

Speech intelligibility as a function of temporal resolution and the number of stimulation channels for signal processing using the sinusoidal speech model

O. Timms*, S. Allegro⁺, V. Kühnel⁺, A. von Buol⁺, N. Dillier*

*Dept. Otorinolaringology University Hospital Zurich Switzerland; ⁺Phonak AG, Stäfa, Switzerland

Abstract

The number of stimulation channels as well as the temporal and spectral resolution required for 100% speech intelligibility were determined for a sinusoidal speech signal processing system. German vowel, consonant and sentence tests were processed with an algorithm based on the sinusoidal speech model of Quatieri and McAulay [1]. The algorithm uses a limited number of stimulation channels (1-5) as well as different temporal and spectral resolutions (FFT frame length 128, 256 and 512 points, sampling frequency 22050 Hz). Speech intelligibility tests in quiet were performed with normal hearing native German speaking adults using conventional diagnostic speech perception tests as well as paired comparison consonant tests.

The results show that different numbers of stimulation channels are required for 100% speech recognition in sentences processed with different temporal and spectral resolution. A tendency of increased speech intelligibility was observed by either increasing temporal resolution or increasing the number of stimulation channels. In addition significant learning effects were observed.

Method

A block diagram of the present sinusoidal speech analysis / synthesis system is depicted in figure 1. The most important parts of the original algorithm of Quatieri and McAulay as FFT and peak picking are preserved. However, due to the high complexity of cubic polynomial phase interpolation, frequency tracking and sine wave generation in the original algorithm, an IFFT with a 75% overlap-add is used for synthesis. These simplifications result in a quality reduction of the processed signal. The most remarkable effect is musical or "bubble"-like noise as known from vocoder systems.

Three FFT lengths (128, 256 and 512 points) at a sampling frequency of 22050 Hz corresponding to time resolutions of 5.8, 11.6, and 23.2 ms were tested with 1 to 5 spectral lines.

For peak selection the lowpass filtered smoothed frequency spectrum as proposed by Kates was used [2].

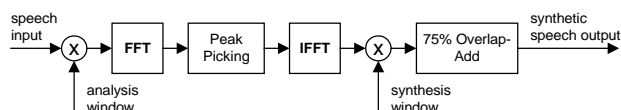


Fig. 1 Sinusoidal analysis/synthesis system block diagram.

Performed Tests

Oldenburg Sentence tests [3] were carried out in quiet with 16 native German speaking adults. The sentences were presented at 65 dB RMS in a sound proof room at a distance of 1.5 m from a Westra Audiometer Box Type LAB - 1001 (playback by a 16 bit PC sound card).

C12 Consonant aCa-Logatome tests were carried out in quiet with 21 native German speaking adults. Logatome syllabics were presented at 65 dB RMS in a sound proof room at a distance of 1.5 m from a Westra Audiometer Box Type LAB - 1001 (playback by a 16 bit PC sound card). The subjects were asked to identify the aCa-syllable from a choice of 12 aCa-syllables on a touch screen (MAC-test [4]).

Vo8 Vowel dV-Logatome tests were carried out with the same test setup as for the consonant tests with 14 native German speaking adults. The subjects were asked to identify the correct dV-syllable from a choice of 8 dV-syllables.

Consonant C12 aCa- Logatome paired comparison tests were carried out in quiet with 24 native German speaking adults. Only those consonant pairs were tested which caused difficulties

in the logatome identification tests. Logatome syllables were presented at most comfortable level by Sennheiser headphones Type HD 570 (playback by a 16 bit PC sound card) in a sound proof room. The subjects were asked to judge discrimination certainty of the two presented aCa-logatome syllables.

