the quality of the output speech signal. The resulting footprint of the compressed speech corpus was just below 2 Mb.

## 4. Evaluation Results

Over recent years, various guidelines have been proposed for evaluating the quality of text-to-speech systems. Yet there are still no existing standards for their evaluation, although a number of different methods have been tried and it has been pointed out that the test results they yielded were often inconsistent /23/.

The adequacy of the resulting concatenation-based TTS system was evaluated in terms of acceptability and intelligibility. The objective of the test was to compare the quality of the small-footprint TTS system to the baseline full-blown large-footprint unit-selection server-based TTS system /8/.

The experiment was performed in laboratory conditions with 51 test subjects. It was designed according to ITU-T Recommendations P.81 and P.85, describing methods for subjective performance assessment of the quality of voice output devices. The evaluators were selected from a wide range of professional backgrounds, and they were in general not familiar with synthetic voice quality. The test was divided into two sessions, neither lasting more than 20 minutes, in order to reduce the fatigue of the evaluators. The synthesis output was directed to a loudspeaker. Each test speech was presented only once.

The first part served to evaluate whether the intelligibility and the quality of the synthetic speech were sufficiently high for a real application of the system in a potential embedded-system application, simulating spoken directions provided by a car-navigation system. The subjects were asked to fill in different application-specific templates based on the information they heard. Each message consisted of a fixed part, which was specific to the task, and a variable part, which was different in all the produced messages. The intelligibility for both systems, when spelling errors are ignored, was nearly 100%. Over 95% of the listeners estimated that both TTS systems were mature enough for deployment in a car-navigation system.

In the second part of the test, the performance of both TTS systems was evaluated by the listeners with grades on a five-point MOS (Mean Opinion Score) scale. For the test, the sentences were synthesized by both TTS systems and presented to the listeners in random order. The listeners were asked to evaluate the overall quality, intelligibility, naturalness, and voice pleasantness. The results are provided in Figure 2. In terms of overall quality MOS grades, the full-size server-based TTS outperformed the reduced-footprint version by only a small margin of 0.15. The majority of the test subjects evaluated the reduced-footprint (as well as the large-footprint) synthetic speech produced by the TTS system as pleasant and quite natural-sounding, sufficiently rapid and not over-articulated.

## 5. Conclusions

The memory and computational resources in TTS applications on embedded portable devices are inherently limited. Various footprint reduction considerations for embedded TTS implementation are discussed in the paper. We concentrated on shrinking the speech corpus while maintaining high coverage of the frequent allophone sequences in a given language: our goal was to extract a sentence subset with high coverage and small size. The sentence subset was selected from a large phonetically transcribed text corpus.

The greedy sentence selection algorithm implementation described in the paper was capable of selecting a rather modest subset of sentences out of the reference text corpus that covers the most frequent collocations, words, quadphones, and triphones in a given language. A total of 299 sentences were selected out of the initial 500,000-sentence text corpus. The phonetic transcription of the selected sentence subset covered all of the preselected most-frequent triphones and quadphones, words, and collocations.

An implementation of the proposed sentence subset selection method for Slovenian has resulted in a small-footprint TTS system yielding intelligible and sufficiently natural-sounding speech, so that the system is ready for deployment in embedded applications. Listening experiments proved that the TTS system gives satisfactory performance in phonetization and speech concatenation quality with considerably reduced memory resources. The system is implemented in ANSI-C and runs on several operating systems. The object code size of the small-footprint TTS system is 98 Kb, while the size of the language resources used by the system is just below 2Mb. Using the designed compact speech database, the program's current version runs 300 times faster than real time on a Pentium 2 GHz personal computer. At run-time, the program code requires 472 Kb, minimum RAM requirement is 3,500 Kb, while the minimum disc or flash memory requirement is 5,560 Kb.

In order to further compress exception lexica, alternative lexicon representation approaches will be examined, such as /11/. An initial implementation of the described methods on an embedded TTS Unix platform built around an ARM9 processor with an AT91RM9200 core (Figure 1) is under development.

## 6. Acknowledgements

The authors wish to thank the Slovenian Ministry of Higher Education, Science, and Technology and the Slovenian Research Agency for co-funding this work under contract no. V2-0896.

## 7. References

- /1/ Black, A.W. and Lenzo, K.A., "Flite: a small fast run-time speech synthesis engine", In Proceedings of the 4<sup>th</sup> ISCA Workshop on Speech Synthesis, 2001, pp. 204-207.
- /2/ Tomokoyo, M.L., Black, W.A. and Lenzo, A.K., "Arabic in my hand: small footprint synthesis of Egyptian Arabic", In Proceedings of the Eurospeech'03, Geneva, Switzerland, 2003, pp. 2049-2052.
- /3/ Campbell, N., "CHATR: a high-definition speech resequencing system", In Proceedings of the 3<sup>rd</sup> ASA/ASJ Joint Meeting, 1996, pp. 1223-1228.
- /4/ Beutnagel, M., Conkie, A., Schroeter, J. and Stylianou, Y., "The AT&T Next-Gen TTS System", in Proceedings of the 137<sup>th</sup> Meeting of the Acoustic Society of America, 2000.
- /5/ Möbius, B. "The Bell Labs German text-to-speech system", Computer Speech and Language, Vol. 13, 1999. pp. 319-358.

