

COMPARISON BETWEEN AUDITIVE AND AUDIO-VISUAL PERCEPTION OF PB-WORDS MASKED WITH WHITE NOISE

Carl Ludvigsen

1. Introduction

The basis of this paper is furnished by some experiments carried out two years ago at The State Hearing Center, Bispebjerg Hospital, Copenhagen. The purpose of these experiments was in broad outline to examine the influence on the discrimination score when PB-words were presented to listeners with a normal hearing belonging to three different age-groups, and the words were presented with and without the possibility of seeing the face of the speaker. A report on this experiment is given by Ewertsen et al. (1970).

The aim of the present paper is to study the answers obtained from the subjects under the two modes of presentation and four different signal to noise ratios (S/N).

1.1. The word material

The stimulus material consisted of four phonetically balanced lists (A, B, C, D) each containing 25 words. Due to the phonetic balancing the lists contained only mono- and disyllabic words, all commonly used Danish words. All disyllabic words had a trochaic stress pattern with the first syllable stressed and the second unstressed.

1.2. The presentation

A recording of the lists read by a male speaker whose dialect was close to standard Copenhagen was made on a video tape

recorder for use on internal television. The speech signal from the sound track on the video tape was fed through an attenuator and then mixed with noise. The noise level was kept constant throughout the experiment and different signal to noise ratios were obtained by attenuating the speech signal. The noise was approximately "white" within the audible frequency range. The sound signal was presented monaurally through a pair of headphones. The visual signal i.e. a frontal picture of the speaker's face could be presented synchronously with the acoustic signal on a TV-screen in front of the subject.

Each of the four lists was presented four times to every subject. The first presentation was given through the headphones, with no picture on the TV-screen, for the second presentation the TV film was added, and finally approximately one month later these two presentations were repeated.¹ The S/N's were different for the four lists: list A: S/N = -20 dB, list B: S/N = -10 dB, list C: S/N = 0 dB, and list D: S/N = +10 dB.

1.3. The subjects

28 subjects participated in the experiment. For the sake of studying the influence of age on the discrimination score these subjects were selected from different age groups. Thus, 9 subjects were from 20-25 years old, 10 from 45-55 years old, and 9 from 65-75 years old. Audiograms were taken of all subjects, and only subjects with audiograms normal for their group of age participated in the experiment. For the present study, however, it was decided to disregard answers from the oldest group mainly because of the pronounced hearing losses at high

1) This was done in order to study the effect of retesting.

frequencies which are typical for their age-group. Furthermore two persons from the group 45-55 years were discarded because of hearing losses of respectively 30 and 35 dB for high frequencies. The remaining 17 subjects had audiograms differing less than 25 dB from the normal threshold in the frequency range from 125 Hz to 8000 Hz.

1.4. Experimental procedure

The subject was placed in a quiet room normally used for audiometry. The experimenter was placed in an adjacent room from where he could watch the subject through a window. Near the subject was placed a microphone which allowed the experimenter to hear the replies from the subject. The subject was instructed in the testing routine and was asked to repeat the words as he heard them. The experimenter registered whether the subject was able to repeat the words and if the subject answered incorrectly, i.e. with a word which was not identical to the stimulus, he noted this word on the list.

2. Analysis of answers

As the subjects were not forced to answer, three alternative types of reply are possible: 1) the subject may not have answered or he may have answered 2) with the correct word or 3) with an incorrect word. In the first case no information is obtained about the perception of the stimulus. In the second case some uncertainty exists concerning the cues used for identification of the word. Thus it is possible that a certain feature of a stimulus may not be detected by the subject although a correct answer is given, since the identification may have been based solely upon other features. The most useful source of information about the mechanism of perception seems to be the incorrect answers. Thus a comparison between an incorrect answer and the stimulus word provides information about which cues are detected and which are not.

2.1. Detection of number of syllables

If we assume that the subjects are unable to detect the number of syllables we shall expect the number of syllables found in the incorrect answers to be distributed approximately in the same way as in a normal word material and independently of the number of syllables in the stimulus word. This hypothesis can clearly be rejected from the material: Even at the most unfavourable S/N in the auditive tests the detection of number of syllables is very accurate. This is shown in TABLE 1 below.