Results

Figure 2 shows the results of the Oldenburg sentence test. For 100 percent sentence recognition at a temporal resolution of 5.8 ms (128 pt FFT) only one stimulation channel is required. With decreasing temporal resolution the number of spectral channels required for 100 percent recognition increases up to four for a temporal resolution of 23.3 ms (512 pt FFT). The test shows that temporal resolution is critical for speech recognition in quiet. Furthermore, remarkable learning effects were observed during the tests (not depicted).

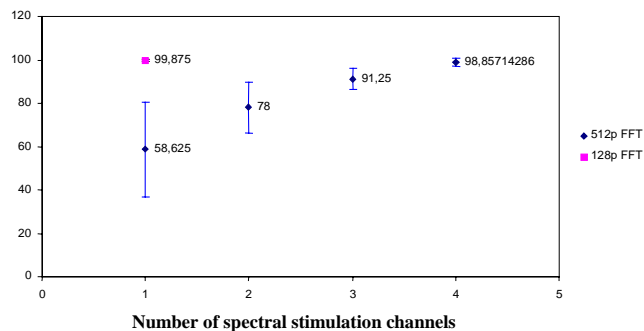


Fig. 2 Results of the Oldenburg sentence test. On the ordinate the score of right answers in percents are given.

In agreement with the results from the sentence tests, the identification scores of both vowels and consonants increase with increasing number of spectral stimulation channels and with increasing temporal resolution.

After being processed, some consonants sound unnatural due to musical noise effects, resulting in low identification scores for the logatome identification tests. In these cases, no dominant logatome confusions are observable in the confusion matrices (figures 3 to 6). For example, for each of the tested temporal resolutions the syllable "aRa" is poorly recognized even with five stimulation channels (see figures 3 and 4). However, in the paired comparison logatome tests the same syllable "aRa" is well separable from the other consonants (figure 7). This might be explained by the unnatural pronunciation of these consonants resulting from the processing artifacts. With training those consonants should become well recognizable. There are also some consonants and vowels, for example German "S" "A", which are well recognizable independent of the temporal or spectral resolution. For those phonemes the temporal resolution is less significant and the spectral resolution of 174 Hz (128 pt FFT) is sufficient.

The paired comparison tests showed that there are some confusions, for example German "N-M" or "F-P" (figures 7 and 8), where a defined number of stimulation channels is required to reach an acceptable rate of discrimination certainty. This consonants are especially interesting for future investigations.

Confusion matrix for 128 Point FFT with 1 stimulation channel based on a total of 33 presentations of each token.

	P	T	K	B	D	G	M	N	L	R	F	S
P	30	30	21	3							9	6
T	3	70	6	3	3						3	12
K		18	76								6	6
B	3		3	42	36	3					6	6
D		3	3	3	76	12						3
G		12		9	3	76						
M							82	15	3			
N	3						9	64	3	3	3	18
L		3						3	24	70		
R		3	3	6	15	9	6		3	48		6
F	24	12	3	6	9	6					24	21
S		15			3						3	79

Fig. 3 Confusions matrix for 5.8 ms temporal resolution with 1 stimulation channel (confusions are given in percents).

Confusion matrix for 128 Point FFT with 5 stimulation channel based on a total of 30 presentations of each token.

	P	T	K	B	D	G	M	N	L	R	F	S	
P	100												
T		100											
K			100										
B				100									
D					100								
G						100							
M							100						
N								3	97				
L										100			
R											93		
F	27		3	3	3							60	
S													100

Fig. 4 Confusions matrix for 5.8 ms temporal resolution with 5 stimulation channels (confusions are given in percents).

Confusion matrix for 512 Point FFT with 1 stimulation channel based on a total of 36 presentations of each token.

	P	T	K	B	D	G	M	N	L	R	F	S
P	28	25	6	8	8				3		17	6
T	14	19	3	14	8	3	3		3		17	17
K	22	31	17	6	11	6					6	3
B	6	3		72	6	6				3	6	
D	3	3	3	44	14	25				8		3
G				44	25	14	3		3	6	3	3
M	3	3	3	6	6	6	25	39	6		8	3
N	3	3	6	3	6	19	36	3	8	8	6	6
L	3	3	3	14	8		11	8	44		3	3
R	8	6	6	6	11	3	11	8	17	8	17	8
F	3	6	3	8	3					3	44	31
S	3	6		3							8	81

Fig. 5 Confusions matrix for 23.2 ms temporal resolution with 1 stimulation channel (confusions are given in percents).

