

## XI. COGNITIVE INFORMATION PROCESSING\*

### Academic and Research Staff

Prof. M. Eden	Prof. T. S. Huang	C. L. Fontaine
Prof. J. Allen	Dr. R. R. Archer	E. R. Jensen
Prof. F. F. Lee	Dr. W. L. Black	A. J. Laurino
Prof. S. J. Mason	Dr. J. E. Green	D. M. Ozonoff
Prof. W. F. Schreiber	Dr. K. R. Ingham	E. Peper
Prof. D. E. Troxel	Dr. O. J. Tretiak	Sandra A. Sommers
	F. X. Carroll	

### Graduate Students

G. B. Anderson	R. V. Harris III	D. S. Prerau
T. P. Barnwell III	D. W. Hartman	G. M. Robbins
W. L. Bass	A. B. Hayes	S. M. Rose
B. A. Blesser	P. D. Henshaw	C. L. Seitz
J. E. Bowie	M. Hubelbank	D. Sheena
A. E. Filip	W-H. Lee	W. W. Stallings
A. Gabrielian	J. I. Makhoul	R. M. Strong
A. M. Gilkes	G. P. Marston III	Y. D. Willems
R. E. Greenwood	D. R. Pepperberg	J. W. Woods
E. G. Guttman	G. F. Pfister	I. T. Young
	R. S. Pindyck	

#### A. COMPUTER SYNTHESIZED HOLOGRAMS – TWO EXPERIMENTS

In a previous report,<sup>1</sup> we derived a new method for recording a complex function in the form of a real non-negative function. The new method has advantages over the conventional holographic method because it does not require a reference signal or a biasing constant in recording the complex wave front of an object. In this report we discuss two experiments in which computer-synthesized holograms are used to display continuous-tone pictures and three-dimensional objects.

##### 1. Fourier Transform Hologram of a Continuous-Tone Picture

In this experiment we want to generate a Fourier transform hologram for a continuous-tone picture by a computer. In an experiment like this the transparency of the picture is first scanned by a flying-spot scanner, and the data of the transmission of the transparency are then stored on magnetic tape. These data constitute the information of the object for the synthesized hologram. The data stored on the magnetic tape will be read by the computer when we execute the program that generates the data of the Fourier transform hologram. After the transmission of the transparency has been

---

\*This work was supported principally by the National Institutes of Health (Grants 5 PO1 GM14940-03 and 5 PO1 GM15006-02), and in part by the Joint Services Electronics Programs (U.S. Army, U.S. Navy, and U.S. Air Force) under Contract DA 28-043-AMC-02536(E) and the National Institutes of Health (Grant 5 TO1 GM-01555-02).

(XI. COGNITIVE INFORMATION PROCESSING)

transferred from the tape to the core memory of the computer, we use the fast Fourier transform program HARM to compute its Fourier transform.<sup>2</sup> The complex Fourier transform of the transparency is then converted to the amplitude transmission of the hologram by our new method. We shall briefly summarize the experimental procedure for obtaining the amplitude transmission of the hologram here.

Assume that the transmission of the picture is represented by a  $N \times N$  matrix. Before we compute the Fourier transform of this two-dimensional array, we put the  $N \times N$  matrix in a  $N \times 4N$  matrix with zeros added to fill in the remaining elements of the larger matrix. By so doing, we change the sampling rate along one of the directions in the Fourier transform domain. With the  $N \times 4N$  matrix we then compute the two-dimensional Fourier transform of the picture. Suppose that  $\{F(n, m)\}$  are the elements of the Fourier transform of the picture, and  $\{H(n, m)\}$  are the elements of the matrix representing the amplitude transmission of the hologram. Then  $H(n, m)$  is given by

$$H(n, 4m-4+k) = \begin{cases} \operatorname{Re} \{(-j)^{k-1} \cdot F(n, 4m-4+k)\} & \text{if } \operatorname{Re} \{(-j)^{k-1} F(n, 4m-4+k)\} > 0 \\ 0 & \text{otherwise} \end{cases}$$
$$k = 1, \dots, 4. \quad n, m = 1, \dots, N \quad (1)$$

In Eq. 1 the symbol  $\operatorname{Re} ( )$  denotes the real part of the enclosed complex number. The detail theory pertaining to the formula in Eq. 1 can be found in the previous report.<sup>1</sup>

There is one difficulty in making a Fourier transform hologram of a continuous-tone picture. Because the transmission of the picture is a real non-negative function, the DC value in its Fourier transform is so large that it usually exceeds the dynamic range that can be displayed by the flying-spot scanner. Unlike the Fourier transform of the image, which has only transparent line segments on an opaque background, the distortion on the low frequencies of the picture will degrade the quality of the image reconstructed from the hologram. To alleviate such a problem, we multiply the transmission of the picture by a pseudorandom phase function before its Fourier transform is computed. The function of the pseudorandom phase function is similar to that of a piece of ground glass in conventional holography. When the transparency of a picture is put in contact with a piece of ground glass, the roughness of the surface of the ground glass introduces a phase variation in the light wave passing through the transparency. As a result, the light distribution in the Fourier transform is more evenly distributed. Since the ground glass only changes the phase of the light in some random manner, we can model the phase variation by a random function. Because neither the photographic material nor the human observer can detect the phase variation in an image, the pseudorandom phase function can be used very conveniently to modify the structure of the Fourier transform of an image.

(XI. COGNITIVE INFORMATION PROCESSING)

Figure XI-1 is the hologram of the continuous-tone picture shown in Fig. XI-2. Figure XI-3 is the image reconstructed from the hologram on an optical bench. The speckles in the reconstructed image are due to three factors. First, when the coherent light is used to form images, the dust particles in the path of illumination will cause

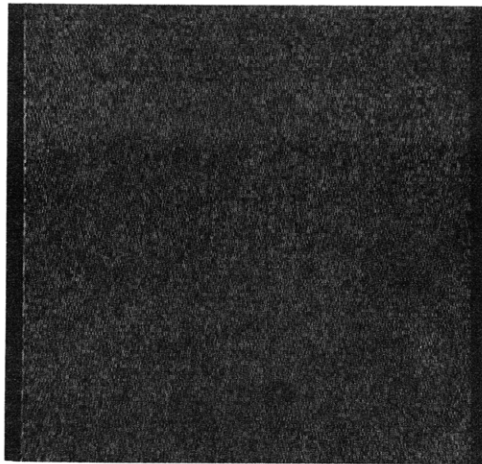


Fig. XI-1. Fourier transform hologram of the picture shown in Fig. XI-2. The transmission of the picture is multiplied by a pseudorandom phase function.



Fig. XI-2. Girl's face.



Fig. XI-3. Image reconstructed from the hologram in Fig. XI-1.

the light on the image plane to interfere to form specklelike patterns. The two other factors are due to the quality of the Fourier transform hologram itself. The multiplication of the pseudorandom phase function to the transmission of the picture gives a Fourier transform that is the convolution of the Fourier transforms of the two functions involved in the product. In most cases, the Fourier transform of the pseudorandom phase function has a very wide spread in its spatial frequencies. Therefore the

convolution of the Fourier transform of the picture with such a function will move some of the high spatial frequency information of the picture outside the aperture of the hologram. This will cause the appearance of the speckles in the reconstructed image. This effect has been observed in the conventional Fourier transform hologram, if the transparency is put in contact with a piece of ground glass and the size of the hologram is small.<sup>3</sup> Last, the quantized noise in the synthesized Fourier transform hologram also adds speckles to the reconstructed image shown in Fig. XI-3. There is a very interesting property of such noise. The study of noise in the Fourier transform domain with regard to spatial filtering has been investigated by Anderson.<sup>4</sup> Here we only wish to point out one observation that is relevant to the speckle noise in the reconstruction from a Fourier transform hologram. Suppose that the quantized noise in the Fourier transform hologram can be regarded as an additive noise. Then in the reconstruction the image of the original picture will also have an added noise term. Since the photographic film can only record the intensities of the light wave, the image recorded on the film is proportional to

$$I(x, y) = |f(x, y) \exp[j\theta(x, y)] + n(x, y)|^2. \quad (2)$$

Here,  $f(x, y)$  is the transmission of the transparency,  $n(x, y)$  is the Fourier transform of the quantized noise, and  $\theta(x, y)$  is the pseudorandom phase. In general,  $n(x, y)$  is a complex function. The function  $I(x, y)$  can be rewritten

$$I(x, y) = f(x, y)^2 + |n(x, y)|^2 + 2f(x, y)|n(x, y)| \cos [\theta(x, y) - \phi(x, y)], \quad (3)$$

where  $\phi(x, y)$  is the phase angle of the complex noise  $n(x, y)$ . The last term in Eq. 3 is a noise term whose strength depends on the signal itself. If we examine the image in Fig. XI-3, we notice similar characteristics of the speckle noise there.

It has been noted that the image reconstructed from a hologram of the type that we discuss here is less sensitive to dust and scratches on the surface of the hologram. The speckle noise, however, in the image reconstructed from such a hologram makes it less attractive for use as a means of transmitting high-quality television pictures.

## 2. Fourier Transform Hologram of a Three-Dimensional Object

In the course of making a synthesized hologram of a three-dimensional object, we need a method to compute the wave front of the object on a particular recording plane. Thus far, the method of finding the wave front of an object has been based on the point radiator approach.<sup>5</sup> That is, in finding the wave front of the object each point of the object is regarded as a point radiator. For a point located at  $(x_n, y_n, z_n)$  in space the wave front of the point can be approximated by

$$W(x, y; x_n, y_n, z_n) = S(x_n, y_n, z_n) \exp\left[j\pi\left(\frac{(x-x_n)^2 + (y-y_n)^2}{\lambda z_n}\right)\right]. \quad (4)$$

The position of the point and the recording plane for obtaining  $W(x, y; x_n, y_n, z_n)$  is illustrated in Fig. XI-4. The expression in Eq. 4 is the Fresnel approximation of a point source. If we sum up all waves of the points that constitute the object, we have the wave front needed to generate a synthesized hologram. If we want to make a synthesized Fourier transform of the object, however, the wavelets that make up the total wave front

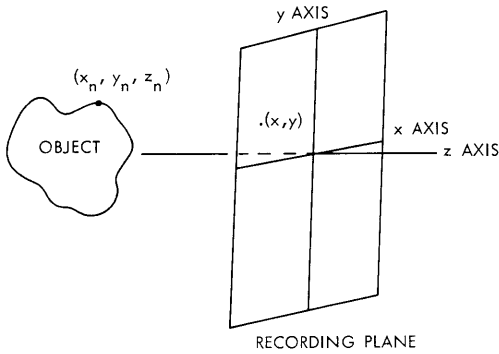


Fig. XI-4. Positions of object and recording plane.

