

AD-785 073

HADAMARD TRANSFORM FOR SPEECH WAVE
ANALYSIS

Hozumi Tanaka

Stanford University

Prepared for:

Advanced Research Projects Agency

August 1972

DISTRIBUTED BY:

NTIS

**National Technical Information Service
U. S. DEPARTMENT OF COMMERCE
5285 Port Royal Road, Springfield Va. 22151**

DISCLAIMER NOTICE

THIS DOCUMENT IS THE BEST
QUALITY AVAILABLE.

COPY FURNISHED CONTAINED
A SIGNIFICANT NUMBER OF
PAGES WHICH DO NOT
REPRODUCE LEGIBLY.

AD 785073

HADAMARD TRANSFORM FOR
SPEECH WAVE ANALYSIS

BY

HOZUMI TANAKA

Reproduced from
best available copy.

SUPPORTED BY

ADVANCED RESEARCH PROJECTS AGENCY

ARPA ORDER NO. 457

AUGUST 1972

COMPUTER SCIENCE DEPARTMENT

School of Humanities and Sciences

STANFORD UNIVERSITY



Reproduced by
NATIONAL TECHNICAL
INFORMATION SERVICE
U S Department of Commerce
Springfield VA 22151



AUGUST 1972

COMPUTER SCIENCE DEPARTMENT
REPORT NO. CS-307

HADAMARD TRANSFORM FOR
SPEECH WAVE ANALYSIS

by

Hozumi Tanaka

Abstract: Two methods of speech wave analysis using the Hadamard transform are discussed. The first method is a direct application of the Hadamard transform for speech waves. The reason this method yields poor results is discussed. The second method is the application of the Hadamard transform to a log-magnitude frequency spectrum. After the application of the Fourier transform the Hadamard transform is applied to detect a pitch period or to get a smoothed spectrum. This method shows some positive aspects of the Hadamard transform for the analysis of a speech wave with regard to the reduction of processing time required for smoothing, but at the cost of precision. A formant tracking program for voiced speech is implemented by using this method and an edge following technique used in scene analysis.

The views and conclusions contained in this document are those of the author and should not be interpreted as necessarily representing the official policies, either expressed or implied, of the Advanced Research Projects Agency.

This research was supported in part by the Advanced Research Projects Agency of the office of the Secretary of Defence (SD-183)

Reproduced in the USA. Available from the National Technical Information Service, Springfield, Virginia 22151.

Acknowledgements

The author would like to express deepest thanks to Dr. A.L. Samuel and Mr. J.K. Sizerz. Dr. A.L. Samuel suggested this research to the author. Mr. K.L. Sizerz developed a speech wave analysis system which has been used for this research as a basic tool. The author is grateful to Dr. R. Thosar for his many valuable suggestions to implement a formant tracking program. The author takes this opportunity of expressing his appreciation to Prof. John McCarthy for the opportunity to work at Artificial Intelligence Project as a visiting scholar from 1971 August to 1972 August. I would also like to thank Mr. and Mrs. Paul and Mr. J.E. Gips for their help in preparing this report.

Table of contents

1.	Introduction.	(1)
2.	Direct application of the Hadamard transform for speech wave analysis.	(3)
2.1	Definition of sequency and sequency power.	
2.2	Strong shift-sensitivity of the Hadamard sequency spectrum.	
2.3	Difficulties in calculating shift invariants for sequency power.	
3.	nasstron technique.	(12)
3.1	Outline.	
3.2	Pitch detection.	
3.3	Smoothing of a spectrum.	
3.4	A formant tracking program as an application of the nasstron technique.	
3.4.1	Logical structure of a formant tracking program.	
3.4.2	example.	
4.	Conclusion.	(26)
5.	Appendix.	(28)
5.1	Filtering in the sequency domain.	
6.	References.	(31)

1 Introduction.

Recently people in various fields have paid much attention to the Hadamard transform and have obtained results from its application in such fields as filter design, voice analyzer/synthesizer and multiplexer equipment [1]. The Hadamard (or discrete Walsh) transform is one of the orthogonal transforms using discrete Walsh functions and has a fast algorithm similar to the Fourier transform [2],[3].

There are many reasons why the Hadamard transform is attractive. Two major reasons are as follows. First, the Fast Hadamard Transform algorithm -FHT- uses only add / subtract operation. Multiplication is not necessary for the FHT. This makes the calculation of the FHT extremely simple and faster than the Fast Fourier Transform -FFT. In the Fourier transform case one needs multiplication for the sine-cosine coefficients, sometimes even with irrational numbers. The FHT offers quite a simple and an appropriate algorithm when using a digital computer.

Secondly, the discrete Walsh functions give us a general basis for signal analysis, namely the concept of sequency rather than that of frequency. The sequency of discrete Walsh functions is defined by one half of the average number of zero crossings per second. This concept enables us to replace the concept of frequency of the sine-cosine functions.

Because of this feature of the Hadamard transform one may well think of the possibility that all problems which have been solved using the Fourier transform might be re-interpreted by the Hadamard transform. Furthermore, one might hope for some interesting new discoveries since the Hadamard transform might reveal some new aspect of the problem concerned.

From this optimistic standpoint, the author has attempted an analysis of the speech wave using the Hadamard transform. Similar attempts have been made in the past [4], and they have suggested some possibilities about the application of the Hadamard transform to the speech wave by showing some correspondence between the frequency spectrum and the sequency spectrum. This report will show two methods of speech wave analysis using the Hadamard transform, the direct and the indirect methods. These two methods show both the advantages and disadvantages of the Hadamard transform for speech wave analysis.

Section 2 will explain the direct method by an application of the Hadamard transform to a speech wave. This method gives a poor result due to the strong shift sensitivity of the Hadamard sequency spectrum. Some shift invariant terms of the sequency power spectrum are known but they are complicated to calculate or too simple to provide enough information. A few experimental results are shown in this section to demonstrate these facts.

section 3 will explain the indirect method named the "hapstrum" technique. The hapstrum technique is a similar technique to the so called cepstrum technique [5] except that the FFT is applied to the log-magnitude frequency spectrum. This technique is indirect in the sense that at first the FFT (not FHT) is applied to a short span of a speech wave and then the FHT is used to detect the pitch period or to get a smoothed spectrum. This technique shows some positive aspect of the Hadamard transform for the analysis of a speech wave with regard to smoothing of a spectrum. Some experimental results will demonstrate this.

A formant tracking program has been implemented using the technique of an edge follower in scene analysis combined with the hapstrum technique. However, such an approach always contains a pitfall, namely the problem of wrong way entrance. This will be discussed in section 3.3.

Finally, in section 4 a tentative evaluation will be made of the Hadamard transform for analyzing speech waves.

Reproduced from
best available copy.

Direct application of the Hadamard transform to speech wave analysis.

In this section the Hadamard transform will be directly applied to a speech wave to get the sequency power spectrum. The existence of some correspondence between frequency spectrum and sequency spectrum has been reported on [4]. As a given vocalic sound can be characterized by the location of its first three formant frequencies, it is worth investigating the existence of formant "sequencies" in the Hadamard sequency spectrum instead of formant frequencies. A few experiments will demonstrate poor results and the reason will be discussed.

2.1 Definition of sequence and sequency.

The definition of sequence was introduced by H. F. Harnuth [3] and it gives a new basis from which to investigate the characteristic of signals. A sequence number of a Walsh function is defined by the number of sign changes per unit time. Let $N = 2^n$ consecutive real numbers $a(j)$, $0 \leq j < N$ be represented by a $1 \times N$ matrix $[a(j)]$. The Hadamard transform of $[a(j)]$ is

$$[A(k)] = (1/N)[a(j)]H(n) \quad (1)$$

where the $N \times N$ Hadamard matrix $H(n)$ is defined recursively in the equation (2).

$$H(n+1) = \begin{vmatrix} H(n) & d(n) \\ H(n) & -d(n) \end{vmatrix} \quad (2)$$

$$H(0) = \begin{vmatrix} 1 \end{vmatrix}$$

Each column of $d(n)$ represents one of discrete Walsh functions [7]. The examples of $H(3)$ and $H(4)$ are shown in the Fig. 2.1.