TABLE 1

STIMULUS	NUMBER OF WORDS	PRESENTATIONS	NUMBER OF			NUMBER OF WRONG	
			CORR. ANSW.	NO ANSW.	WRONG ANSW.	ANSW. WITH 1 SYLL.	2 SYLL.
2 SYLL. WORDS	12	408	52	298	58	0	58
1 SYLL. WORDS	13	442	33	360	49	47	2

TABLE 1. Answers from all 17 subjects pooled (list A, auditive test, S/N = -20 dB)

Table 1 shows that although more than 50 % of the answers are incorrect almost all of these contain the correct number of syllables. This finding also indicates that addition of the visual signal will not improve the detection of syllables appreciably.

2.2. Detection of the unstressed vowels in the disyllabic words

A cursory glance at the incorrect answers tells that the subjects rarely fail to detect the vowel in the second, unstressed syllable of disyllabic words. This impression holds true even for the tests with the smallest S/N and with auditive presentation only. In list A the second syllable of 11 of the 12 disyllabic words were of the type (C)ə(C) (zero or one consonant + schwa + zero or one consonant) and only one ended in a different vowel, viz. ʌ (the word was "tænder" (teeth)). To this word were given 10 answers; six of these were incorrect but all ended in unstressed ʌ. To the remaining 11 words 104 answers were obtained 52 of which were incorrect. 51 of these ended in Cə(C) and only one in a different vowel (unstressed, short i).

This very accurate detection of the unstressed vowels was to be expected: since approximately 90 % of Danish disyllabic words with stressed first syllable have schwa as unstressed vowel, the a priori uncertainty about the identity of the unstressed vowel is relatively small. On the other hand a reduction in intensity or duration of unstressed vowels might be expected to make the identification difficult.²

2.3. Confusions between stressed vowels

A larger percentual number of errors occur among the stressed vowels. And, of course, the nature of the errors depends heavily on the mode of presentation, auditive or audio-visual, and the S/N. In order to study the confusions among stressed vowels the words were grouped with respect to the acoustic quality of the stressed vowel. Obviously, this grouping can be done

-
- 2) Intensity curves of the stimulus material showed no reduced intensity for the unstressed vowels.

in more or less detail. For the present purpose the grouping was based on the frequencies of the two lowest formants, F_1 and F_2 , as found from a frequency analysis of a tape recording of the stimulus material. The formant frequencies of vowels with marked transitions of F_1 and F_2 were measured at the most steady part of the vowel (generally in the middle) or, if no such part could be found (as in diphthongs), in the first part of the vowel. The letters used for transcription are given below with a few examples.

i	(fine, tit, bi)	[fɪ:ne, t ^h ɪt, bɪ·ʔ]
e	(dele, fedt, sne)	[de:le, fet, sne·ʔ]
ɛ	(sæbe, mælk, æg)	[se:bə, mɛɪʔk, ɛ·ʔk]
æ	(stave, værst, sal)	[sdæ:və, vænst, sæ·ʔɪ]
a	(nat, trække)	[nat, trage]
ɑ	(aften, mig, leg)	[ɑfdən, mɑɪ, ɑɪʔ]
ɑ	(varme, brand, barn)	[va:mə, branʔ, ba·ʔn]
y	(lyve, ny)	[ly:və, ny·ʔ]
ø	(løbe, dyppe, sø)	[lø:bə, døbə, sø·ʔ]
œ	(køn)	[kœnʔ]
œ	(gør, tørstig)	[gœr t ^h œrsdɪ]
u	(bule, skulder, jul)	[bu:le, sɡulɔ, ju·ʔɪ]
o	(hoste, sko)	[ho:sdə, sɡo·ʔ]
ɔ	(kåbe, stå)	[k ^h ɔ:bə, sdɔ·ʔ]
ʌ	(slot, øjne)	[slʌt, ʌne]
ɒ	(går, får)	[ɡɒ·ʔɜ], [fɒ·ʔɜ]

After the grouping confusion matrices were formed: one for each age-group and mode of presentation. As the confusions were distributed in the same manner for the two age-groups the answers from these were pooled for the subsequent examinations. From

these matrices it followed that with a few exceptions the vowels in the incorrect answers were either identical to the stimulus vowel or (especially for a low S/N) with approximately the same first formant frequency. This observation is illustrated in Figs. 1-6.