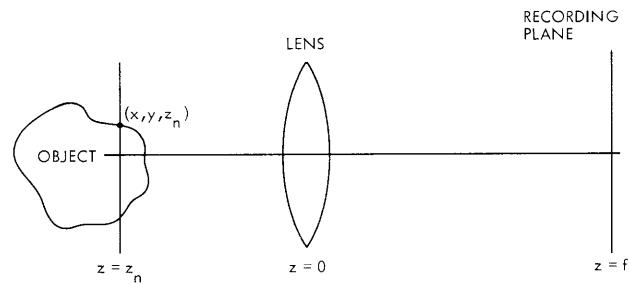


Fig. XI-5. Positions of object and recording plane for making a Fourier transform hologram.

will be slightly different from the expression in Eq. 4. Suppose that in making the Fourier transform hologram the object is located on one side of a lens. The arrangement for recording the hologram is shown in Fig. XI-5. The wave front on the back focal plane of the lens produced by a point of the object is approximately equal to

$$\begin{aligned} V(u, v; x, y, z_n) &= S(x, y, z_n) \iint_A \exp\left[j\pi\left(\frac{(t-x)^2 + (s-y)^2}{\lambda z_n} - \frac{j\pi(t^2 + s^2)}{\lambda f}\right)\right] \\ &\quad \cdot \exp\left[j\pi\left(\frac{(t-u)^2 + (s-v)^2}{\lambda f}\right)\right] dt ds \\ &= P(u, v; x, y, z_n) \cdot S(x, y, z_n) \\ &\quad \cdot \exp\left[j\pi\left(1 - \frac{z_n}{f}\right)\frac{(u^2 + v^2)}{\lambda f} - \frac{j2\pi(xu + yv)}{\lambda f}\right] \end{aligned} \quad (5)$$

In Eq. 5,  $f$  is the focal length of the lens, and the area of integration,  $A$ , depends on the size of the lens. If the object is small in comparison with the aperture of the lens, the function  $P(u, v; x, y, z_n)$  is approximately a constant. Equation 5 is the Fresnel-Kirchhoff diffraction formula with the Fresnel approximation. When there are more than one points at the same distance from the lens, the wave front resulting from all

## (XI. COGNITIVE INFORMATION PROCESSING)

of the points on that plane is given by

$$R(u, v; z_n) = \exp\left[j\pi(1-z_n/f)(u^2+v^2)/\lambda f\right] \sum_x \sum_y S(x, y, z_n) \exp(-2j\pi(xu+yv)/\lambda f). \quad (6)$$

The summation in Eq. 6 is the discrete Fourier transform of the two-dimensional pattern formed by the points lying on the same plane. The total wave front of the object is then equal to

$$F(u, v) = \sum_n R(u, v; z_n). \quad (7)$$

To show that the theoretical result derived for the wave front of the object is valid, we have performed a simple experiment. The three-dimensional object comprises two images located at two different planes. If we use the formula in Eq. 6 to find the wave front of the object, the wave front for this simple experiment is equal to

$$F(u, v) = R(u, v; z_1) + R(u, v; z_2). \quad (8)$$

In the actual experiment the values of  $z_1$  and  $z_2$  are  $f$  and  $0.99f$ , respectively. If we substitute these values in Eq. 7, the equation becomes

$$F(u, v) = \sum_x \sum_y S(x, y; z_1) \exp(-j2\pi(xu+yv)/\lambda f) + \exp[j\pi(u^2+v^2)/(100\lambda f)] \\ \cdot \sum_x \sum_y S(x, y; z_2) \exp(j2\pi(ux+vy)/\lambda f), \quad (9)$$

where  $f$  is equal to 500 mm and the wavelength  $\lambda$  is chosen to be  $6328 \text{ \AA}$ . The nonlinear phase term associated with one of the Fourier transforms in Eq. 9 is the same as the transmission of a thin lens having a focal length of  $100f$ . In the reconstruction process the nonlinear phase function will combine with the reconstruction lens to give an effective focal length of  $0.99f$ . This means that we shall have the two images in the reconstruction process. One is located on the back focal plane of the reconstruction lens, and the other is separated from the first image by a distance equal to  $0.01f$ .

Figure XI-6 shows the hologram corresponding to the wave front in Eq. 9. Figure XI-7 illustrates the arrangement used to reconstruct the images from the hologram and the three planes from which we can obtain sharp images. The images reconstructed from the hologram on the three planes are shown in Fig. XI-8. The letter 'E' has the Fourier transform, which is not modified by the nonlinear phase function. Therefore its image appears on the back focal plane of the lens. The symbol '+' has two images

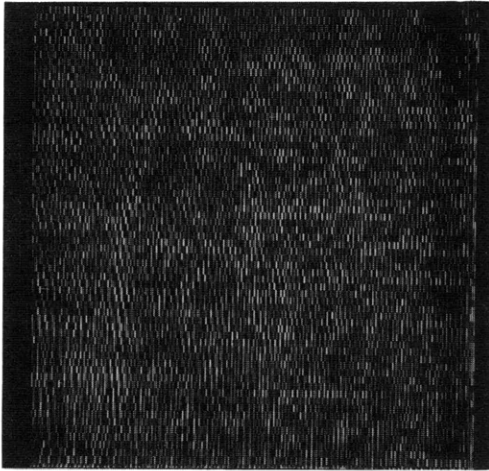


Fig. XI-6.  
Holograms of a three-dimensional object.

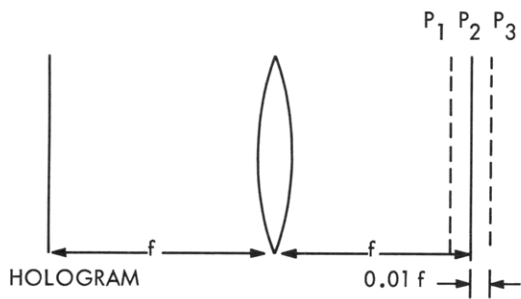


Fig. XI-7.  
Optical arrangement for reconstructing the image from a Fourier transform hologram.

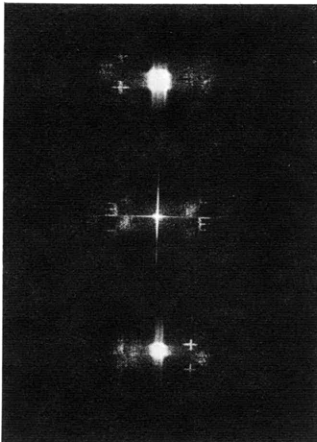


Fig. XI-8.  
Images reconstructed from the 3-D hologram. From top to bottom: image on plane P3; image on back focal plane of the lens P2; and image on P1.

## (XI. COGNITIVE INFORMATION PROCESSING)

in the reconstruction. They are the conjugate images of the symbol '+'.

The images reconstructed from the hologram have one interesting property. Because the technique that we use to record the depth information of the object is equivalent to putting a small thin lens in contact with the Fourier transform of each of the layers of the object, the effective focal length of the reconstruction for the different layers is different. In using a lens and coherent light to obtain the inverse Fourier transform, the result is proportional to

$$f(x, y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} F(u, v) \exp(j2\pi(ux+vy)/\lambda f) \, du \, dv. \quad (10)$$

Therefore the reconstructed image of the different layer will have a scale factor that is a function of the effective focal length. This property of the hologram can be observed in the reconstructed images of the symbol '+'. The image on plane P1 has an effective focal length shorter than the effective focal length of the image on plane P3. Hence the image on plane P1 is smaller than the image on plane P3. Generally, if we want to use the formula in Eq. 9, we must select a scale factor for the images in the different planes to compensate for the magnification in the reconstruction process.

W. H. Lee

### References

1. W-H. Lee, "Sampled Fraunhofer Hologram Generated by Computer," Quarterly Progress Report No. 88, Research Laboratory of Electronics, M. I. T., January 15, 1968, pp. 310-315.
2. System/360 Scientific Subroutine Package (360A-CM-03X) Version III Programmer's Manual. (276 p.)
3. H. J. Gerritsen, W. J. Hannen, and E. G. Ramberg, "Elimination of Speckle Noise with Redundancy," *Appl. Opt.* 7, 2301-2311 (1968).
4. G. B. Anderson and T. S. Huang, "Errors in Frequency Domain Processing of Image," Spring Joint Computer Conference, May 1969.
5. J. P. Waters, "Three-Dimensional Fourier-Transform Method for Synthesizing Binary Hologram," *J. Opt. Soc. Am.* 58, 1284 (1968).

## B. PICTURE BANDWIDTH COMPRESSION USING FOURIER ANALYSIS

1. Basic System Design
  - a. Thresholding

Two main factors led to the bandwidth compression achieved by the scheme described in this report. The reduction in bits used for image representation is mainly due to



thresholding of coefficients in the frequency domain. Detection of areas of little importance and subsequent simplified representation for these areas is also used in reducing bandwidth.

The method of data compression begins by subdividing an image into square subsections. The variance of each subsection is then measured to determine whether anything significant occurs in the subsection. If the variance is less than a threshold value, the subsection is represented only by its average value and its variance. If the variance exceeds the threshold, the subsection is expanded in a two-dimensional Fourier series, and the Fourier coefficients are subjected to an adaptive thresholding procedure. The adaptive threshold for the Fourier coefficients is linearly proportional to the subsection variance. The energy of each coefficient is tested against the adaptive threshold. The number of coefficients passing the adaptive threshold test is compared with an integer specifying the maximum allowable number of coefficients to be retained for transmission. If the maximum allowable number of coefficients is exceeded, the adaptive threshold is raised and the Fourier coefficients are tested against the new threshold. The process is repeated until an allowable number of Fourier coefficients is obtained.

Image reconstruction from subsection data is simple. Subsections represented only by a variance and an average value parameter are filled in by statistically independent Gaussian random variables with variance and mean equal to the transmitted parameters. Fourier inversion is used in subsections having Fourier coefficient data.

#### b. Quantization and Coding

Linear adaptive quantization of the Fourier coefficients is performed with the maximum Fourier coefficient, aside from the average value, determining the position of the largest quantization level for the magnitudes of the coefficients. The phases of the Fourier coefficients are linearly quantized on a scale from 0 to  $2\pi$ . Partitioning of the subsections represented by Fourier coefficients is based upon subsection variance. The subsections with the smaller variances are quantized with fewer levels than the higher energy subsections.

Data are needed to specify the positions in the frequency plane of the Fourier coefficients passing the adaptive energy threshold. Run-length coding is used, with the positions of suprathreshold coefficients denoted by ones and other positions denoted by zeros.