The sequency is defined by the average number of zero crossings per unit time divided by 2. Let $p(k)$ be the number of sign changes (zero related to frequency or sequency) it is desirable to calculate defined by eq.(3)

$$q(k) = \lfloor (p(k)+1)/2 \rfloor \quad (3)$$

where $\lfloor x \rfloor$ represents the largest integer which does not exceed x (see $H(3)$ of the Fig. 2.1). It is known that $p(k)$ takes on all values between zero and $N-1$ and $q(k)$ takes all values between zero and $N/2$.

$$H(3) = \begin{pmatrix}
1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
1 & -1 & 1 & -1 & 1 & -1 & 1 & -1 \\
1 & 1 & -1 & -1 & 1 & 1 & -1 & -1 \\
1 & -1 & -1 & 1 & 1 & -1 & -1 & 1 \\
1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 \\
1 & -1 & 1 & -1 & -1 & 1 & -1 & 1 \\
1 & 1 & -1 & -1 & -1 & -1 & 1 & 1 \\
1 & -1 & -1 & 1 & -1 & 1 & 1 & -1
\end{pmatrix}$$

↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓
0 7 3 4 1 6 2 5 --- sequence of each column: p
0 4 2 2 1 3 1 3 --- sequence of each column: q,

$$H(4) = \begin{pmatrix}
1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
1 & -1 & 1 & -1 & 1 & -1 & 1 & -1 & 1 & -1 & 1 & -1 & 1 & -1 \\
1 & 1 & -1 & -1 & 1 & 1 & -1 & -1 & 1 & 1 & -1 & -1 & 1 & 1 \\
1 & -1 & -1 & 1 & 1 & -1 & -1 & 1 & 1 & -1 & -1 & 1 & -1 & 1 \\
1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & 1 & 1 & 1 & 1 & -1 & -1 \\
1 & -1 & 1 & -1 & -1 & 1 & -1 & 1 & 1 & -1 & -1 & 1 & -1 & 1 \\
1 & 1 & -1 & -1 & -1 & -1 & 1 & 1 & 1 & 1 & -1 & -1 & 1 & 1 \\
1 & -1 & -1 & 1 & -1 & 1 & 1 & -1 & 1 & -1 & 1 & 1 & -1 & -1 \\
1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & -1 & -1 \\
1 & -1 & 1 & -1 & 1 & -1 & 1 & -1 & -1 & 1 & -1 & 1 & -1 & 1 \\
1 & 1 & -1 & -1 & 1 & 1 & -1 & -1 & -1 & 1 & 1 & 1 & 1 & 1 \\
1 & -1 & -1 & 1 & 1 & -1 & -1 & 1 & 1 & -1 & -1 & 1 & -1 & 1 \\
1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & -1 & 1 & 1 & 1 & 1 & 1 \\
1 & -1 & 1 & -1 & -1 & 1 & -1 & 1 & -1 & -1 & 1 & -1 & 1 & -1 \\
1 & 1 & -1 & -1 & -1 & -1 & 1 & 1 & -1 & -1 & -1 & 1 & -1 & -1 \\
1 & -1 & -1 & 1 & -1 & 1 & 1 & -1 & 1 & -1 & -1 & 1 & -1 & 1
\end{pmatrix}$$

Fig. 2.1 Examples of the Hadamard matrix.

Let us introduce two notations, $A(c, p(k))$ and $A(s, p(k))$ for $A(k)$.

$$A(k) = \begin{cases} A(c, q(k)) & \text{if } p(k) \text{ is even} \\ A(s, q(k)) & \text{if } p(k) \text{ is odd.} \end{cases}$$

In analogy of frequency power spectrum, sequency power spectrum is defined as follows.

$$\begin{aligned} & \sum_{q=0}^2 A(c, q) \\ & \sum_{q=0}^2 A(c, q) + \sum_{q=1}^2 A(s, q) \quad 0 < q < N/2 \\ & \sum_{q=1}^2 A(s, N/2) \end{aligned} \quad (4)$$

The Parseval's relation is preserved on the coefficients $A(k)$ and $a(k)$.

$$(1/N) \sum_{k=0}^{N-1} a^2(k) = \sum_{q=0}^2 [A(c, q) + \sum_{q=1}^{(N/2)-1} [A(c, q) + A(s, q)] + A(s, N/2)] \quad (5)$$

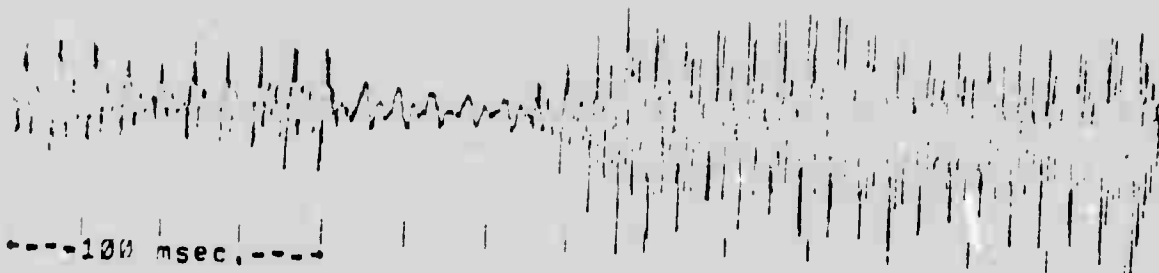
2.2 Strong shift-sensitivity of the Hadamard sequency spectrum.

It is interesting to investigate the formant structure in the sequency spectrum of a speech wave. The most convenient way is to change many consecutive spectra into a visual form, that is the sonogram of sequencies. A short time span (12.8 ms.) of a digitized speech wave (sample rate = 22000 Hz.) is directly transformed into sequency spectrum. Then the log-magnitude of this spectrum is taken. Many short time sequency spectra are calculated in this way, are accumulated, and eventually output to a video screen.

Experimental results are shown in Fig. 2.2. The upper part shows a speech wave to be analyzed, the middle part a sonogram of frequency spectra of this speech wave and the lower part a sonogram of sequency spectra. It is easy to see that the sonogram of sequency spectra (the lowest one) is rougher than that of the frequencies (the middle one). The formant sequency structure is not clear and it appears to be very difficult to build a speech wave analysis system based on the extraction of formant components using the Hadamard sequency spectrum.

The reason why the sonogram of sequency spectra becomes so rough and irregular is made clear by the following experiment. The Hadamard sequency spectrum is calculated for a fixed time span (12.8 msec. long) of a speech wave. The time span is shifted right by 100 microseconds for each successive calculation of the sequency spectrum. In other words, calculation of a sequency spectrum is made each 100 microsecond time-shift. A Frequency spectrum of the Fourier transform is calculated in the same way to make comparison with sequency spectrum. The results are shown in Fig. 2.3.

Reproduced from
best available copy.



Speech wave



Frequency spectrum



Sonogram spectrum

Fig. 2.2 Sonogram of sequency and frequency spectra.

Frequency spectrum
Frame No. 1



Sequency spectrum
Frame No. 1



0 1000 2000 3000 (hz)

No. 2

0 1000 2000 3000 (cps)

No. 2



No. 3



No. 3



Fig. 2,3 Strong shift-sensitivity of the Hadamard sequency spectrum. Each frame is calculated each 100 microsecond time-shift.

From Fig. 2.3 we can easily understand that although the time-shift is limited to this small value, the shape of consecutive sequency spectra changes rapidly. The location of a peak which appears to represent a formant component changes drastically in the next sequency spectrum. One cannot expect these rapid changes from observation of the original speech wave since the speech wave does not appreciably change its shape during 100 microseconds. In contrast, in the Fourier case, a frequency spectrum does not change its shape so much during 100 microseconds. This strong shift-sensitivity of the Hadamard sequency spectrum causes the irregularity or rough pattern of a sequency sonogram and makes impossible the application of the pitch-synchronous method.

The strong time-shift sensitivity of a sequency spectrum also can be explained theoretically. Pichler [6] shows the Hadamard sequency spectrum is invariant under the dyadic time-shift:

$[b(j)]$ is obtained by the dyadic time-shift t

$$[b(j)] = [a(j \circledast t)]$$

where $j \circledast t$ stands for component-wise modulo two addition (no carry) for the binary representation of j and t . Pichler's result is written as follows.

$$B^2(c, q) + B^2(s, q) = A^2(c, q) + A^2(s, q) \quad (6)$$

Unfortunately the Hadamard sequency spectrum is not invariant under circular time-shift of the input $[a(j)]$. If $[a(j)]$ is shifted by t circularly forming $[c(j)]$ we obtain:

$$[c(j)] = [a((j + t))]$$

where $((j + t))$ is the principal value of $j + t$ modulo N . In general

$$C^2(c, q) + C^2(s, q) \neq A^2(c, q) + A^2(s, q) \quad (7)$$

The experiment shown in Fig. 2.3 is not the case of circular time-shift but one can easily understand that the relation of eq (7) causes the strong shift sensitivity in the Hadamard sequency spectrum. Note that in contrast to the Hadamard sequency spectrum a frequency spectrum of the discrete Fourier transform is invariant under circular time-shift since absolute value of a shift operator is one.