Confusion matrix for 512 Point FFT with 5 stimulation channel based on a total of 33 presentations of each token.

	P	T	K	B	D	G	M	N	L	R	F	S
P	85	3	3								9	
T	3	91				6						
K			97									3
B				100								
D	3				97							
G						100						
M							97		3			
N								100				
L									100			
R	3			3	3	3	3			79	3	3
F	21			3							67	9
S												100

Fig. 6 Confusions matrix for 23.2 ms temporal resolution with 5 stimulation channels (confusions are given in percents).

In the paired comparison tests there are also some consonant confusions which show a decrease of discrimination certainty with growing number of spectral stimulation channels (see e.g. "F-P" in figure 8). These effects can be explained by the uncertainty of the observations.

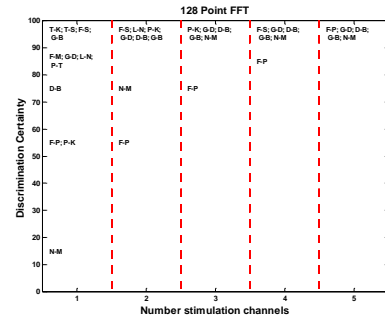


Fig. 7 Results of the Logatome aCa tests for 5.8 ms temporal resolution

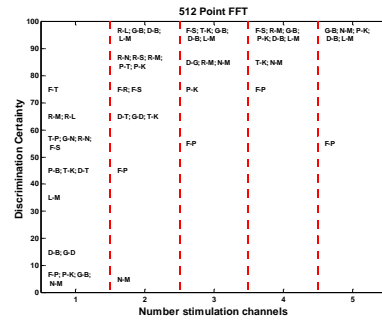


Fig. 8 Results of the Logatome aCa tests for 23.2 ms temporal resolution.

Conclusions

The use of three spectral stimulation channels seems to be sufficient for the recognition and discrimination of German vowels and consonants for each of the three tested temporal and spectral resolutions. With temporal resolutions of 5.8 ms and 11.6 ms, only two spectral stimulation channels are required for good vowel and consonant recognition. With a temporal resolution of 5.8 ms and the use of only one spectral stimulation channel, recognition difficulties occur with the four German consonants "P", "B", "R", and "F" and the three vowels, "I", "Ä", and "Ü".

It was shown that for speech recognition in quiet with a reduced number of spectral stimulation channels, temporal resolution is more important than spectral resolution. Remarkable learning effects during the sentence- and Logatome tests are can be seen. Paired comparison tests are necessary to get meaningful results.

References

- [1] T.F. Quatieri and R.J. McAulay, "Audio signal processing based on sinusoidal analysis / synthesis", chapter 9 of "Applications of Digital signal processing to audio and acoustics", edited by M. Kahrs and K. Brandenburg, Kluwer Academic Publishers, 1998.
- [2] J.M. Kates, "Speech Enhancement Based on a Sinusoidal Model," Journal of Speech and Hearing Research, vol. 37 pp. 449-464, 1994.
- [3] Wagener, K., Kühnel V. and Kollmeier, B. "Entwicklung und Evaluation eines Satztests für die deutsche Sprache I: Design des Oldenburger Satztests", Zeitschrift für Audiologie, 1, 4-15, 1999.
- [4] N. Dillier, T. Spillmann, „Deutsche Version der Minimal Auditory capability (MAC) – Test- Batterie“ in Kollmeier, B. Moderne Verfahren der Sprachaudiometrie, Median Verlag, Heidelberg, 1992.