Each subsection of the image, including those with Fourier coefficient data, has quantized variance and average value parameters associated with it. A bit is also transmitted for each subsection to separate subsections without Fourier coefficient representation from those passing the variance test and having a Fourier coefficient representation. This additional bit permits the variance and average value parameters to be quantized with a differing number of levels for the two types of subsections.

## (XI. COGNITIVE INFORMATION PROCESSING)

### 2. Subjective Considerations

#### a. Contrast Sensitivity<sup>1</sup>

The Weber fraction  $\Delta B/B$ , which relates the just-noticeable difference in brightness  $\Delta B$  to the background brightness  $B$ , is nearly constant for a wide range of values of  $B$  at approximately 2 per cent. This contrast sensitivity of the eye causes errors in images to tend to be more noticeable in the darker areas of images. If an image is to be subjected to an error environment, then the subjective quality of the resulting picture normally improves if the effects of the error can be more evenly spread throughout the image, rather than being more pronounced in darker areas. Taking the logarithm of image brightness before processing and exponentiating after processing compensates for the eye's contrast sensitivity. This feature is incorporated in the design of our compression system.

#### b. Subsection Size

Choice of subsection size is influenced by different things. The larger subsection size is made, the more likely will areas of little energy be included with areas of large variance. This means that the total area of the picture represented by only average values and variances will decrease. Also, low-energy data, which need good reproduction, are more likely to be lost, because of the presence of higher energy data in the same subsection. This comes about because of the adaptive frequency-domain thresholding used in our compression scheme. Although it is not clear, it is also probably true that smaller subsection size results in more peaked Fourier spectra, because of simplification in data structure. The smaller the subsection, the less likely will different parts of the subsection cause frequency-domain averaging to result in a flattened spectrum. The peaked spectrum can result in a representation by fewer Fourier coefficients than the flat spectrum. The subsection size cannot become too small, however, because the data-compression ratio will then be severely limited. Each subsection is represented by at least a mean value and a variance. The total number of these parameters that is used to represent the entire image is inversely proportional to subsection size.

Subsection size was chosen to be  $16 \times 16$  samples. This choice divided our  $256 \times 256$  test picture into an integral number of subsections. It is interesting to note that when the picture appears at a distance from the eye of four times the picture height (normal viewing distance), the width of a subsection is roughly the period of the maximum spatial frequency response of the eye.

#### c. Variance Threshold

A grid was constructed on the test image of Fig. XI-9, partitioning the image into its subsections of size  $16 \times 16$  samples. Variances were calculated for the logarithm

of brightness in each subsection. The image with superimposed grid was then viewed to determine which subsections appeared as random noise and could be represented by

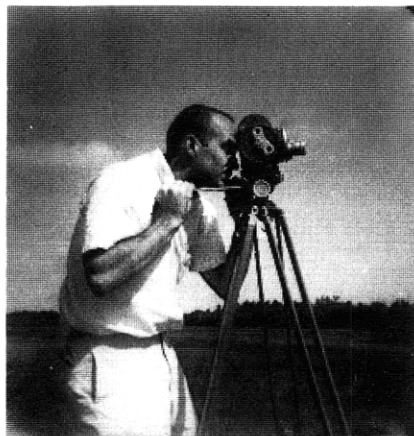


Fig. XI-9. Test picture.

only a mean value and a variance. We found that a variance threshold could separate all subsections that looked like random noise from other subsections. There are 89 such subsections represented in Fig. XI-9.

d. Padding

Frequency-domain thresholding can cause very noticeable errors. Strong discontinuities often exist between brightness along an edge of a subsection and the opposite edge of the same subsection. The subsection is considered as a period of a periodic two-dimensional function when the subsection is expanded in a Fourier series. This periodic two-dimensional function has discontinuities corresponding to the discontinuities between opposite edges of the subsection. This causes interference in the frequency domain with spectra resulting from features within the subsection. Even worse, the error along a subsection edge in the reconstructed image tends to be correlated and, therefore, most displeasing. The edge error tends to be correlated because discontinuities between opposite edges in the original subsection tend to be correlated.

To eliminate the edge-to-edge discontinuity problem, "padding" of the subsection can be performed. Strips of additional samples can be added to two adjacent edges of the subsection before Fourier expansion. The samples in a strip will be interpolated from samples along the edge to which the strip is connected and samples along the opposite edge.

e. Frequency-Domain Quantization

Quantization of an image causes artificial contours to appear if too few quantization levels are used.<sup>2</sup> Conventional television probably requires 6 or 7 bits of linear

## (XI. COGNITIVE INFORMATION PROCESSING)

brightness quantization for acceptable elimination of artificial contours. Random noise requires a higher energy to produce the same noticeability.<sup>3</sup> Therefore, frequency-domain quantization, which acts to produce randomized noise in the image, should be more acceptable than brightness quantization of the same mean-square error. This will be true if one pitfall is avoided. The mean values of the subsections must be very finely quantized; otherwise, contours, which are easily detected by the eye, will appear along the boundaries of subsections in the reconstructed image. A bandwidth saving occurs with coarse quantizing of the other Fourier series coefficients.

### f. Artificial Noise in Low-Variance Subsections

Random noise is generated in subsections of the reconstructed image, failing the variance threshold. If the noise in the input image is not too great, this should act to produce a more natural appearance in the system output. (Most areas of images have a granular or noisy appearance, because of a combination of texture in the original scene and film noise.) Also, the introduction of noise in these subsections tends to disguise defects that appear along subsection edges.

### 3. Experimental Results

The compression system described was simulated on the IBM 360/65 digital computer. A flying-spot scanner built by W. F. Schreiber served as an interface between the digital computer and picture hard copy.

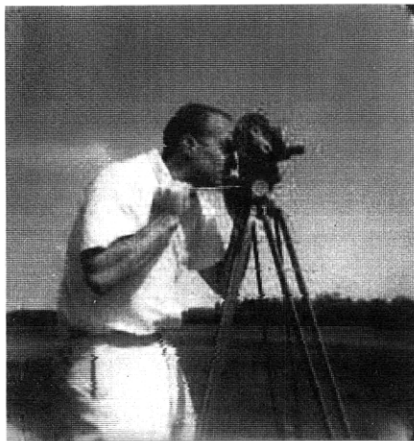


Fig. XI-10. Output picture.

The picture in Fig. XI-10 is the result of the system acting on our test picture of Fig. XI-9. A total of 46,632 bits was expended to "transmit" this image. This corresponds to a compression of 8.43 to 1 in bandwidth when comparison is made with conventional 6-bit PCM

G. B. Anderson

## References

1. W. F. Schreiber, "Picture Coding," Proc. IEEE 55, 320-330 (1967).
2. T. S. Huang, O. J. Tretiak, B. Prasada, and Y. Yamaguchi, "Design Considerations in PCM Transmission of Low-resolution Monochrome Still Pictures," Proc. IEEE 55, 331-335 (1967).
3. L. G. Roberts, "Picture Coding Using Pseudorandom Noise," IRE Trans., Vol. IT-8, pp. 145-154, February 1962.

## C. IMPULSE RESPONSE FOR A LENS WITH SEIDEL ABERRATIONS

For an optical system with axial symmetry, the retention of the first-order term in the series expansion of the angle characteristic leads to Gaussian optics. By including the next higher order in the expansion, which includes five terms of the same order, the five Seidel aberrations may be treated. Each of these aberrations – spherical aberrations, coma, astigmatism, curvature of field, and distortion – arises because its corresponding coefficient is nonzero.

For a lens satisfying Gaussian optics, a point source of light in the object plane produces a point source of light in the Gaussian image plane. Thus the impulse response, or point source response, of a system obeying Gaussian optics is again an impulse, but attenuated and displaced by the magnification factor. Let  $\underline{x}_o$  denote object plane coordinates,  $\underline{x}_i$  the image plane coordinates, and  $M$  the lateral magnification. Then an input brightness of  $\delta(\underline{x}_o - \underline{x})$  will produce an output brightness of  $\frac{\kappa(\underline{x})}{|M|} \delta(\underline{x}_i - M\underline{x})$ . The function  $\kappa(\underline{x})$  accounts for the fraction of light intercepted by the lens and the transmission loss through the lens. To determine the image intensity distribution attributable to a surface in the object plane, besides knowledge of the impulse response of the lens, it is necessary to know the photometric nature of the object surface. For example, it may radiate according to Lambert's law.

To account for aberrations in the lens, expressions for the ray aberration may be applied to find the impulse response. The ray aberration is a vector in the image plane from the Gaussian image point to the point where the actual ray intersects the image plane. For a Gaussian image point the ray aberration depends upon which ray from the object point is chosen.

Suppose the coordinate systems are as shown in Fig. XI-11. Let the following functions be given

$$x_i = Mx_o + \Delta_x(u, v; x_o, y_o)$$

$$y_i = My_o + \Delta_y(u, v; x_o, y_o)$$

(XI. COGNITIVE INFORMATION PROCESSING)

The ray aberration is given by the vector  $(\Delta_x, \Delta_y)$ . That is, the ray originating at object point  $(x_o, y_o)$  and passing through exit pupil point  $(u, v)$  will deviate from the Gaussian

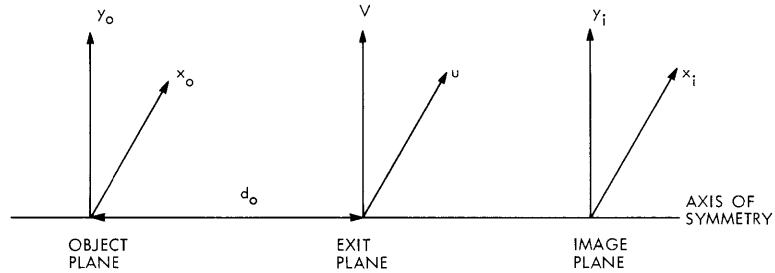


Fig. XI-11. Coordinate system.

image point  $(Mx_o, My_o)$  by  $(\Delta_x, \Delta_y)$ . In general, for a fixed object point several values of  $(u, v)$  may yield the same ray aberration. Let  $B_e$  be the brightness (power/unit area) at exit pupil point  $(u, v)$ , and  $B_i$  the brightness at the corresponding point  $(x_i, y_i)$ . A small