Results obtained at different S/N cannot immediately be compared since they come from different stimulus words. In order to make such a comparison meaningful, more information about the stimulus material is given in TABLE 2. This table shows for the three lists A, B, and C separately the total number of times a word containing a specific vowel was presented to a subject (e.g. in list A two words have /i/ as stressed vowel; these words are presented twice to seventeen subjects; consequently, the total number of presentations is 68). TABLE 2 also shows the total number of answers containing the correct vowel. The results are given for auditive (upper figures for each vowel) as well as audio-visual presentation (lower figures).

In Figs. 1-6 the number of mistakes between vowels is indicated by different types of lines according to the signature. An arrow on the line points towards the incorrect vowel.

Fig. 1 shows confusions observed at the auditive presentation of list A (S/N = -20 dB). It follows that mistakes occur mainly between vowels with approximately identical first formant frequencies.

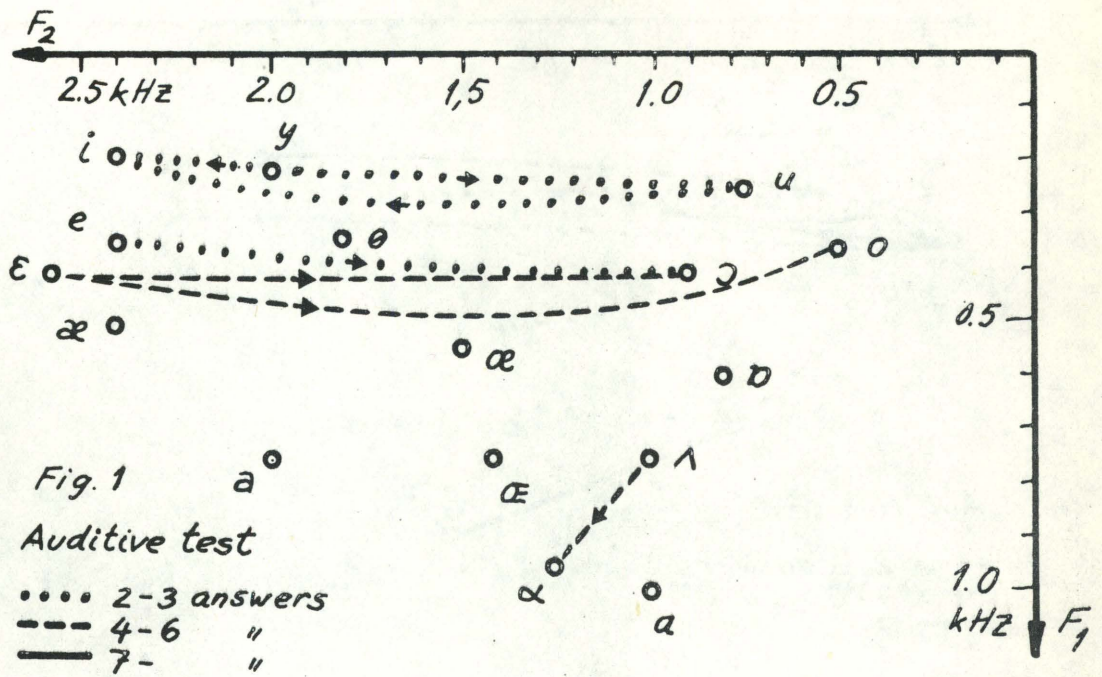
Fig. 2 shows that for the same words and the same S/N but with audio-visual presentation no confusions occur between rounded and unrounded vowels. This is in part surprising since no marked difference in lip configurations is present when pronouncing Λ and α . The explanation seems to be that all words containing Λ in list A also contain a bilabial stop consonant. This is an easily detected visual cue and therefore the wrong answers observed in the auditive test will not appear in the audio-visual test as they contain no such consonants.