2.3 Difficulties in calculating shift invariants for the Hadamard transform.

Some attempts have been made to define circular time-shift invariants for the Hadamard transform. Ohnson [7] has defined a complete set of circular time-shift invariants of the Hadamard transform and also has

shown intermediate forms which are invariant to both circular time-shift and dyadic time-shift. For more detailed derivation of a complete set of circular time-shift invariants and its intermediate forms see [7].

As a first step, consider intermediate forms, a set $\{P(k)\}$ which is a sum of groups of components in $[A(k)]$ squared such that

$$\begin{aligned} P^2(0) &= A^2(0) \\ P^2(1) &= A^2(1) \\ P^2(2) &= A^2(2) + A^2(3) \\ &\dots\dots\dots \end{aligned} \tag{8}$$

In general

$$P^2(m) = \sum_k A^2(k)$$

where $2^{m-1} \leq k < 2^m$ for $1 \leq m \leq n$.

Examples of calculations of a set $\{P\}$ for various input waves are shown in Fig. 2.4. In the figure the short time span of the speech wave for the Hadamard transform is fixed to 12.8 msec. Each component of a set $\{P\}$ is shown as a function of time in the Fig. 2.4. Overlap of the time span for the next Hadamard transform is 6.4 msec. The case of a sinusoidal wave indicates the filtering characteristic of a set $\{P\}$ because the position of each peak moves to the left as k increases in $P(k)$. In other words, the smaller the value of k in eq (8), the more likely it is that the component $P(k)$ will pass the higher frequency component since frequency increases with time passing in the original input wave. However, as the band of each filter is determined by the number N , which is the dimension of an array $[A(k)]$, we lose flexibility. Although the calculation of a set $\{P\}$ from N components of $[A(k)]$ is straightforward, we can get only $1 + n (= \log_2 N)$ components of P . For instance, if $N = 256$ one can get only 9 components of P and one of them is d.c. component. This means a great deal of information reduction is made and it is doubtful if a set $\{P\}$ contains enough information to perform speech wave analysis.

Ohnsorge has defined another complete set of the Hadamard transform which has exactly $(N/2) + 1$ invariants for a circular time-shift. (The discrete Fourier transform -DFT- gives a $(N/2) + 1$ point spectrum.) However it is not a straightforward way to calculate the invariants since it includes many matrix multiplications. According to [7] if we let $\{J\}$ be a quadratic invariant set of the Hadamard transform, then

In the case when $N = 8$

$$\begin{aligned} J^2(0) &= A^2(0) \\ J^2(1) &= A^2(1) \\ J^2(2) &= A^2(2) + A^2(3) \\ J^2(3) &= A^2(4) + A^2(6) - A(4)A(7) + A(5)A(6) \\ J^2(4) &= A^2(5) + A^2(7) + A(4)A(7) - A(5)A(6) \end{aligned} \tag{9}$$

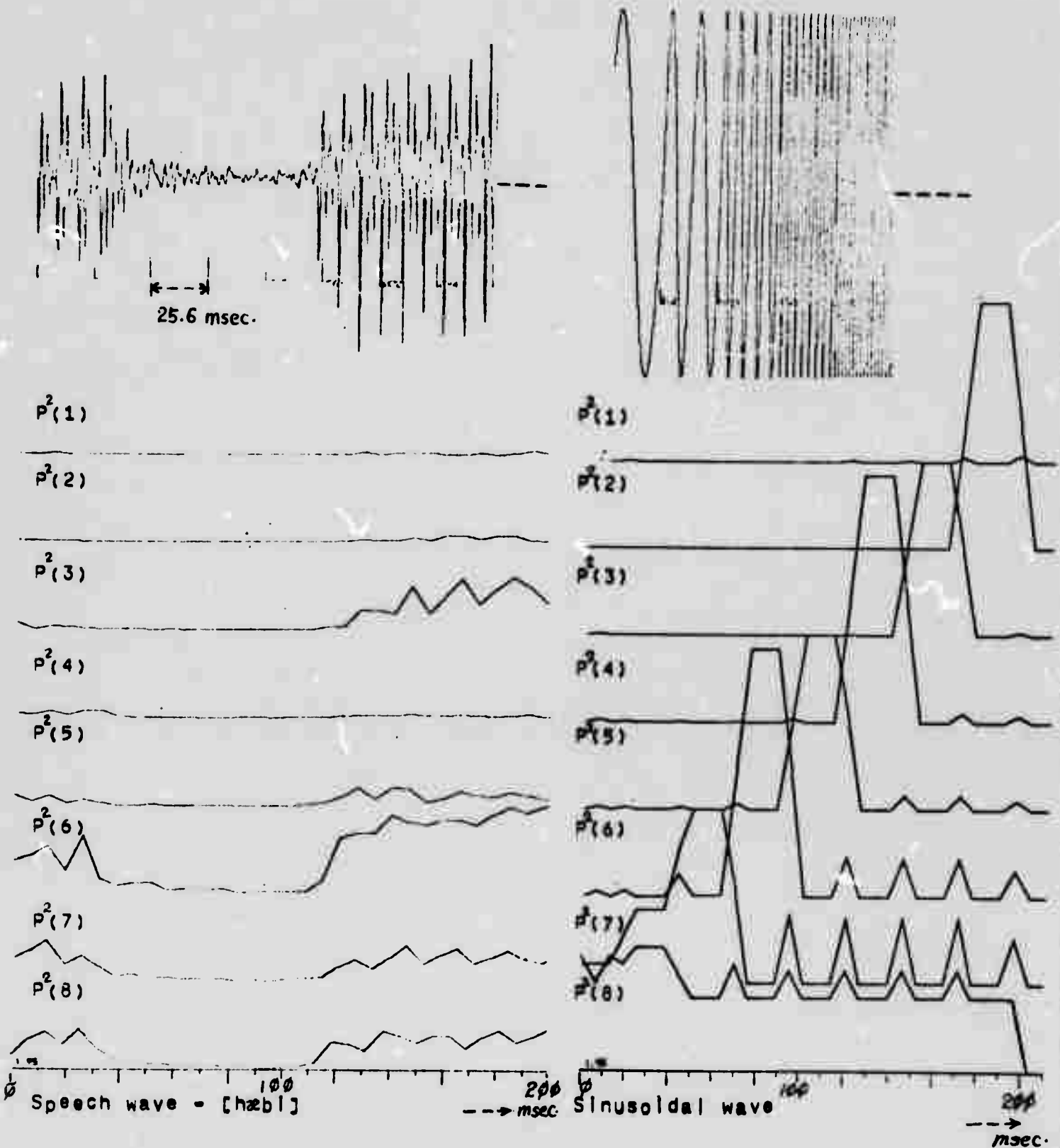


Fig. 2.4 Calculation of eq (8) for various input waves.

Although there is no explanation about how these terms (J) are related to frequency or sequU Ohnsorg's [7] Invariants. As Ohnsorg suggests that the prominent energy lines of the discrete Fourier spectrum tend to be exaggerated in the quadratic spectrum (J).

N. Ahmed et al [8] found an efficient algorithm to calculate these terms. However multiplication by an irrational number is included in the algorithm and it is more complicated than that of Hadamard transform.

3 cepstrum technique.

In this section the "hapstrum" technique is introduced. The hapstrum technique is a similar technique to the cepstrum technique except that the inverse fast Hadamard transform -IFHT- is applied to the log-magnitude frequency spectrum and the output is called "hapstrum." This technique is indirect in the sense that at first the FFT (not FHT) is applied to a short time span of a speech wave to obtain the spectrum and then the FHT is used to extract pitch period or to get smoothed spectrum. The strong time-shift sensitivity of the Hadamard transform is removed by the first application of Fourier transform to speech waves.

This technique illustrates a positive aspect of the Hadamard transform for the analysis of a speech wave, especially with regard to the smoothing of a spectrum. A formant tracking program has been implemented using this technique.

3.1 Outline.

To show both the advantages and disadvantages of the hapstrum technique we will depict the outline of both the cepstrum and the hapstrum techniques. Although there is more than one definition of the cepstrum technique we give a typical application in the upper part of Fig. 3.1. The hapstrum technique is shown in the lower part of Fig. 3.1.

From Fig. 3.1, one can easily understand the difference between both techniques. The frequency spectrum of a short time span of a speech wave filtered by a Hamming window is obtained by the discrete Fourier transform -DFT. Then the log-magnitude of this spectrum is taken. After the processing, in the case of the cepstrum technique the inverse discrete Fourier transform -IDFT- and DFT are applied to get pitch period and smoothed spectrum. On the other hand, in the case of the hapstrum technique the IDFT and DFT are replaced by the IFHT and FHT, respectively. A hapstrum, which is ordered in sequence (not sequency), is obtained by the IFHT of a log-magnitude spectrum. From the replacements one gets the advantage of the fast calculation of the Hadamard transform. Due to the elimination of linear filtering, computing cost is even further reduced by the method.