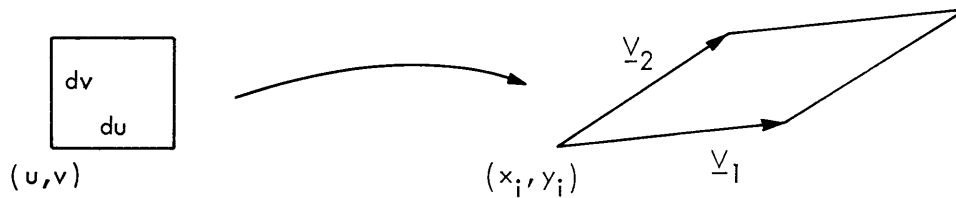


Fig. XI-12. Exit to image plane mapping.

rectangle in the exit pupil will be mapped into a small parallelogram in the image plane as shown in Fig. XI-12. Under the assumption that the power transferred down this ray tube is conserved,

$$B_e |dudv| = B_i |\underline{V}_1 \times \underline{V}_2|$$

$$B_i = B_e \left| \frac{dudv}{\underline{V}_1 \times \underline{V}_2} \right|$$

$$\underline{V}_1 = \left( \frac{\partial \Delta_x}{\partial u}, \frac{\partial \Delta_y}{\partial u} \right) du$$

$$\underline{V}_2 = \left( \frac{\partial \Delta_x}{\partial v}, \frac{\partial \Delta_y}{\partial v} \right) dv$$

$$B_i = B_e(u, v; x_o, y_o) \Big/ \det \begin{bmatrix} \frac{\partial \Delta_x}{\partial u} & \frac{\partial \Delta_y}{\partial u} \\ \frac{\partial \Delta_x}{\partial v} & \frac{\partial \Delta_y}{\partial v} \end{bmatrix}$$

$$= B_e(u, v; x_o, y_o) (J(u, v; x_o, y_o))^{-1}.$$

Let  $n$  points in the exit pupil  $(u_1, v_1) \dots (u_n, v_n)$  be mapped into the point  $(x_i, y_i)$ . Then the brightness at  $(x_i, y_i)$  is

$$B_i(x_i, y_i; x_o, y_o) = \sum_{j=1}^n B_e(u_j, v_j; x_o, y_o) (J(u_j, v_j; x_o, y_o))^{-1}. \quad (1)$$

Now it is simply an application of the formula (1) to find the impulse response for the five aberrations. All the ray aberration formulas that follow have been taken from Born and Wolf.<sup>1</sup> Also, the object point will be at  $(0, y_o)$ , and polar coordinates  $(p, \theta)$  in the exit pupil will be used to specify the ray aberration functions. The angle  $\theta$  is measured positive counterclockwise from the V-axis of Fig. XI-11. It will be assumed throughout that the optical system is simply a lens of radius  $R_L$ .

### 1. Spherical Aberration

The ray aberrations for spherical aberration are

$$\Delta_x = Bp^3 \sin \theta$$

$$\Delta_y = Bp^3 \cos \theta$$

and the Jacobian is

$$J(u, v; x_o, y_o) = 3B^2 p^4.$$

Let  $(r_o, \theta_o)$  be polar coordinates of a unit point source in the object plane, and  $(r_i, \psi_i)$  be polar coordinates measured from the Gaussian image point. Both  $\theta_o$  and  $\psi_i$  should be measured from the  $x$  axes.

$$B_e(p, \theta; r_o, \theta_o) = \frac{d_o}{4\pi} \left[ r_o^2 + p^2 + d_o^2 + 2pr_o \sin(\theta - \theta_o) \right]^{-3/2}$$

Also,  $r_i = Bp^3$ , and  $\psi_i = \theta + \frac{\pi}{2}$ . The variable  $d_o$  is the distance between the object plane

(XI. COGNITIVE INFORMATION PROCESSING)

and the exit plane.

The output brightness distribution is

$$B_i(r_i, \psi_i; r_o, \theta_o) = \frac{1}{12\pi B^2 d_o^6} \left( \frac{B d_o^3}{r_i} \right)^{4/3} \left[ 1 + \left( \frac{r_o}{d_o} \right)^2 + \left( \frac{r_i}{B d_o^3} \right)^{2/3} - 2 \left( \frac{r_i}{B d_o^3} \right)^{1/3} \frac{r_o}{d_o} \cos(\theta_o - \psi_i) \right]^{-3/2} \quad \text{for } r_i \leq BR_L^3$$

$$= 0 \quad \text{for } r_i > BR_L^3. \quad (2)$$

Regardless of the position of the point source in the object plane, a spot with nonuniform brightness of radius  $BR^3$  is the image. In general, the spot will have the same shape as

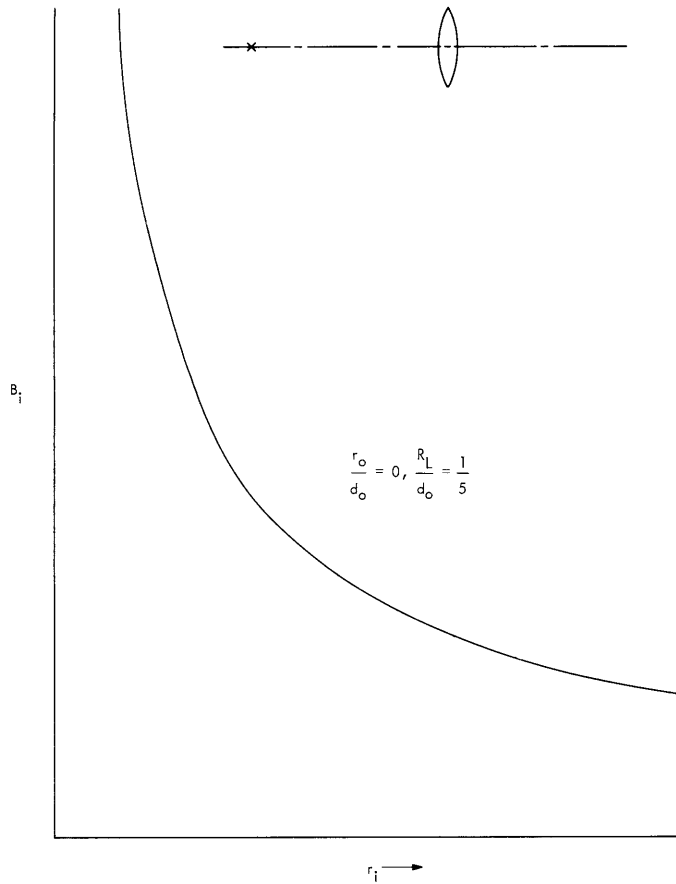


Fig. XI-13. Spherical aberration: Impulse response for an on-axis object.



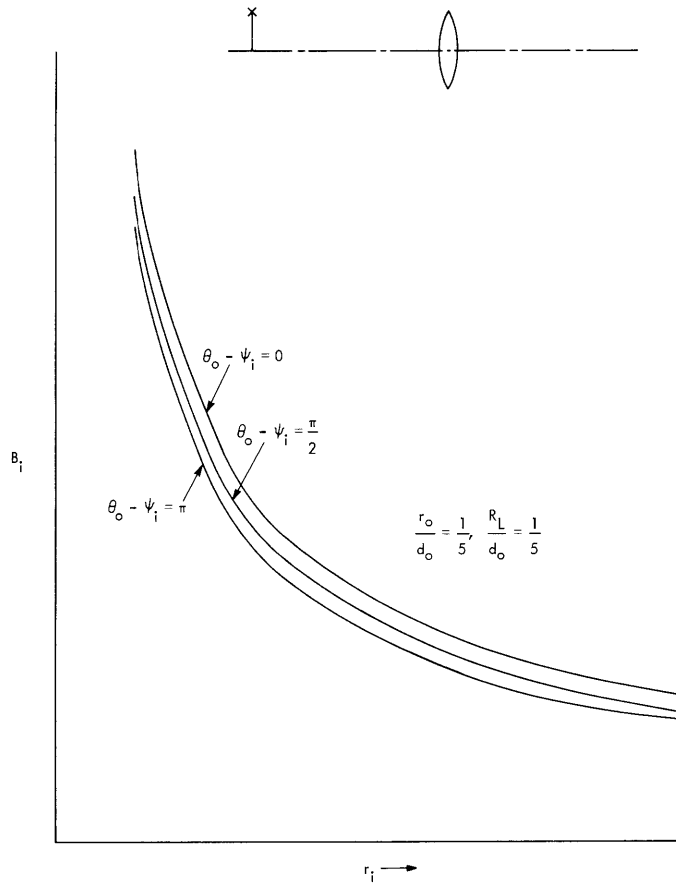


Fig. XI-14. Spherical aberration: Impulse response for an off-axis object.

the exit pupil aperture, and a brightness distribution dependent upon the object point position. Even though the general formula for the impulse response indicates that it is not object position invariant, Figs. XI-13 and XI-14 show, for "reasonably small" values of  $\frac{r_o}{d_o}$  and  $\frac{R_L}{d_o}$ , that the response is practically

$$B_e = \begin{cases} \frac{1}{12\pi B^{2/3} d_o^2} \frac{1}{r_i^{4/3}}, & r_i \leq BR_L^3 \\ 0, & r_i > BR_L^3 \end{cases}$$

## 2. Coma

The ray aberrations for coma are

$$\Delta_x = -Fy_o p^2 \sin 2\theta$$

(XI. COGNITIVE INFORMATION PROCESSING)

$$\Delta_y = -Fy_0 p^2 (2 + \cos 2\theta).$$

The boundary for the impulse response is shown in Fig. XI-15. For a fixed object point  $(0, y_0)$ , 4 points of the exit pupil map into each point in region I, and 2 points

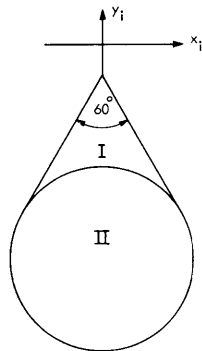


Fig. XI-15. Boundary of the coma pattern.

map into each point in region II. Let  $(x, y)$  be measured from the Gaussian image point.

$$x = x_i$$

$$y = y_i - My_0.$$

The brightness distribution for  $x_i, y_i$  in region I is

$$B_i(x, y; 0, y_0) = \frac{B_e(p_1, \theta_{11}; 0, y_0)}{4F^2 y_0^2 p_1^2 |1 - 2 \cos 2\theta_{11}|} + \frac{B_e(p_1, \theta_{12}; 0, y_0)}{4F^2 y_0^2 p_1^2 |1 - 2 \cos \theta_{12}|} \\ + \frac{B_e(p_2, \theta_{21}; 0, y_0)}{4F^2 y_0^2 p_2^2 |1 - 2 \cos 2\theta_{21}|} + \frac{B_e(p_2, \theta_{22}; 0, y_0)}{4F^2 y_0^2 p_2^2 |1 - 2 \cos 2\theta_{22}|}. \quad (3)$$

For  $x_i, y_i$  in region II the brightness distribution includes only the first two terms. The variables  $p_1, \theta_{11}, \theta_{12}, p_2, \theta_{21}, \theta_{22}$  are functions of  $x, y$  as follows:

$$p_1^2 = \frac{1}{Fy_0} \left( -\frac{2}{3}y - \frac{1}{3}\sqrt{y^2 - 3x^2} \right)$$

$$p_2^2 = \frac{1}{Fy_0} \left( -\frac{2}{3}y + \frac{1}{3}\sqrt{y^2 - 3x^2} \right)$$

$$\theta_{j1} = \frac{3\pi}{2} - \frac{1}{2}\psi_j$$

$$j = 1, 2$$

$$\theta_{j2} = \begin{cases} \frac{\pi}{2} - \frac{1}{2} \psi_j, & 0 \leq \psi_j \leq \pi \\ \frac{5\pi}{2} - \frac{1}{2} \psi_j, & \pi < \psi_j \leq 2\pi \end{cases}$$

The angles  $\psi_1$  and  $\psi_2$  are defined by Fig. XI-16.