TABLE 2

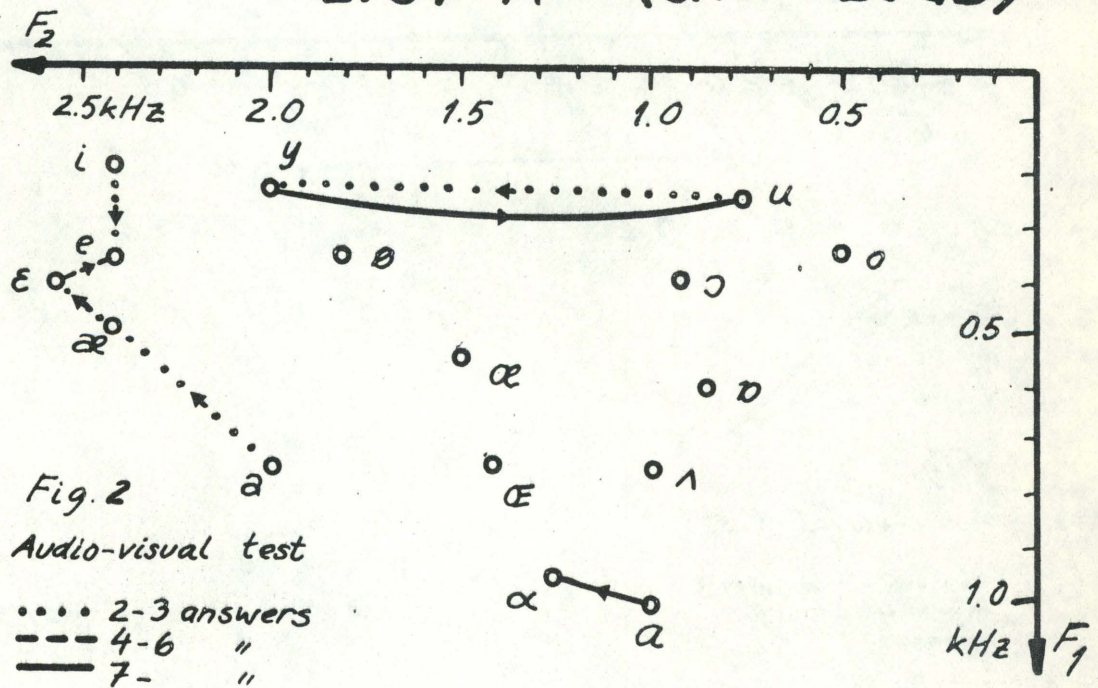
LIST A (S/N = -20 dB) LIST B (S/N = -10 dB) LIST C (S/N = 0 dB)

	NUMBER OF: PRES. ANSW.		NUMBER OF: CORR. CORR. ANSW. VOWEL		NUMBER OF: PRES. ANSW.		NUMBER OF: CORR. CORR. ANSW. VOWEL		NUMBER OF: PRES. ANSW.		NUMBER OF: CORR. CORR. ANSW. VOWEL	
i	68	8	1	6	34	16	1	7	68	66	58	65
	68	50	22	48	34	24	6	24	68	68	68	68
e	68	19	8	15	136	99	55	77	136	114	67	98
	68	55	37	53	136	131	114	127	136	131	114	122
ɛ	170	34	14	23	136	84	42	62	102	97	87	91
	170	144	124	139	136	127	108	124	102	102	101	102
æ	68	17	5	14	34	21	1	20	68	63	50	63
	68	58	40	55	34	30	11	20	68	67	58	67
a	34	5	1	4	68	39	7	20	34	30	29	29
	34	24	4	16	68	62	54	58	34	34	34	34
α	102	35	29	32	34	29	6	17	68	65	62	65
	102	96	92	94	34	28	15	28	68	68	66	68
ɑ	34	3	2	2	34	22	17	22	68	67	62	66
	34	16	5	6	34	31	30	31	68	68	68	68
ɣ	34	8	2	3	-	-	-	-	34	34	34	34
	34	25	10	14	-	-	-	-	34	34	34	34
ø	-	-	-	-	34	20	6	15	-	-	-	-
	-	-	-	-	34	30	23	30	-	-	-	-
œ	-	-	-	-	-	-	-	-	34	31	28	28
	-	-	-	-	-	-	-	-	34	33	32	32
œ	-	-	-	-	-	-	-	-	-	-	-	-
	-	-	-	-	34	27	16	24	-	-	-	-
u	68	5	0	0	68	53	46	48	34	23	17	19
	68	38	27	31	68	66	60	61	34	32	31	32
o	34	8	2	5	68	59	47	57	34	33	31	33
	34	28	21	25	68	67	62	67	34	34	34	34
ɔ	34	16	5	12	34	30	23	30	34	32	21	30
	34	29	19	28	34	34	34	34	34	32	27	31
ʌ	102	26	13	18	102	74	58	67	102	99	81	90
	102	94	76	93	102	101	96	101	102	99	91	99
ɐ	34	13	8	12	34	29	21	26	34	33	33	33
	34	28	22	27	34	32	31	32	34	34	34	34