Let us note that in the cepstrum case after the application of the inverse discrete Fourier transform we need low-pass filtering of the log-magnitude of the discrete Fourier transform. By means of low-pass filtering a smoothed spectrum is obtained due to the elimination of the fine structure of the spectrum. This is accomplished by multiplying the cepstrum by a low-pass filter function.

In contrast to the cepstrum technique, the hapstrum technique uses an ideal filter as a low-pass filter in the sequency domain of the hapstrum. Therefore one needs no multiplication to cut higher

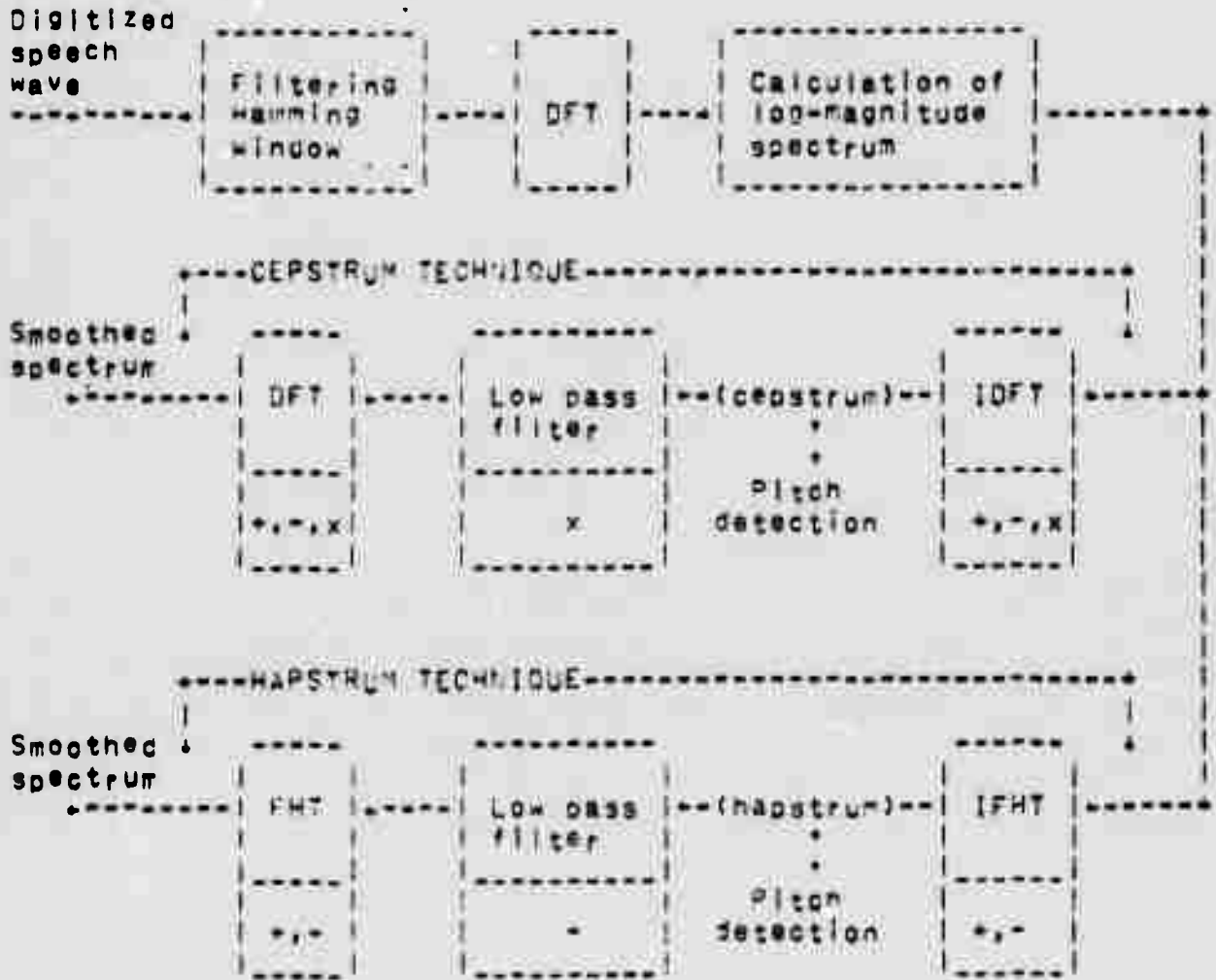


Fig. 3.1 The outline of the cepstrum technique and the harstrum technique.

sequence components. The higher sequence components are simply made zero. This also reduces computing cost (the symbols +/- and x in the figure indicate the necessity of add/subtract operations or multiplications).

From the author's experience the calculation of the FHT is ten times as fast as that of the FFT. This suggests that by using the hapstrum technique we can make the calculation of spectrum smoothing at most three times as fast as that of spectrum smoothing using the cepstrum technique.

However, we should be aware that smoothing by the cepstrum gives us a better approximation for an original log-magnitude spectrum in the sense of least-square error criterion and that smoothing by the hapstrum degrades resolution of peak position of log-magnitude spectrum. The theoretical reason for this will be discussed in section 3.3.

3.2 Pitch detection.

To extract a pitch period we have to take a sufficient time-span of a speech wave to calculate a log-magnitude spectrum, namely long enough to include at least two glottal pulses.

In our experiments the duration is taken to be 25.6 msec corresponding to 512 samples of a digitized speech wave since the sampling rate of a speech wave is 20704 Hz.

Fig. 3.2-a shows a series of cepstrum plots. A series of cepstrum are calculated for each consecutive segment of speech wave one half of which overlaps the previous segment. In the case of the cepstrum, to get a higher resolution 512 zeros are added to the next 512 samples of a digitized speech wave. This means the IDFT and DFT are calculated on 1024 points.

Fig. 3.2-b shows a series of hapstrum plots. The hapstrum is calculated under the same condition as the cepstrum of Fig. 9. To calculate a hapstrum we do not add zero to the next 512 samples of a speech wave, since one cannot get higher resolution of the hapstrum by adding zeros (see 3.3 in this section). If 512 zeros are added to the next 512 samples of a digitized speech wave one will get a hapstrum such that the component of the sequence (not sequency) $2i$ and $2i + 1$ becomes the same value, where i is a positive integer. In other words a hapstrum of a speech wave segment with added zeros is easily calculated from one without added zeros. This special feature of the Hadamard transform is utilized by the smoothing of the log-magnitude spectrum in the next section. The proof is shown in the APPENDIX in more a generalized form.

Comparing Fig. 3.2-a with Fig. 3.2-b, we observe that in the cepstrum a sharp peak appears at approximately 4.5 msec but in the case of the