### 3. Astigmatism and Curvature of Field

The ray aberrations in this case are

$$\Delta_x = Dpy_0^2 \sin \theta$$

$$\Delta_y = (2C+D)py_0^2 \cos \theta.$$

The elliptical boundary for the impulse response is shown in Fig. XI-17. Fortunately, this is a one-to-one mapping between the exit pupil and image plane, as was the mapping for spherical aberration. The brightness distribution is

$$B_i(x, y; 0, y_0) = \begin{cases} \frac{1}{|D(2C+D)|y_0^4} B_e(p, \theta; 0, y_0), & p \leq R_L \\ 0, & p \geq R_L \end{cases}$$

The inversion of the mapping is

$$p^2 = \frac{x^2}{4D^2y_0^4} + \frac{y^2}{4(2C+D)^2y_0^4}$$

and  $\theta = \psi$ , where

$$\sin \psi = -\frac{x}{\sqrt{x^2 + y^2}} \quad \text{and} \quad \cos \psi = \frac{y}{\sqrt{x^2 + y^2}}.$$

Besides the elliptical, instead of circular, boundary, this brightness distribution is significantly different from the one for spherical aberration because there is no singularity in the brightness distribution at the Gaussian image point.

### 4. Distortion

The ray aberrations for distortion are

$$\Delta_x = 0$$

$$\Delta_y = -Ey_0^3.$$

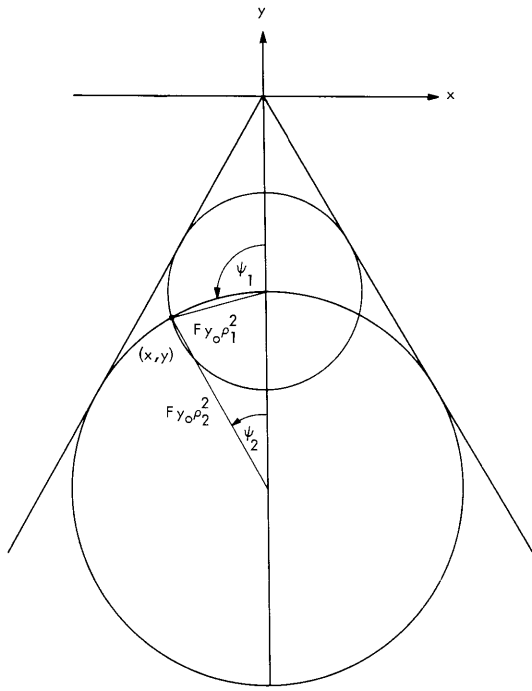


Fig. XI-16.  
 Definition of  $\psi_1$  and  $\psi_2$  for inverting  
 the mapping.

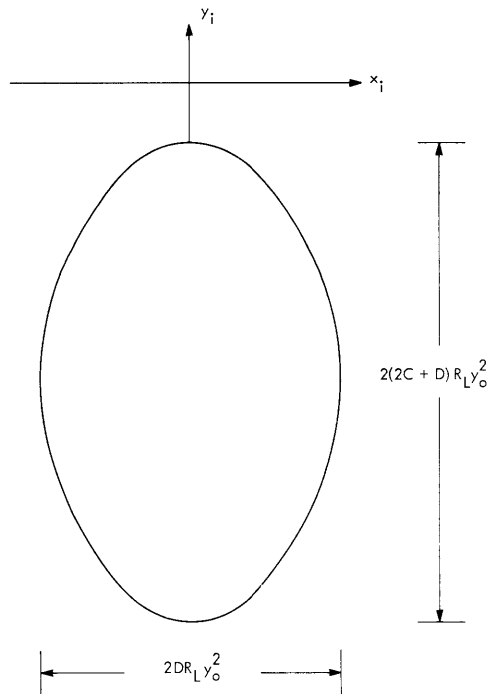


Fig. XI-17.  
 Boundary for the astigmatism and  
 curvature of field pattern.

Since  $(\Delta_x, \Delta_y)$  are not functions of the exit pupil point, the brightness distribution is simply another impulse, but displaced.

G. M. Robbins

#### References

1. M. Born and E. Wolf, Principles of Optics (Pergamon Press, London, 3rd edition, 1965).

#### D. SPECIFICATION OF THE PROSODIC FEATURES OF SPEECH

For several years, many research workers have been developing speech synthesis algorithms that operate directly on a string of phonemes representing the sound representation of the sentence under consideration. This phoneme sequence has usually been derived manually, but recently a scheme has been invented<sup>1</sup> to transform written words to their phonemic representation. Speech based on word-by-word phonemic representation is usually intelligible, but not suitable for long-term use. Several problems remain, apart from those concerned directly with speech synthesis by rule from phonemic specifications.

First, many words can be nouns or verbs, depending on context [refuse, incline, survey], and proper stress cannot be specified until the intended syntactic form class is known. Second, punctuation and phrase boundaries may be used to specify pauses that help to make the complete sentence understandable. Third, more complicated stress contours over phrases can be specified which facilitate sentence perception. Finally, intonation contours, or "tunes," are important for designating statements, questions, exclamations, and continuing or terminal juncture. These features (stress, intonation, and pauses) comprise the main prosodic or suprasegmental features of speech.

Several experiments<sup>2-4</sup> have shown that we tend to perceive sentences in chunks or phrasal units, and that the grammatical structure of these phrases is important for the correct perception of the sentence. It is also known<sup>5</sup> that the syntactic structure of the sentence is sufficient to specify much of the stress, pauses, and intonation of natural speech. Hence prosodic features of speech can be used in a limited fashion to help point out the intended surface syntactic structure of synthesized speech, which is then useful in the perception of the phonetic shape of the sentence.

It is desirable, then, to perform a limited parsing of sentences in order to allow the specification of the prosodic features of speech by rule. Such a scheme requires, first, the determination of the parts of speech of the words of the sentence; second, a phrase-level parsing of the sentence; and finally, execution of a phonological algorithm that computes the prosodic features from the previously derived sentence structure.

## (XI. COGNITIVE INFORMATION PROCESSING)

First, we consider determination of the parts of speech of words. The starting point for this procedure is the given word converted into one or more "morphs," each of which is either a prefix, base word, or suffix. Thus [grasshopper] → [grass] + [hop] + [er], [browbeat] → [brow] + [beat], and [unfit] → [un] + [fit]. Each of these morphs corresponds to a dictionary entry that contains, in addition to phonemic specifications, parts-of-speech information. In the case of base word morphs, this information consists of a set of parts of speech for that word, called the "grammatical homographs" of the word. For prefixes and suffixes, information is given indicating the resultant part of speech when the prefix or suffix is concatenated with a base word. Thus [-ness] always forms a noun, as in [goodness].

Other researchers<sup>6,7</sup> have used a computational dictionary to compute parts of speech, relying on the prevalence of function words (determiners, prepositions, conjunctions, and auxiliaries), together with suffix rules of the type just described and their accompanying exception lists. This procedure, of course, keeps the lexicon small, but results in arbitrary parts-of-speech classification when the word is not a function word, and does not have a recognizable suffix. Furthermore, ambiguous suffixes such as [-s] (implying plural noun or singular verb) carry over their ambiguity to the entire word, whereas if the root word has a unique part of speech, like [cat], our procedure gives a unique result; [cats] (plural noun). Hence the morph lexicon can often be used to advantage, especially in the prevalent noun/verb ambiguities.

The parts-of-speech algorithm considers each morph of the word and its relation with its left neighbor, starting from the right end of the word. If there are two or more suffixes [commendables], they are entered into a last-in first-out pushdown list. Then the top suffix is joined to the root morph, and the additional suffixes are concatenated until the list is empty. Compounding is done next, and finally any prefixes are attached. Prefixes generally do not affect the part of speech of the base word, except for [em-, en-, be-] which imply verb as the resultant part of speech. Compounds can occur in English in any of three ways, and there appears to be no reliable method for distinguishing these classes. There can be two separate words [bus stop], two words hyphenated [hand-cuff], or two words concatenated directly, as in [sandpaper]. The parts-of-speech algorithm treats the last two cases, leaving the two-word case for the parser to handle. The algorithm ignores the presence of a hyphen, and then processes hyphenated and one-word compounds as though they were both single words. The parts of speech of the two elements of the compound are considered as row and column entries to a matrix whose cells yield the resulting part of speech. Thus Adverb · Noun → Noun, as in [underworld]. Combinations of suffixes with compounds [handwriting] can be accommodated, as well as one-word compounds containing more than two morphs.

A special routine is provided to handle troublesome suffixes such as [-er, -es, -s], which often reduces the resultant number of parts of speech to a minimum.

In this way, the algorithm makes use of the parts-of-speech information of the individual morphs to compute the parts-of-speech set for the word formed by these morphs. These sets then serve as input to the parser, after having first been ordered to suit the principles of the parser.