- /6/ Tian, J., Nurminen, J. and Kiss, I., "Optimal subset selection from text databases", In Proceedings of the ICASSP'05, PA, USA, 2005.
- /7/ Meron, J. and Veprek, P., "Compression of exception lexicons for small footprint grapheme-to-phoneme conversion", In Proceedings of the ICASSP'05, PA, USA, 2005.
- /8/ Žganec Gros, J., Mihelič, A., Pavešić, N., Žganec, M., Gruden, S., "AlpSynth - concatenation-based speech synthesis for the Slovenian language", In Proceedings of ELMAR'05, Zadar, Croatia, 2005, pp. 213-216.
- /9/ Šef, T. "A two level lexical stress assignment model for highly inflected Slovenian language", In Proceedings of the International Conference on Information Technology and Applications, Sydney, Australia, 2005, pp. 347-351.
- /10/ Mohri, M., "On some applications of finite/state automata theory to natural language processing", Natural Language Engineering I, Cambridge University Press, 1996.
- /11/ Rojc, M., Kačič, Z., Kramberger, I., "Hardware implementation of language resources for embedded systems". Inf. MIDEM, Vol. 32, No. 3, 2002, pp. 199-203.
- /12/ Gros, J., Pavešić, N. and Mihelič, F., "Speech timing in Slovenian TTS", Proceedings of the Eurospeech'97, Rhodes, Greece, 1997, pp. 323-326.
- /13/ Van Santen, J.P.H., "Methods for optimal text selection", In Proceedings of the Eurospeech'97, Rhodes, Greece, 1997, pp. 553-556.
- /14/ Kawai, H., Yamamoto and Shimizu, T., "A design method of speech corpus for text-to-speech synthesis taking into account prosody", in Proceedings of the ICSLP'00, 2000, pp. 420-425.
- /15/ Kuo, C. and Huang, J., "Efficient and scalable methods for text script generation in corpus-based TTS design", in Proceedings of the ICSLP'02, 2002, pp. 121-124.
- /16/ Bozkurt, B., Ozturk, O. and Dutoit, T., "Text design for TTS speech corpus building using a modified greedy selection", in Proceedings of the Eurospeech'05, Geneva, Switzerland, 2003, pp. 277-180.
- /17/ Isogai, M., Mizuno, M. and Mano, K., "Recording script design for corpus-based TTS system based on coverage of various phonetic elements", In Proceedings of the ICASSP'05, PA, USA, March 18-23, 2005.
- /18/ Rojc, M. and Kačič, Z., "Design of optimal Slovenian speech corpus for use in the concatenative speech synthesis system", In Proceedings of the LREC'00, Athens, Greece, 2000, pp. 321-325.
- /19/ Malfrère, F. and Dutoit, T., "High quality speech synthesis for phonetic speech segmentation", In Proceedings of the Eurospeech'97, Rhodes, Greece, 1997, pp. 2631-2634.
- /20/ Mihelič, F., Gros, J., Dobrišek, S., Žibert, J. and Pavešić, N., "Spoken language resources at LUKS of the University of Ljubljana", International Journal on Speech Technologies, Vol. 6, No. 3, 2003, pp. 221-232.
- /21/ Xydas, G. and Kouroupetroglou, G., "An intonation model for embedded devices based on natural F0 samples", In Proceedings of the Interspeech'04, Korea, 2004, pp. 801-804.
- /22/ Hoffmann, J., Jokisch, O., Hirschfeld, D., Strecha, G., Kruschke, G., Kordon, U. and Koloska, U., "A multilingual TTS system with less than 1 Mbyte footprint for embedded applications", In Proceedings of the ICASSP'03, Hong Kong, 2003.
- /23/ Alvarez, Y. and Huckvale, M., "The reliability of the ITU-T P.85 standard for the evaluation of text-to-speech systems", In Proceedings of the ICSLP'02, Denver, CO, 2002, pp. 329-332.

*mag. Aleš Mihelič, dr. Jerneja Žganec Gros,  
dr. Mario Žganec  
Alpineon, Ulica Iga Grudna 15, SI-1000 Ljubljana,  
Slovenia  
info@alpineon.com  
tel +386 1 423 9440  
tel +386 1 423 9445*

*prof. dr. Nikola Pavešić  
Faculty of Electrical Engineering, University of Ljubljana  
Tržaška 25, SI-1000 Ljubljana, Slovenia  
nikolap@fe.uni-lj.si  
tel +386 1 476 8840  
tel +386 1 476 8319*

*Prispelo (Arrived): 03. 01. 2006; Sprejeto (Accepted): 30. 01. 2006*