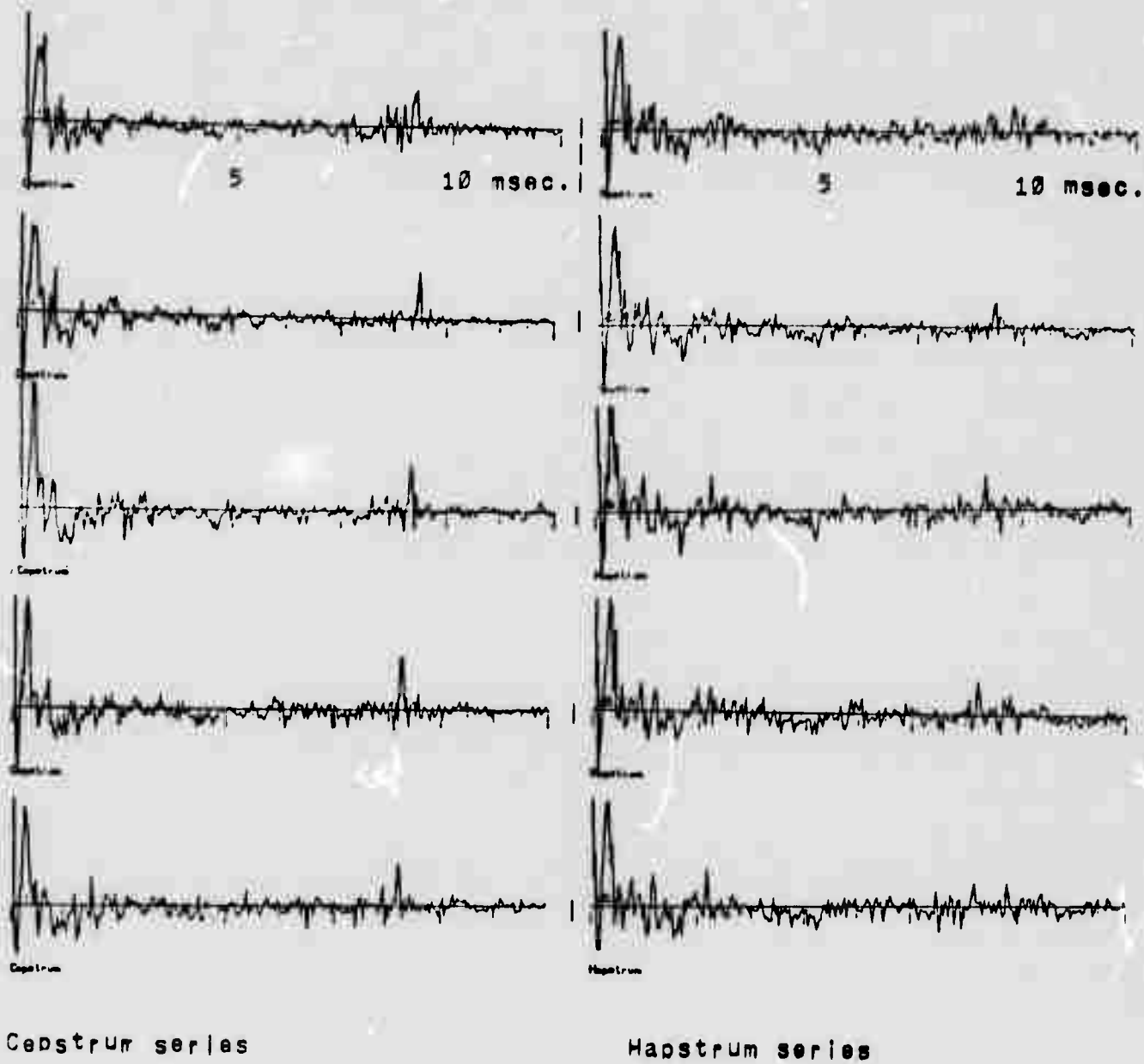


Fig. 3.2 An example of cepstrum series and the hapstrum series.

hapstrum the peak is not so sharp.

As pitch period is determined by the location of this sharp peak it will be more difficult to extract pitch period from the hapstrum. The cepstrum is superior to the hapstrum in so far as pitch detection is concerned.

3.3 Smoothing of a spectrum.

The resonant frequencies of a vocal tract are called formant frequencies, and the first three characterize a given vocalic sound. Therefore, for the analysis of a speech wave it is very important to extract these three frequencies. The procedure is called a formant tracker.

There exist two methods to extract formant frequencies. One is the linear prediction method which extracts these frequencies directly from a given speech wave. In other words the formant extraction is performed on the time domain. Atal et al [9] reported good results from the method. The other method is based on peak detection in the frequency domain of a speech wave [10].

In Fig. 3.3, an example of a log-magnitude spectrum of a short-time speech wave is shown. Fig. 3.3 suggests that a log-magnitude spectrum is composed of its spectral envelope and the spectral fine-structure. Roughly speaking, the spectral fine-structure has equidistant peaks at the pitch (fundamental) frequency and harmonics.

As formant frequencies are represented by several prominent peaks in a spectral envelope, smoothing or elimination of fine-structure is important. The cepstrum technique is one of the prominent methods for it [5],[10] but its computational speed is rather slow since it includes three FFT calculations. With respect to this point the hapstrum is faster at the cost of degradation of accuracy of peak positions in smoothed log-magnitude spectrum (see eq (11) and (12)).

The hapstrum technique is based on filtering in the sequence domain as shown in Fig. 3.1. The output of the IFHT, which is in sequence order, is a hapstrum. After the detection of a pitch period from the hapstrum, all hapstrum components with more than a fixed sequence number are set to zero. This is accomplished by an ideal filter on the sequence domain. With the higher sequence components cut from the hapstrum the FHT is used to recover an original log-magnitude spectrum.

The determination of the cut-off sequence number is as follows. Let an hapstrum be represented by an array $[h(j)]$ of dimension $N (= 2^i)$ and the location of a peak caused by pitch frequency be an index number k of $[h(j)]$. If $N/(2^{i-1}) \leq k < N/(2^i)$ then cut-off sequence number r is

$$r = N/(2^{(i-1)}) \quad (10)$$

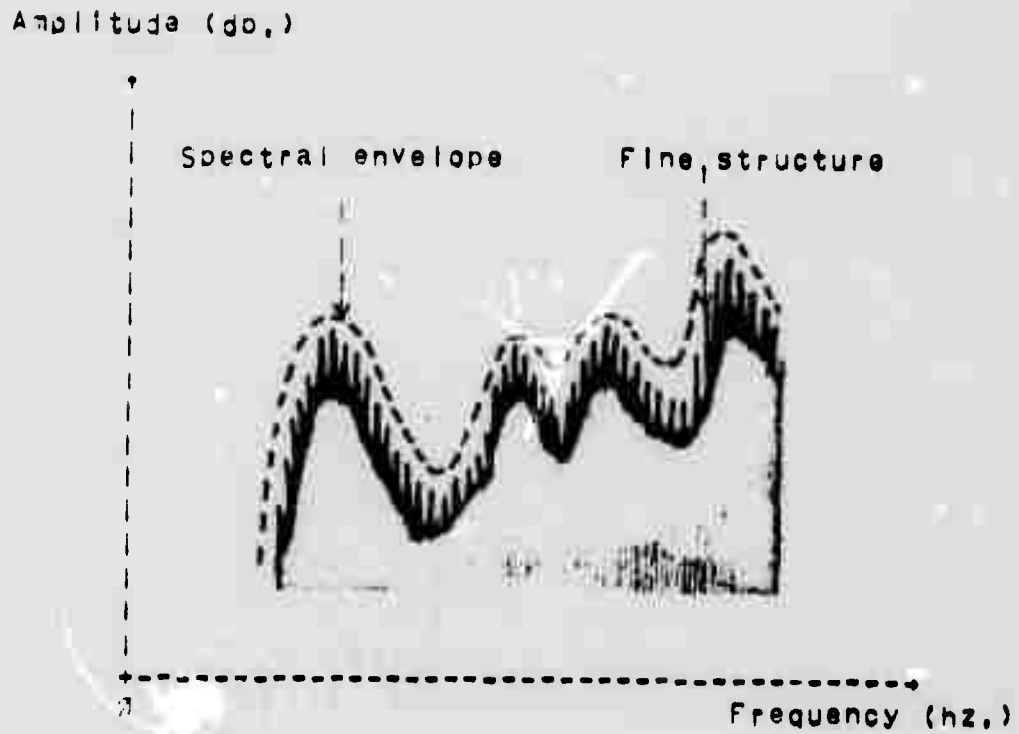


Fig. 3.3 Spectral envelope and spectral fine structure of log-magnitude spectrum of a speech wave.

Consider the meanings of filtering by an ideal filter in the sequence domain. Let an array $[a(j)]$ of dimension $N (= 2^n)$ be a digitized signal in which all components such that $N/2 \leq j < N$ are set to zero. By the application of the FHT including sequence ordering the array $[a(j)]$ is transformed into an array $[B(k)]$ such that each adjacent component becomes the same, namely:

$$\begin{aligned}
 B(0) &= B(1) \\
 B(2) &= B(3) \\
 &\dots\dots\dots \\
 B(N-2) &= B(N-1)
 \end{aligned}
 \tag{11}$$

Furthermore, when all components such that $N/(2^2) \leq j < N$ are set to zero the array $[a(j)]$ is transformed into such an array $[B(k)]$ by the application of the FHT including sequence ordering

$$\begin{aligned}
 B(0) &= B(1) = B(2) = B(3) \\
 B(4) &= B(5) = B(6) = B(7) \\
 &\dots\dots\dots \\
 B(N-4) &= B(N-3) = B(N-2) = B(N-1)
 \end{aligned}
 \tag{12}$$