### 1. Parsing

Our goal is to determine the syntactic structure that is sufficient to specify the prosodic features of the sentence, which can then serve as cues to the perception of the intended text. Since we are trying to provide only a limited number of such cues (enough to allow the structure to be deduced), we have designed a limited parser that reveals the syntax of only a portion of the sentence. We have tried to find the simplest parser consistent with these phonological goals that would also use minimum core storage and run fast enough to allow for a realistic speaking rate, say, 150-180 words per minute. Because the absence, or incorrect implementation of prosodics in a small percentage of the output sentences is not likely to be catastrophic, we can tolerate occasional mistakes by the parser, but we have tried to achieve 90 per cent accuracy. These requirements, for a limited, phrase-level parser operating in real-time at comfortable speaking rates within restricted core storage, are indeed severe, and many features found in other parsers are absent here. We do not use a large number of parts-of-speech classifications, nor do we exhaustively cycle through all the homographs of the words of a sentence to find all possible parsings. Inherent syntactic ambiguity ([They are washing machines]) is ignored, the resulting phrase structures being biased toward noun phrases and prepositional phrases. No deep-structure "trees" are obtained, since these are not needed in the phonological algorithm, and only noun phrases and prepositional phrases are detected, so that no sentencehood or clause-level tests are made. We do, however, compute a bracketed structure within each detected phrase, such as [the [old house]] and [in [[brightly lighted] windows]], since this structure is required by the phonological algorithm. The result is a context-sensitive parser that avoids time-consuming enumerative procedures, and consults alternative homographs only when some condition is detected (such as [to] used to introduce either an infinitive or a prepositional phrase) which requires such a search.

The parser makes two passes (left-to-right) over a given input sentence. The first pass computes a tentative bracketing of noun phrases and prepositional phrases. Inasmuch as this initial bracketing makes no clause-level checks and does not directly examine the frequently occurring noun/verb ambiguities, it is followed by a special routine designed to resolve these ambiguities by means of local context and grammatical number agreement tests. These last tests are also designed to resolve noun/verb ambiguities that do not occur in bracketed phrases, as [refuse] in [They refuse to leave]. As a result

## (XI. COGNITIVE INFORMATION PROCESSING)

of these two passes, a limited phrase bracketing of the sentence is obtained, and some ambiguous words have been assigned a unique part of speech, yet several words remain as unbracketed constituents.

The first pass is designed to quickly set up tentative noun phrase and prepositional phrase boundaries. This process may be thought of as operating in three parts. The program scans the sentence from left to right looking for potential phrase openers. For example, determiners, adjectives, participles, and nouns may introduce noun phrases, and prepositional phrases always start with a preposition. In the case of some introducers, such as present participles, words further along in the sentence are examined, as well as previous words, to determine the grammatical function of the participle, as in [Wiring circuits is fun.]. Once a phrase opener has been found, very quick relational tests between neighboring words are made to determine whether the right phrase boundary has been reached. These checks are possible because English relies heavily on word order in its structure. Having found a tentative right phrase boundary, right context checks are made to determine whether or not this boundary should be accepted. After completion of these checks, the phrase is closed and a new phrase introducer is looked for. The procedure continues until the end of the sentence is reached.

When the bracketing is complete, further tests are made to check for errors in bracketing caused by frequent noun/verb ambiguities. For example, the sentence [That old man lives in the gray house.] would be initially bracketed

[That old man lives]<sub>NP</sub> [in the gray house]<sub>Prep P</sub>.

Notice that sentencehood tests (although not performed by the parser) would immediately reveal that the sentence lacks a verb, and further routines could deduce that [lives], which can be either a plural noun or a third person singular verb, is functioning as a verb, although the bracketing routine, since it is biased toward noun homographs, made [lives] part of the noun phrase. We also note the importance of this error for the phonetic shape of the sentence, since [lives] changes its phonemic structure according to its grammatic function in the sentence. An agreement test, however, compares the rightmost "noun" with any determiners that may reflect grammatical number. In this case, [that] is a singular demonstrative pronoun, so we know that [lives] does not agree with it, and hence must be a verb. After the agreement test has been made for each noun phrase, local context checks are used in an attempt to remove noun/verb ambiguities that are important for the phonological implementation, and yet have not been bracketed into phrases containing more than one word. Thus in the sentence [They produce and develop many different machines.], the algorithm would note that [produce] is immediately preceded by a personal pronoun in the nominative case, and hence the word is functioning as a verb. Such knowledge can then be used to put stress on the second syllable of the word in accordance with its function.



At the conclusion of the parsing process described above, phrase boundaries for noun phrases and prepositional phrases have been marked, but the structure within the phrase is not known. In order to apply the rules that are used for computing stress patterns within the phrase, however, internal bracketing must be provided. For this reason, determiner-adjective-noun sequences are given a "progressive" bracketing, as [the[long[red barn]]], whereas noun phrases beginning with adverbials are given "regressive" bracketings, as [[[very brightly] projected] pictures]. A preposition beginning a prepositional phrase always has a progressive relation to the remaining noun phrase, so that we have [in [the [long [red barn]]]] and [ of [[[ very brightly] projected] pictures]]. Furthermore, two nouns together, as in [the local bus stop], are marked as a compound for use by the phonological algorithm.

The procedure described above is thus able to detect noun phrases and prepositional phrases and to compute the internal structure of these phrases. The grammar and parsing logic are intertwined in this procedure, so that an explicit statement of the grammar is difficult. Nevertheless, the rules are easily modified. At present, however, the provision of prosodics is supplied for noun phrases and prepositional phrases only.

## 2. Phonological Algorithm

Once noun phrases and prepositional phrases have been detected, the phonological algorithm uses the surface syntactic bracketing, plus punctuation and clause-marker words, to deduce the pattern for stress, pauses, and intonation related to the detected phrases.

Stress is implemented within the detected phrases by iterative use of the stress cycle rules, described by Chomsky and Halle.<sup>5</sup> These rules operate on the two constituents within the innermost brackets to specify where main stress should be placed. All other stresses are then "pushed down" by one. (Here, "one" is the highest stress.) The innermost brackets are then "erased," and the rules applied to the next pair of constituents. The cycle is then continued until the phrase boundaries are reached. For compounds, the rules specify main stress on the leftmost element (compound rule), whereas for all other syntactic units (e. g. , phrases) main stress goes on the rightmost unit (nuclear stress rule). For example, we have

[the[long[red barn]]]  
                   2   1  
 4   2   3   1

where initially stress is 1 on all units except the article the, and two cycles of the phrase rule are used. The parser has, of course, provided the bracketing of the phrase. Also,

(XI. COGNITIVE INFORMATION PROCESSING)

[ in [[[ very brightly] lighted] rooms]]

	2	1		
	3	2	1	
4	4	3	2	1

requires three applications of the rules, and

[ the [new [bus stop]]]

	1	2	
2	1	3	
4	2	1	3

which contains a compound, requires two iterations. Each morph in the lexicon is given lexical stress, so that the phrase-level stress numbers described above provide the framework for the local stress variations of individual morphs or words.

Pauses are provided in a definite hierarchy throughout each sentence. Although we have not determined optimal values for these durations, pauses are used at phrase boundaries, where punctuation appears, and before clause-marker words such as [that, since, which]. Finally, terminal pauses are used between sentences. Much work remains to be done on the acoustic correlates of juncture and stress, and several experiments are currently contributing to our understanding of these phenomena.

The provision of intonational  $f_0$  contours by rule has been described by Mattingly,<sup>8</sup> and our technique is similar to his. The slope of the  $f_0$  contour is controlled by the specific phonemes encountered in the sentence, and by the nature of the pause at the end of the phrasal unit. Rising terminal contours are specified at the end of interrogative clauses just preceding the question mark, except when the clause starts with a [wh-] word, as [where is the station?]. In the absence of a question mark, the intonation  $f_0$  contour is falling with a slope determined by rules like those of Mattingly.

Work is continuing on improvements for the parser. We are focusing particularly on adverbial functions and the use of prepositions in verb phrases. We hope that several additional classes of phrases can be detected within the present parser framework. Also, the parsing scheme is being implemented to run in real time on a PDP-9 computer. We have developed a list structure to expedite this task, and the algorithms are being coded by using these low-level list handlers to provide maximum speed in minimum core storage.

J. Allen

References

1. F. F. Lee, "A Study of Grapheme to Phoneme Translation of English," Ph.D. Thesis, M. I. T., 1965.
2. G. A. Miller, "Decision Units in the Perception of Speech," IRE Trans., Vol. IT-8, No. 2, p. 81, February 1962.

## (XI. COGNITIVE INFORMATION PROCESSING)

3. G. A. Miller, G. A. Heise, and W. Lichten, "The Intelligibility of Speech as a Function of the Context of the Test Materials," *J. Exptl. Psychol.* 41, 329 (1951).
4. G. A. Miller and S. Isard, "Some Perceptual Consequences of Linguistic Rules," *J. Verb. Learn. and Verb. Behav.* 2, 217 (1963).
5. N. Chomsky and M. Halle, Sound Patterns of English (Harper and Row Publishers, Inc., New York, 1968).
6. S. Klein and R. F. Simmons, "A Computational Approach to the Grammatical Coding of English Words," *J. Acoust. Computing Machinery* 10, 334 (1963).
7. D. C. Clarke and R. E. Wall, "An Economical Program for the Limited Parsing of English," AFIPS Conference Proceedings, 1965, FJCC, p. 307.
8. J. G. Mattingly, "Synthesis by Rule of Prosodic Features," *Language & Speech* 9, 1 (1966).

### E. "OPTICAL BENCH" SIMULATOR

A program has been developed for the simulation of one-dimensional optical systems.

Analysis of an optical system is performed with the use of results derived by Vander Lugt<sup>1</sup> and Goodman.<sup>2</sup> The approach is to characterize each element of an optical system, for example, lens or stop, by its impulse response, or equivalently, by its system function in the frequency domain. The light distribution may then be followed through the system and sequentially processed by each element.

A Fortran program has been written to accomplish the procedure outlined above. The equations were put in a discrete form, and a Fast Fourier Transform package was utilized. The sample spacing in the field was scaled to 1  $\mu$ . Maximum field size is limited only by available core memory. Test programs have been run on the M. I. T. Computation Center's IBM System 360 with field lengths of a few centimeters.

The user's program consists in calls to various subroutines with appropriate arguments. Each subroutine corresponds to a basic element of the optical system, or to a certain output function. The available elements are briefly described as follows.

INPUT – sets up a monochromatic plane wave across the field.

STOP – places a stop in the optical path whose transmittance function can be square or sinusoidal over variable limits.

LENS – places a lens in the optical path whose radius and focal length are variable.

SPACE – gives a desired separation to any of the elements above.

These elements may be cascaded indefinitely, and objects may be called consecutively to form compound elements. Two output routines, OUTPUT and CCOUT, which plot the light intensity as a function of position for any point on the axis, are available. OUTPUT plots the intensity on the computer print-out, while CCOUT utilizes the Cal Comp plotter. The subroutines are called in the order of their position on the

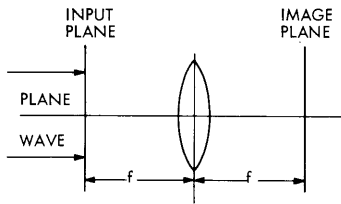


Fig. XI-18. Configuration of a Fourier-transforming system.