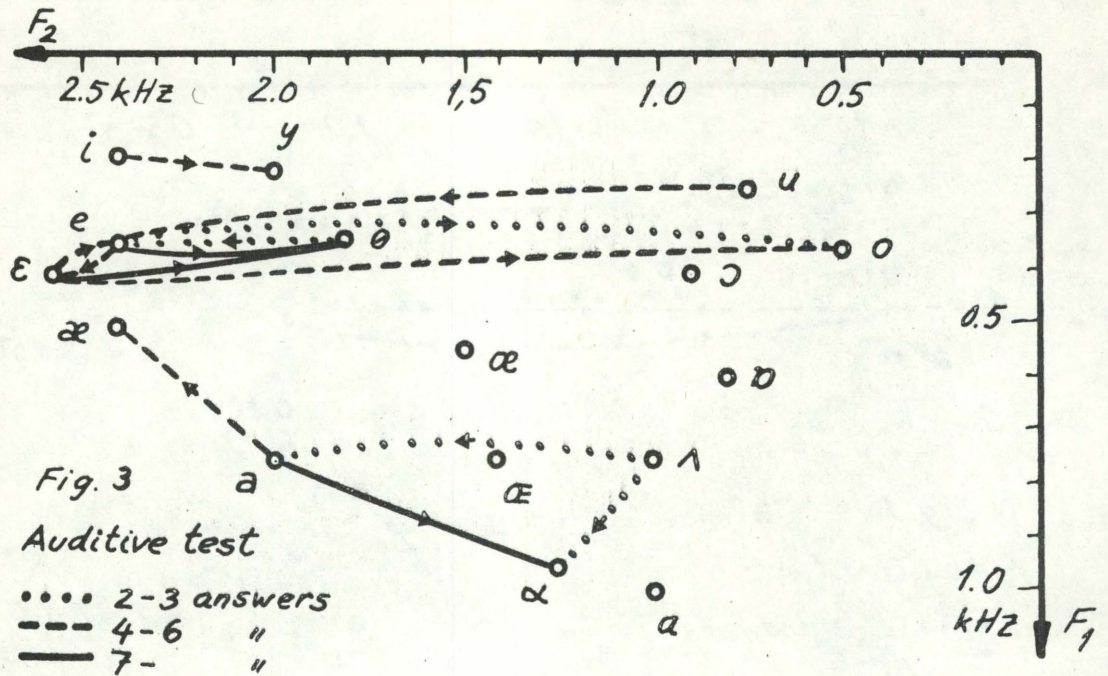
LIST A (S/N = -20 dB)



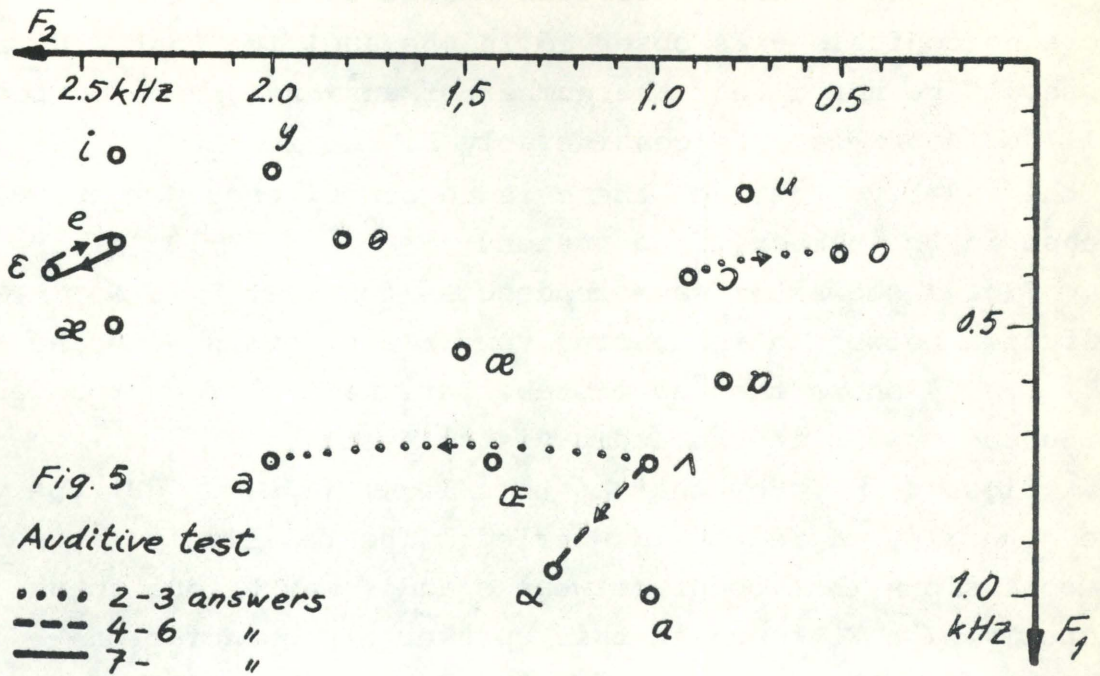
LIST A (S/N = -20 dB)