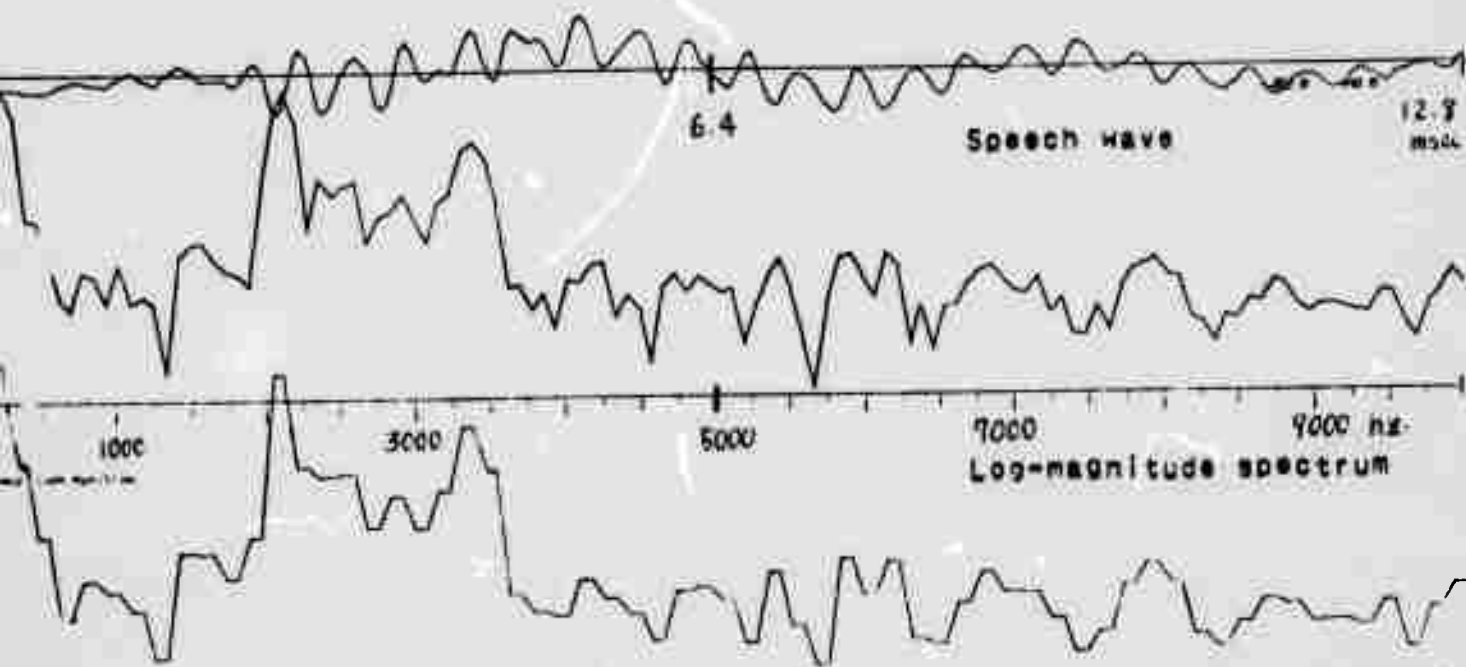
Eq (11) and (12) are generalized in (A) and (B) of APPENDIX.

Both equations suggest that if $[B(k)]$ is plotted as a function of array index k the curve becomes flat as the value of each adjacent component is the same. Because of this flattening effect it degrades resolution of peak positions in the array $[B(k)]$. To demonstrate this an example is shown in Fig. 3.4. A segment of a speech wave is shown and is analyzed by the hapstrum technique. The two lower curves represent the log-magnitude spectrum and the smoothed result by the hapstrum technique.

The smoothed log-magnitude spectrum in Fig. 3.4 demonstrates the smoothing effect stated before by eq (11). Many sharp maxima and minima caused by glottal pulses (or pitch frequency) in the original log-magnitude spectrum are diminished in the smoothed spectrum. From the author's experience the number of peaks is decreased to one half that of the original.

The important question is whether or not the prominent peaks caused by resonance of a vocal tract are preserved by the smoothing. From the example shown in Fig. 3.4 we can see that the hapstrum smoothing technique gives good smoothing with regard to preserving the first three formants.

Fig. 3.5 gives another example which suggests the formant components are preserved after the smoothing by the hapstrum technique. The upper is a sonogram of spectra without smoothing and the lower is a sonogram of smoothed results using the hapstrum technique.



Smoothed spectrum by
the hapstrum technique

Fig. 3.4 Speech wave, the log-magnitude spectrum and
the smoothed spectrum by the hapstrum technique.

a)



b)



Fig. 3.5 Sonograms of log-magnitude spectra and their smoothed spectra. The upper is a sonogram of log-magnitude spectra and the lower is that of smoothed spectra.

3.4 A formant tracking program as an application of the hapstrum technique.

A formant tracking program has been implemented using the hapstrum technique and an edge follower technique as used in scene analysis [11]. In principle, the formant tracking program presented here accepts any kind of smoothing technique such as cepstrum or inverse filtering [12].

Edge followers were first implemented to recognize objects in a scene. An edge follower detects a position where sharp change of contrast occurs and follows it successively. A sonogram is just such a scene with formant trajectories represented as dark stripes. By detecting dark stripes we find the locations of peaks in a spectrum since a sonogram is represented as a sequence of spectra.

There are many difficulties in implementing a formant tracking program based on an edge follower. One problem is that a formant trajectory is not a straight line, but is curved. Some of the edge followers have treated objects composed only of straight lines, such as cubes. This limitation can be of use to an edge follower. For instance we can prevent the following of the wrong path by using the criterion of curvature. We also can forecast the existence of edge, which is hard to detect because of noise, by using straight line interpolation methods. As the production of a speech wave is a dynamic and stochastic process, the human speech wave contains much noise.

A second problem is that it is very difficult to decide a formant frequency from local information. A wide range of overlap exists between the region of the first formant frequency and that of the second, also between the second formant frequency and the third. In the case of a male voice, the first formant frequency ranges from 220 hz to 900 hz., the second from 550 hz. to 2700 hz, and the third from 1100 hz. to 3000 hz.

A third problem is that if we see a sonogram in a microscopic way there exist too many peaks to discriminate the formant components. It is desirable to have a technique to eliminate trivial peaks while preserving prominent peaks caused by the first three formants. Cepstrum is such a technique. Rabiner and Schafer [10] have implemented a formant tracking program based on the cepstrum technique. As their method makes frame-by-frame decisions for the first three formant frequencies, they use only local information in a sonogram. It is desirable to utilize more global information.

Markel [12] has developed a very good technique for getting a smoothed spectrum based on the idea of linear prediction method. He calls it inverse filtering and has developed a formant tracking program [13] which uses information from the previous frame when it is difficult to determine the first three formant frequencies.

The formant tracking program explained here follows Markel's approach but with a backtracking mechanism to recover if a wrong path is followed. If decisions are made frame by frame there is no wrong way entrance problem. Even if we make a wrong decision in a frame, the effect does not propagate to the next. However, if we use the information from just the previous frame the effect of a wrong decision will propagate. To cope with this situation it is necessary to have a recovery technique which utilizes more global information.

3.4.1 Logical structure of a formant tracking program.

Our formant tracking program is composed of four modules named PEAK DETECTOR, CANDIDATE SELECTOR, TRACKER, and RECOVERY. General flow of the program is shown in figure 3.6.

A. PEAK DETECTOR.

PEAK DETECTOR accepts a digitized speech wave of a vocalic sound, calculates a smoothed spectrum by using the hamstrum technique, and determines peaks. It should be noted that the hamstrum technique is used to decrease the processing time required for smoothing. It can easily be replaced by another technique such as inverse filtering or the cepstrum technique.

B. CANDIDATE SELECTOR.

For each region of the first three formant frequencies, CANDIDATE SELECTOR selects at most three candidates from many peaks detected by PEAK DETECTOR and orders them by amplitude of peaks. The third candidate whose amplitude is 7.5 db less than that of the second candidate is removed by the ordering process. These candidates selected are accumulated and are used by TRACKER and RECOVERY. This routine reduces the search space.

C. TRACKER.

TRACKER takes the results from CANDIDATE SELECTOR and makes a tentative decision for the first three formant frequencies. At first TRACKER looks for a reasonable place to track. There exists a region within which an overlap of two formant components never occurs. In the case of a male voice, only the first formant exists between 220 Hz. and 300 Hz. and only the second and the third formant exist between 520 Hz. and 1100 Hz. and between 2700 Hz. and 3000 Hz. If the first candidate for a formant frequency is within the first region it is reasonable to assume that this is the peak caused by the formant. After making an initial selection TRACKER begins tracking forward or backward.