Fig. XI-19. (a) Input transmittance of a simple spot.  
(b) Expected output image for the input of (a).

```

COMPLEX IMJ(16384)
REAL OUTAX(10001),OUT(10001)
CALL NEWPLT('M0000','I0000','WHITE','BLACK')
CALL INPUT(IMJ,16384,1,0.0,-50,-5,5,50)
CALL SPACE(IMJ,16384,2.0,5000.)
CALL LENS(IMJ,16384,2.0,5000.,8192)
CALL SPACE(IMJ,16384,2.0,5000.)
CALL STOP(IMJ,16384,1,1.0,-8192,-8192,-500,500)
CALL CCOUT(IMJ,16384,24.,10001,2,OUTAX,OUT)
CALL INPUT(IMJ,16384,2,0.0,-50,50,1,0)
CALL STOP(IMJ,16384,1,1.0,-8192,-8192,-5,5)
CALL SPACE(IMJ,16384,2.0,5000.)
CALL LENS(IMJ,16384,2.0,5000.,8192)
CALL SPACE(IMJ,16384,2.0,5000.)
CALL STOP(IMJ,16384,1,1.0,-8192,-8192,-500,500)
CALL CCOUT(IMJ,16384,24.,10001,3,OUTAX,OUT)
CALL ENDPLT
CALL EXIT
END

```

Fig. XI-20. Simulator program.

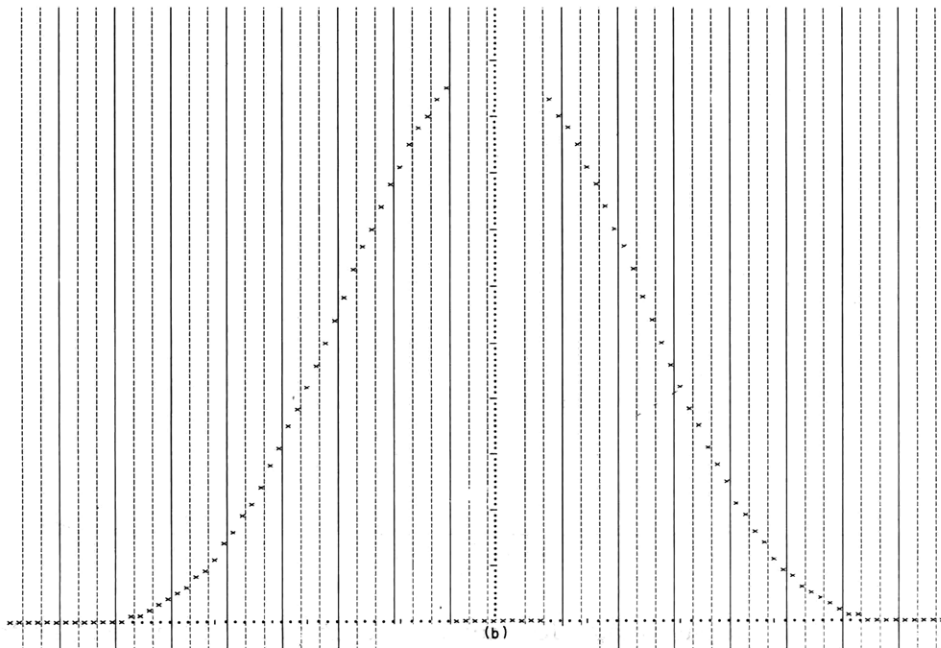
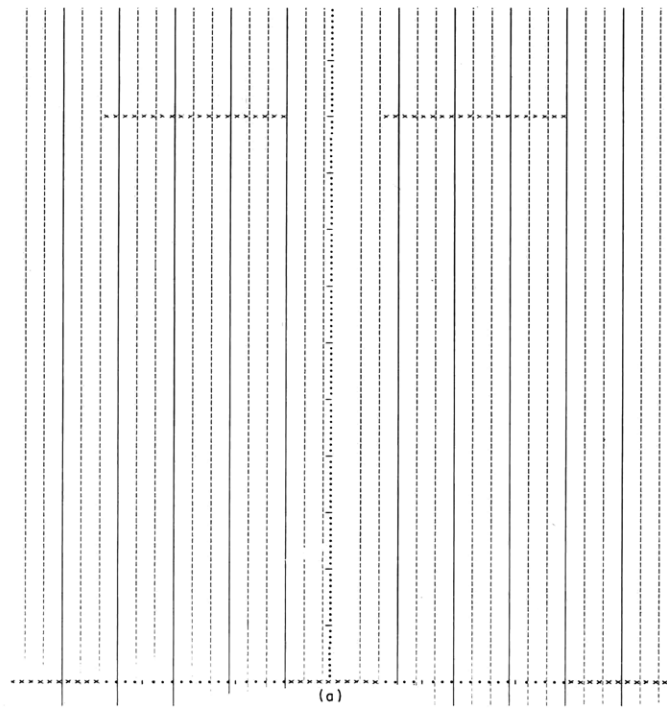


Fig. XI-21. (a) Input intensity for square transmittance function.  
 (b) Input intensity for sinusoidal transmittance function.

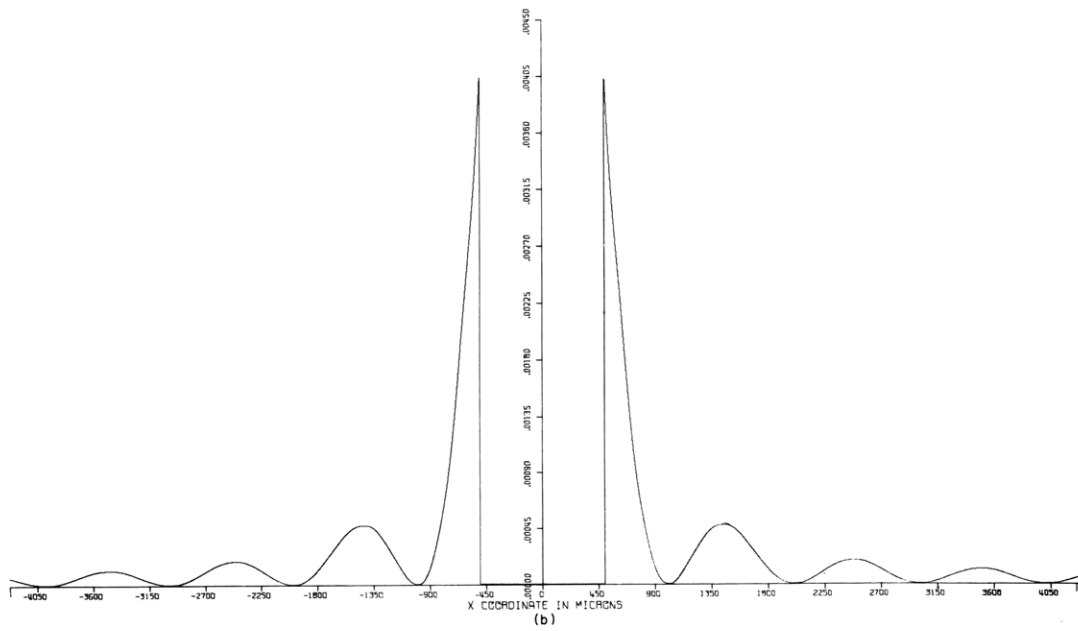
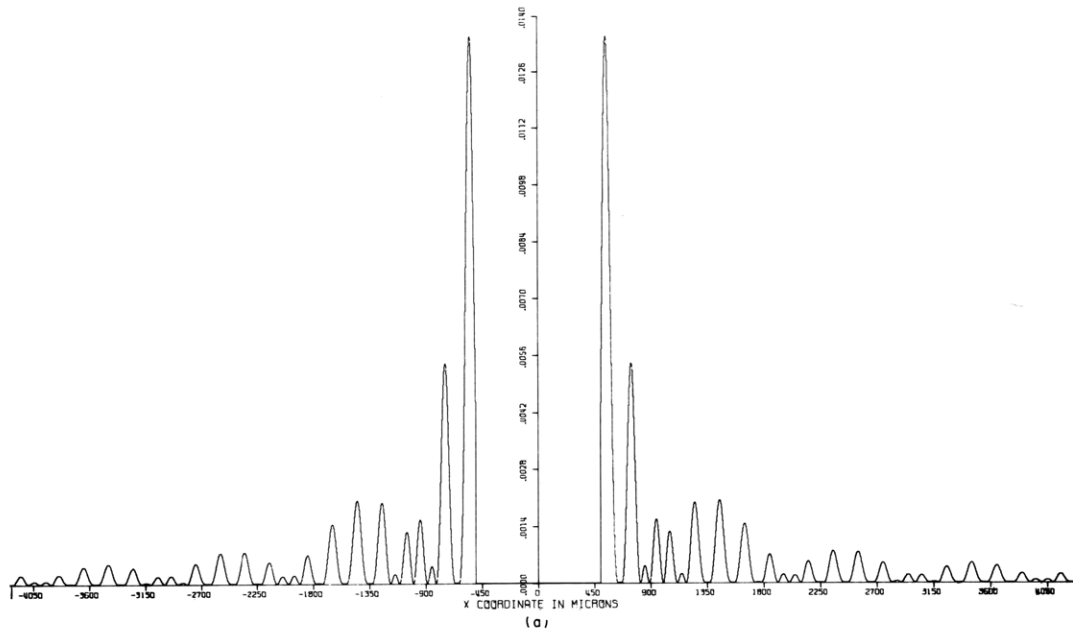


Fig. XI-22. (a) Output image for square input transmittance function.  
 (b) Output image for sinusoidal input transmittance function.

axis or their desired execution. Plots may be called at any point without affecting the following operations.

One of the advantages of the system is the ease with which the experimenter can study the effects of details such as lens size on his desired image. The effects of stops at various positions and with various sizes can be easily judged by direct comparison of the output images. The following example demonstrates the use of the simulator for comparing the effects of two different types of stops.

Figure XI-18 shows the configuration of a Fourier-transforming optical system. The transform of the input field, appropriately scaled, appears at the back focal plane of the lens. We wish to identify the size of a small spot on the input plane from the observed Fourier transform image. The spot transmittance, for a 10- $\mu$  spot, and its expected transform in the image plane, are shown in Fig. XI-19. We also wish to cut down on the incident light by placing the spot in an aperture of approximately 50  $\mu$ . Two such apertures are considered. One has a square transmittance function, while the other has a gradual, sinusoidal transmittance. (The transmittances are shown in Fig. XI-21.) The simulation program for these configurations is shown in Fig. XI-20. The field width is 16,384 points. It should be noted that INPUT has the capability of producing a stop after initialization. In fact, STOP is simply an entry to that part of INPUT. The calls to OUTPUT produce the intensity graphs shown in Fig. XI-21. These are plots of the input intensity and are plotted on the print-out sheet. Stops were introduced at the image plane which block the high intensity on the optic axis. Since the Cal Comp routines automatically scale the output graph, this blocking was necessary to expand the size of the tails of the image. The final images from the Cal Comp plotter are shown in Fig. XI-22. It can be seen that the sinusoidal aperture gives a much clearer image which would be more easily recognized as the transform of a 10- $\mu$  spot.