LIST B (S/N = -10 dB)



LIST C (S/N = 0 dB)



LIST C (S/N = 0 dB)

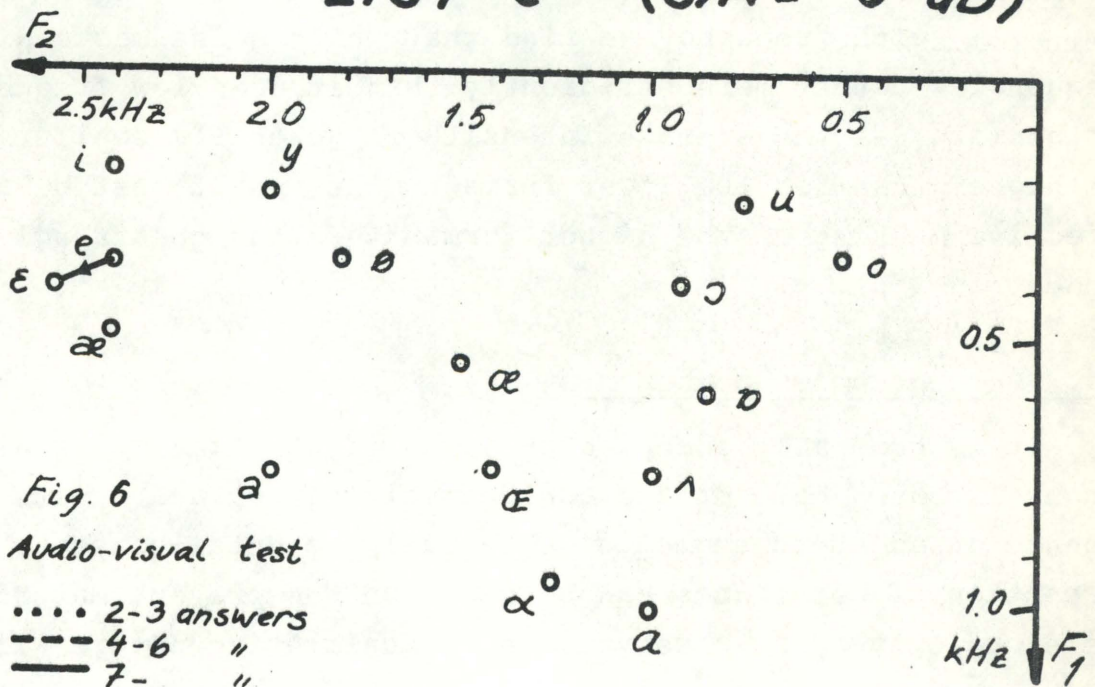


Fig. 2 also shows a certain amount of confusion among vowels where no confusion was observed in the auditive test. However, it should be noted that the number of answers given (correct as well as incorrect) are considerably higher in the A-V test than in the auditive test and there is no significant percentual increase in the number of confusions.

Fig. 3 shows the same tendencies for list B ($S/N = -10$ dB): confusion between neighbouring vowels and vowels with the same F_1 .

Fig. 4 shows that no mistakes are made between rounded and unrounded vowels in the audio-visual test.

Figs. 5-6 show that at this level ($S/N = 0$ dB) the vowels are generally correctly identified. The relatively large number of confusions that occur between e and ϵ may be due to an exceptionally high F_1 found in this speaker's pronunciation of e.

No diagrams are shown for the test with list D ($S/N = +10$ dB) since no confusion between vowels was observed.

The results obtained here fit with the well-known masking properties of white noise. As masking is roughly determined by the intensity level within a critical band and the critical bandwidth grows with frequency, we find that white noise masks high frequencies rather more efficiently than it does low frequencies. Furthermore, as the average intensity is generally smaller for the higher than for the lower formants the result must be an effective masking of the higher formants (although strongly dependent on S/N).