TRACKER uses two criteria to determine formant frequencies of the next frame. Basically, TRACKER uses a criterion of minimum shift of peak position from one frame to the next. This nearest neighbour

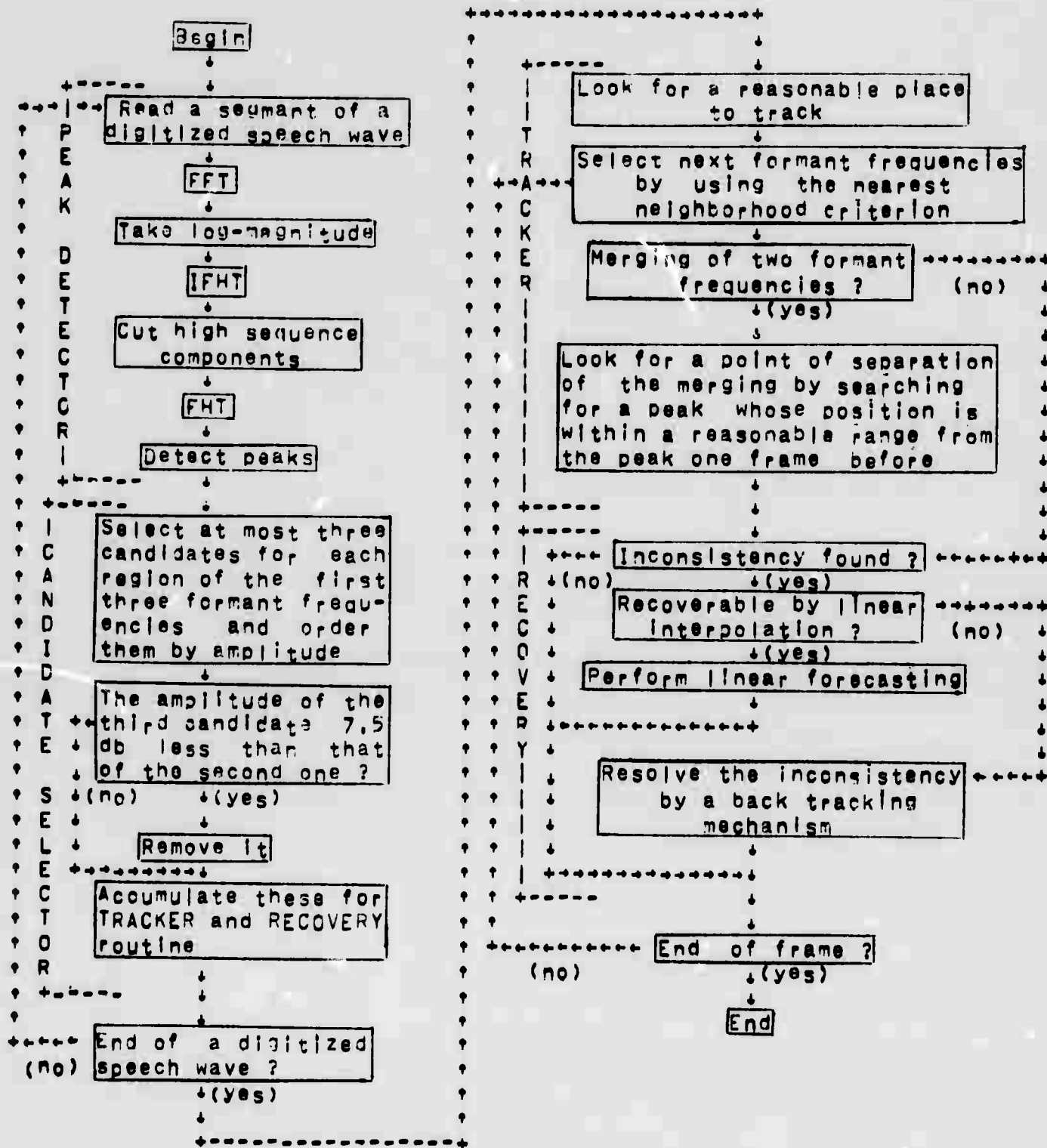


Fig. 3.6 General flow diagram of a formant tracking program.

criterion is used as long as merging of two formant frequencies does not occur. As soon as merging of two formant frequencies occurs, TRACKER makes use of the other criterion to look for a point of separation after merging. After a tentative selection of the next formant components using the first criterion, TRACKER looks for a peak whose position is within a reasonable range from the peak one frame before. If TRACKER can find such a peak it will select the peak as the point of separation of the two trajectories, otherwise the two trajectories remain merged. A wrong decision by TRACKER is corrected by the RECOVERY routine.

D. RECOVERY.

RECOVERY works when some inconsistency is recognized by a formant tracking program. An inconsistency is a discontinuity or a sharp change in following a formant trajectory.

There are two major reasons why TRACKER follows a path that has a sharp change from one frame to the next. The first reason is that a prominent peak caused by a formant component is often lost in a spectrum because of the stochastic movement of glottal pulses. This results in a discontinuity in a formant trajectory if the trajectory is selected in a microscopic way. This can be resolved by using the neighborhood information. A FORECASTER works in this case.

The second reason is that a wrong decision has been made in the past by TRACKER and a wrong path has been followed as a formant trajectory. For example a formant tracking program has mistaken the first formant for the second and the trajectory suddenly enters into the region where the second formant does not exist; namely the region between 220 hz. and 550 hz. The other typical example is the following of a wrong path which is not a formant trajectory and eventually disappears. These are corrected by using a backtracking mechanism in the RECOVERY routine. In the previous example, after the RECOVERY routine has recognized an error, a trajectory followed as the second formant is replaced by the first formant trajectory. Then another peak is selected for the second formant frequency by RECOVERY. The routine extends the new second formant trajectory backward by using TRACKER, and FORECASTER. This is an example of how the back tracking mechanism works. We can see that that a recovery process by back tracking has to have a recursive structure, but in our case the depth of recovery is limited to one.

3.4.2 Example.

An example of formant tracking obtained from the program is shown in Fig. 3.7. A spoken sentence is "we were away." more than 90% of running time is devoted to PEAK DETECTOR and CANDIDATE SELECTOR.

Reproduced from
best available copy.

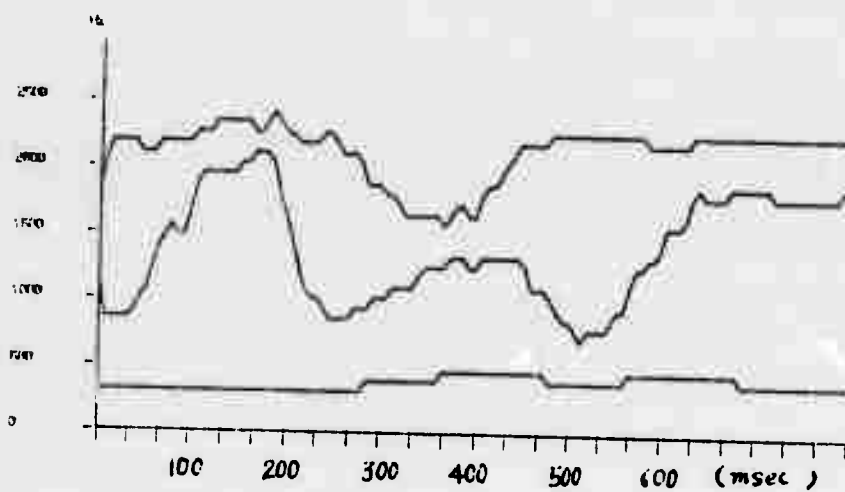


Fig. 3.7 An example of the first three formant trajectories for a sentence of "We were away".

Conclusion.

In the paper we have discussed both the advantages and disadvantages of the Hadamard transform, as compared to the Fourier transform, for a speech wave analysis. The experiments in section 2 reveal that application of the Hadamard transform directly to speech waves yields poor results, as it fails to extract important features. The smaller the number of features necessary to accurately represent a speech wave, the better it is. In the Fourier case, in almost all cases for a vocalic sound a speech wave is represented by the first three formant frequencies and the pitch (fundamental) frequency. Only four parameters are needed. However in the Hadamard frequency spectrum, we cannot observe any typical features because of the strong time-shift sensitivity which makes it impossible to apply even a pitch synchronous method. In other words, typical features which are recognizable in the Fourier case are averaged and are scattered away in a wide range of a frequency power spectrum. Some of the experiments in section 2.2 demonstrate it.

Time-shift invariants for the Hadamard frequency spectrum are known. One of these defined by eq. (8) does not bear enough information to perform a speech wave analysis since from a digitized speech wave composed of 256 points, we get only 9 components. Although each of these components has near relationship with an output from a filter bank, its frequency band is determined by the number of points transformed. Ohnsorg [7] has defined another complete set of the Hadamard transform which has exactly the same number of components as a Fourier frequency spectrum and is invariant under a circular time-shift. Ahned et al [8] found an algorithm to calculate these terms, however multiplication by an irrational number is included and is more complicated than that of the fast Hadamard transform. As Ohnsorg suggests that the prominent energy line of the Fourier spectrum tends to be exaggerated, it is desirable to calculate Ohnsorg's invariants for a speech wave.