2.4. Perception of consonants

It is generally accepted that the identification of a consonant is based both on the consonantal segment and on its influence on adjacent segments. Obviously no detailed study of the perception of consonants can be based on the present material. Therefore, only a few observations of qualitative nature will be given here.

2.4.1. Labials

The most notable observation is the difference in discrimination of bilabials (p, b, m) as well as labio-dentals (f, v) found for the two types of presentation. While the discrimination of these consonants (especially that of f and notably for $S/N \leq -10$ dB) is remarkably poor in the auditive tests, the discrimination in the audio-visual tests is very high, even for $S/N = -20$ dB. Generally the audio-visual detection of f, v, and m (word initially) is almost perfect, i.e. even in the incorrect answers these consonants are always found in the correct positions when they were present in the stimulus word. The consonants p and b are often mutually confused for $S/N \leq -10$ dB but if a bilabial stop occurs in the stimulus word then a bilabial stop will be found even in incorrect answers.

This finding agrees well with earlier observations. The auditive discrimination between voiced and voiceless consonants in white noise is rather good (see e.g. Miller and Nicely 1954) and so is the visual detection of bilabials and labiodentals (see e.g. Woodward and Barber 1960). A generalization of these findings to audio-visual perception agrees with the above mentioned observations (remembering that Danish b is voiceless).

2.4.2. Voiceless fricatives

Another interesting observation is that voiceless fricatives are very hard to detect in white noise. This is not very surprising, but furthermore voiceless fricatives (especially s) as well as the affricated stop t occur in the incorrect answers where the stimulus word had no such consonants. This must be due to the pronounced similarity in acoustic quality between white noise and fricative sounds.

2.4.3. Voiced consonants

At a $S/N = -20$ dB wrong answers to stimulus words containing voiced consonants (m, n, ŋ, l, v, j, ð, r, ɣ) generally con-

tain voiced consonants but with several confusions between them. At more favourable S/N the consonants are discriminated more accurately. The material gives no support to the theory that nasals are detected as a separate group as found by Miller and Nicely (1955). One reason for this could be that in American English vowels are generally strongly nasalized in nasal surroundings.

2.4.5. Stops

The stops in Danish are all voiceless and the difference between b, d, g and p, t, k is mainly one of aspiration (t is also somewhat affricated). The material from the auditive tests gives no indication that confusion within these groups³ should be more likely than confusion between the groups or even with other consonants.

3. Conclusion

The results obtained in this paper show in agreement with earlier findings (see e.g. O'Neill 1954, Sumbly and Pollack 1954) that the visual signal of a speaker's face considerably improves the detection of certain speech segments especially when the signal to noise ratio is unfavourable. The improvement is particularly conspicuous in the detection of bilabials and labio-dentals but also in separating rounded vowels from unrounded. The results are obtained from a discrimination test of isolated words and it may be expected that the influence of the visual signal is less pronounced in the perception of running speech, since e.g. the syntactical structure of preceding strings will make detection of certain segments redundant. And furthermore the articulation will generally be less distinct.

3) viz. the group b, d, g and the group p, t, k.

It is also shown that the masking of white noise is not uniformly distributed over the frequency range and that the masking is one of higher frequency components mainly. Thus if the background noise is to give approximately the same decrease in redundancy for all components of the speech signal, another type of masking sound with less intensity in the higher frequency region must be used.

References

- ✓ Ewertsen, H.W., Nielsen, H.B.,
and Nielsen, S.S. 1970: "Audio-visual Speech Perception",
Acta Otolaryngologica 263, p. 229-
230.
- O'Neill, John J. 1954: "Contributions of the Visual Com-
ponents of Oral Symbols to Speech
Comprehension", Journal of Speech
and Hearing Disorders 19, p. 429-
439.
- Miller, G.A. and
Nicely, P.E. 1955: "An Analysis of Perceptual Confu-
sions Among Some English Conso-
nants", Journal of the Acoustical
Society of America 27, p. 338-352.
- Sumby, W.H. and
Pollack, I. 1954: "Visual Contribution to Speech In-
telligibility in Noise", Journal
of the Acoustical Society of Ame-
rica 26, p. 212-215.
- Woodward, M.F. and
Barber, C.G. 1960: "Phoneme Perception in Lipreading",
Journal of Speech and Hearing Re-
search 3, p. 212-222.