In section 3 the hapstrum technique is introduced. This technique is similar to the so called ceostrum technique except that the FHT is applied to the log-magnitude frequency spectrum. After the application of the Fourier transform the Hadamard transform is applied to detect a pitch period or to get a smoothed spectrum. This technique shows some positive aspect of the Hadamard transform for the analysis of a speech wave with regard to the reduction of the processing time required for smoothing. Good smoothing makes it easy to extract the first three formant frequencies in a spectrum. We should note that smoothing by the hapstrum is obtained at the cost of accuracy in determining peak position of a smoothed spectrum. This is explained by eq (11) or (12) in section 3.3. We can conclude that precise formant frequencies are obtained by the ceostrum technique at the cost of processing time, while a reduction of processing time is obtained by the hapstrum technique at the cost of accuracy in determining the formant frequencies. However, it is often true that

to detect a peak caused by a pitch period, is difficult even in the case of a male voice. The author's original optimistic standpoint was that the Hadamard transform might reveal some new aspect of speech waves. However the only gain found from using the Hadamard transform was the reduction of processing time required for smoothing, and this was obtained at the cost of precision.

A formant tracking program using an edge follower has been described in section 3.4. While the algorithm is rather sophisticated, most of the time is still devoted to the smoothing and peak selection procedures.

Reproduced from
best available copy.

APPENDIX.

Let us define a few of the functions used here.

Function G(U) is defined as follows:

Let a binary representation of U or G(U) be

$$\begin{aligned}
 G(U) &= G = G_{n-1}G_{n-2} \dots G_1G_0 \\
 U &= U_{n-1}U_{n-2} \dots U_1U_0 \\
 U_i \text{ and } U_i &\in \{1,2\} \text{ (for } 1 \leq i \leq n-1) \\
 \text{where} \\
 G_{n-1} &= U_n \text{ XOR } U_1 \\
 G_{n-2} &= U_1 \text{ XOR } U_2 \\
 &\dots \dots \dots \\
 &\dots \dots \dots \\
 G_1 &= U_{n-2} \text{ XOR } U_{n-1} \\
 G_0 &= U_{n-1} \\
 \text{(XOR stands for exclusive-or)}
 \end{aligned}
 \tag{A-1}$$

Let an array [a(j)] be [a, f] such that
 [e, f] = [e0, e1, ..., em-1, f0, f1, ..., fm-1]
 and [a] = [a0, e1, ..., em-1]
 [f] = [f0, f1, ..., fm-1].

From the definition of the Hadamard transform

$$[A(j)] = (1/N)[e, f] \begin{vmatrix} H(n-1) & H(n-1) \\ H(n-1) & -H(n-1) \end{vmatrix}
 \tag{A-2}$$

where $N = 2^n$ and $n = l/2$.

5.1 Filtering on the sequence domain.

(A) If array [f] is $[f_0, \dots, f_{m-1}]$ then $A(k) = A(l)$ for $l = k + m$, where $0 \leq k \leq (N/2) - 1$, and the difference of sequence number between A(l) and A(k) is one.

Proof:

Suppose the sequence number of the k-th or l-th element of array [A(j)] is s or t, respectively. Then if the binary representation of k, l, s or t is

$$\begin{aligned}
 k &= k_{n-1} k_{n-2} \dots k_1 k_0 \\
 l &= l_{n-1} l_{n-2} \dots l_1 l_0 \\
 s &= s_{n-1} s_{n-2} \dots s_1 s_0 \\
 t &= t_{n-1} t_{n-2} \dots t_1 t_0
 \end{aligned}
 \tag{A-3}$$

$$\text{then } k = G(s) \text{ and } l = G(t) \quad (\text{see [3]}) \tag{A-4}$$

Since $0 \leq k \leq (N/2) - 1$ and $l = k + (N/2)$
 most significant binary digit k_{n-1} and l_{n-1} are
 $k_{n-1} = 0$
 $l_{n-1} = 1$ and
 $k_i = l_i$ for $i \neq n - 1$ (A-5)

From eq (A-4) and (A-5)

$$\begin{aligned} l_{n-1} &= t_0 \quad \text{XOR} \quad t_1 = 1 \\ k_{n-1} &= s_0 \quad \text{XOR} \quad s_1 = 0 \\ k_{n-2} &= s_1 \quad \text{XOR} \quad s_2 = t_1 \quad \text{XOR} \quad t_2 \\ &\dots\dots\dots \\ k_1 &= s_{n-2} \text{ XOR } s_{n-1} = t_{n-2} \text{ XOR } t_{n-1} \\ k_0 &= s_{n-1} = t_{n-1} \end{aligned} \quad (A-6)$$

We obtain the following relation from eq (A-6),

$s_l = t_l$ for $1 \leq l \leq n-1$, and

$$\begin{aligned} s_0 &= 0 \text{ and } t_0 = 1 \text{ (if } s_1 = 0) \\ s_0 &= 1 \text{ and } t_0 = 0 \text{ (if } s_1 = 1) \end{aligned} \quad (A-7)$$

Eq (A-7) implies that a s or t is in sequence.
 In other words the difference of sequence number between
 $A(k)$ and $A(l)$ is one.

Let $[A(j)]$ be $[E, F]$ where

$$[E, F] = [E_0, E_1, \dots, E_{m-1}, F_0, F_1, \dots, F_{m-1}] \quad (A-8)$$

From eq (A-2)

$$\begin{aligned} E_k &= [e](h(k)) + [f](h(k)) \\ F_k &= [e](h(k)) - [f](h(k)) \end{aligned} \quad (A-9)$$

where $(h(k))$ is the k -th column of matrix $H(n-1)$.

Since in our case $[f] = \overbrace{[0, 0, \dots, 0]}^m$

$E_k = F_k$, namely $A(l) = A(k)$ for $l = k + (N/2)$ Q.E.D.

(B) We can generalize the result of (A) further,
 Zero all components of array $[a(j)]$ such that
 $2^k \leq j \leq N - 1$ where $1 \leq k \leq n-1$, then

$$\begin{aligned} A(0) &= A(2^k) = A(2 \cdot (2^k)) = A(3 \cdot (2^k)) = \dots = A((2^{(n-k)} - 1) \cdot 2^k) \\ A(1) &= A(2^k + 1) = A(2 \cdot (2^k) + 1) = A(3 \cdot (2^k) + 1) = \dots = A((2^{(n-k)} - 1) \cdot 2^k + 1) \\ &\dots\dots\dots \end{aligned}$$

$$A(1) = A(2^k + 1) = A(2 \cdot (2^k) + 1) = A(3 \cdot (2^k) + 1) = \dots = A((2^{(n-k)} - 1) \cdot 2^k + 1)$$

.....

$$A(2^k) = A(2 \cdot 2^k - 1) = A(3 \cdot (2^k) - 1) = A(4 \cdot (2^k) - 1) = \dots = A(2^n - 1)$$

and in each group, for example $(A(1), A(2^k + 1), A(3 \cdot (2^k) + 1), \dots, A((2^{(n-k)} - 1) \cdot 2^k + 1))$, $2^{(n-k)}$ consecutive sequence numbers are included.

proof:

It is apparent from the recursive definition of the Hadamard transform matrix given in eq (1) and the proof given in (A).

References.

- [1] H.F. Harmuth: Application of WALSH functions in communication, IEEE Spectrum, Nov., 82-91, 1969.
- [2] H.F. Harmuth: TRANSMISSION OF INFORMATION BY ORTHOGONAL FUNCTIONS, Springer-Verlag, 1970.
- [3] H.C. Andrews: COMPUTER TECHNIQUES IN IMAGE PROCESSING, Academic Press, 1970.
- [4] S.J. Campanella and G.S. Robinson: Digital frequency decomposition of voice signals, Walsh Function Symp., Naval Res. Lab., 230-237, 1970.
- [5] A.M. Noll: Cepstrum pitch determination, J. Acoust. Soc. Amer., 41, 2, 293-309, 1967.
- [6] F. Pichler: WALSH functions and optimal linear systems, Walsh Function Symp., Naval Res. Lab., 17-22, 1968.
- [7] F.R. Ohnsorg: Spectral modes of the WALSH-HADAMARD transform, Walsh Function Symp., Naval Res. Lab., 55-59, 1971.
- [8] N. Ahmed, A.L. Abdussattar and K.R. Rao: Efficient computation of the WALSH-HADAMARD transform spectral modes, Walsh Function Symp., Naval Res. Lab., 276-279, 1972.
- [9] B.S. Atal and S.L. Hanauer: Speech analysis and synthesis by linear prediction of the speech wave, J. Acoust. Soc. Amer., 47, 2, 637-655, 1971.
- [10] L.R. Rabiner and R.W. Schafer: System for automatic formant analysis of voiced speech, J. Acoust. Soc. Amer., 47, 2, 634-648, 1970.
- [11] A. Herskovits and T.O. Elnford: On boundary detection, MIT Project MAC Artificial Intelligence Memo 183, July 1970.
- [12] J.D. Markel: Formant trajectory estimation from a linear least-squares inverse filter formulation, SCRL-Monograph No.7, Oct., 1971.
- [13] J.D. Markel: Automatic formant and fundamental frequency extraction from a digital inverse filter formulation, Conf. on Speech Comm. and Processing, 81-84, 1972.