Optical Bench is available from the Biological Image Processing section of C. I. P. G.

J. E. Bowie

#### References

1. A. Vander Lugt, "Operational Notation for the Analysis and Synthesis of Optical Data Processing Systems," Proc. IEEE 54, 1055-1063 (1966).
2. J. Goodman, Introduction to Fourier Optics (McGraw-Hill Publishing Company, New York, 1968).

## (XI. COGNITIVE INFORMATION PROCESSING)

### F. INITIAL STUDIES ON THE ACOUSTIC CORRELATES OF PROSODIC FEATURES FOR A READING MACHINE

#### 1. Introduction

For several years, a reading machine has been under development in the Cognitive Information Processing Group of the Research Laboratory of Electronics. The final purpose of this machine is to take printed text as input and to give as output an intelligible spoken equivalent. To realize these aims, pattern recognition routines have been developed to gather the text information,<sup>1</sup> a grapheme-to-phoneme conversion routine has been written to transform the text information into its phonemic equivalent,<sup>2</sup> a limited sentence parser has been produced to supply structural information,<sup>3</sup> and phonemic synthesis programs have been written to produce the output speech.<sup>4, 5</sup>

Some important problems remain to be solved before the reading machine can really become a viable system. Among these is the problem of correctly using the structural information, made available by the parser, to control the acoustic correlates for the prosodic features (stress, juncture, and intonation) of the speech signal. This is important for two reasons. First, a great deal of information is available through the sentence structure, and needs to be transmitted to the listeners. Second, the correct generation of the acoustic correlates for the various prosodic features tends to make the speech more "natural." This, in turn, increases intelligibility and the "ease" with which a listener can understand what is being spoken.

Since the parser for the reading machine is a phrase-level parser, the problem of the acoustic correlates of stress, intonation, and juncture is currently being approached on a phrase level. The acoustic correlates under study, at present, are fundamental frequency, duration (particularly vowel durations), intensity, and pauses. It is the purpose of this report to describe briefly the tact being used in studying these acoustic correlates of prosodic features as they apply to the reading machine, to discuss in some detail the results of one recent experiment, and to indicate the direction of other work in this area.

#### 2. Problems of Describing the Acoustic Correlates of Prosodic Features in the Reading Machine Environment

Any work associated with the reading machine is subject to several constraints, because of the character and goals of the over-all project itself. One of these constraints is, of course, that all processing must be done automatically, but this alone cannot make the reading machine a usable product. Another constraint is that any rules written for use with the reading machine must also conform to the "real time" goals of the project, and, likewise, since many computer programs must be used in processing a sentence from recognition to speech, any written computer code must be as compact as possible.



Hence, the first thing that must be remembered in writing any rules for the reading machine is that they must function as well as possible and still not require an unreasonable amount of processing time or computer space.

In this context, then, the problem of describing the acoustic correlates for stress, intonation, and juncture can be restated as follows. Given as an input the phonemic strings and word-level stress supplied by the grapheme-to-phoneme conversion routines, and given the phrase-level structure supplied by the parser, find a set of rules for controlling fundamental frequency, phoneme duration, phonemic intensity, pause durations, and phonemic quality which are sufficient to transfer to a listener the available structural information in a way that he perceives as natural. Notice here that there is no requirement to describe all of the acoustic phenomena associated with prosodic features in real speech. Clearly, the mapping from derived constituent structure to the acoustic correlates of prosodic feature in real speech is many to one, being effected by dialectic variation, the idiosyncrasies of individual subjects, and semantic and emotional considerations. What is sought, rather, are rules describing one such mapping that is acceptable from a perceptual viewpoint. There is hence no requirement that the rules reflect any fundamental truths about speech or speech production, but only that they give one legitimate set of output acoustic correlates as described above.

The approach now being used in solving this problem can be stated as follows. First, study the effects of various lexical (word) stress patterns on fundamental frequency, phoneme duration, etc. in words spoken in isolation (one-word phrases or sentences). Then observe these words in various structural positions in phrases, and try to describe systematically the variations in the various acoustic correlates attributable to structure. In this way, we hope that a hierarchy of rules, working from the word level to the phrase level, may be devised.

Before going on to describe one of the initial experiments, a word on the collection of data is in order. First, it is a well-known fact that there are a great many individual differences between the speech produced by different subjects. Hence, if consistent results are desired, it is important to analyze the data from each subject individually, and then look for general patterns in the data. Second, it is also known that not all variations in the various acoustic correlates that are being studied are related to structure. For example, it is known that vowel duration is affected by the following consonant, and, likewise, phonemes in general vary in length because of their syllabic position. Hence, where the effects of structure alone are sought, it is desirable, whenever possible, to describe the other known effects and "normalize" them out. Third, there is the problem of describing "phonemic duration" in a systematic way. A "phoneme" is, in fact, an abstract linguistic entity, and does not inherently have a "duration." Hence, a set of rules for consistently placing phonemic boundaries were devised which were meaningful from the point of view of the speech synthesis routines. These rules, however, do not

(XI. COGNITIVE INFORMATION PROCESSING)

include glide-vowel junctures, and, for this reason, the problem of glide durations is being postponed until the other phonemic durations are better understood.

3. Design of an Initial Experiment

The first experiment was designed to study the effect of structure on fundamental frequency and phoneme durations in single-syllable words. All words used in this experiment were chosen to have only one vowel, in particular /æ/. The reason for this is that previous experiments have shown<sup>6</sup> that different vowels have different "inherent durations." Hence, using only one vowel prevents the need to normalize the durations for different vowels.

Single Words

- |          |          |
|----------|----------|
| 1. sad   | 26. dad  |
| 2. sack  | 27. damp |
| 3. fat   | 28. gad  |
| 4. fad   | 29. gap  |
| 5. bat   | 30. mam  |
| 6. back  | 31. tag  |
| 7. bad   | 32. fag  |
| 8. bap   | 33. pad  |
| 9. bam   | 34. pass |
| 10. ban  | 35. pack |
| 11. dam  | 36. sap  |
| 12. sam  | 37. tap  |
| 13. nam  | 38. tan  |
| 14. pan  | 39. tat  |
| 15. fan  | 40. cap  |
| 16. can  | 41. cat  |
| 17. sass | 42. cab  |
| 18. bass | 43. jazz |
| 19. mass | 44. zap  |
| 20. man  | 45. zam  |
| 21. mad  | 46. pad  |
| 22. nab  | 47. maz  |
| 23. nap  | 48. chad |
| 24. bab  | 49. shad |
| 25. nack | 50. gas  |

Fig. XI-23. Test words for the initial experiment.

(XI. COGNITIVE INFORMATION PROCESSING)

Three subjects, all speakers of the "General American" dialect, were used for this experiment. They were each given lists of words and phrases and were told to read them "clearly and carefully!" Recordings were made of these utterances, and each were analyzed by using a Kay Sonograph. For each utterance, a narrow-band and a wideband spectrogram were made, and from these, phoneme durations were tabulated and fundamental frequency contours were plotted.

The material given to the subjects to read consisted of two lists. The first list (see Fig. XI-23) consisted of fifty single-syllable words all having the vowel /æ/. The purpose of this list was to determine individual differences for the various subjects in their

- Group 1
- [ [X]<sub>NVA</sub>[[Y][Z]]<sub>NVA</sub> ]
1. Sad sad sack
  2. Fat cat man
  3. Mad bad man
  4. Damp mad man
  5. Tan back sack
- Group 2
- NVA [ NVA [ [X][Y] ]<sub>NVA</sub> [Z] ]<sub>NVA</sub>
1. Mad man jazz
  2. Sad sad sack
  3. Fat cat man
  4. Damp mad man
  5. Back sack gap
- Group 3
- [ [X][[Y][Z]] ]
1. Sad sad sack
  2. Sad fat bat
  3. Bad tan cat
  4. Damp tan cap
  5. Mad bad man

Fig. XI-24. Three-word structural groups.

## (XI. COGNITIVE INFORMATION PROCESSING)

inherent vowel durations and to observe the effect of the following consonant on vowel durations.

The second list consisted of three groups of 5 three-word phrases (see Fig. XI-24). The groups were chosen so as to be representative of the three most common stress patterns that might be recognized by the parser, in particular the 213, 132, and 231 stress patterns (stress numbers as used by Halle and Chomsky<sup>7</sup>). A fourth common stress pattern, the 321 pattern, was not used, since, in general, it required adverbs. Adverbs with one syllable having the vowel /æ/ are not readily available in English.

### 4. Results

The data available from this experiment were analyzed in two ways. First, for the single-word lists, the vowel durations were plotted as a function of the following consonants for the individual subjects. The results for Subject 1 are shown in Fig. XI-25. Two things are immediately evident from these data. First, the vowel durations seem to be divided into 5 different groups, in particular vowels before voiced-stop consonants, unvoiced-stop consonants, voiced fricatives, unvoiced fricatives, and nasals. Second, the variation of vowel durations within these groups is fairly small. These effects were observable for all three subjects.

Another interesting result is observed if the average vowel durations before the various consonant groups are compared (see Fig. XI-26). Here it is clear that, although all three subjects showed the grouping effect, the actual vowel durations for the individual subjects are quite different. This serves to point out that there exist considerable individual differences, even among speakers of the same dialect.

The data for the durations of the vowels in the three word groups were analyzed in several ways. It was desired that the effects on vowel durations by following consonants and sentence position (prepausal lengthening) be removed so as to be able to observe the results of structure alone. To do this, we assumed that vowel duration could be considered as a function of structural position, location in the phrase (last vowel or not last vowel), and the following consonant group. Several models for describing vowel duration were tried, and in each the criterion for goodness was how well they served to predict the observed data. The model that appeared to give the best results was the one describing vowel duration, V.D. as

$$\underline{V.D.} = \underline{B.D.} + \underline{X} \cdot \underline{B.D.} + \underline{PPL},$$

where B.D. is a basic vowel duration determined by the consonant group following the vowel in question, PPL is a prepausal lengthening predicate that is zero unless the vowel is the last vowel in the utterance and is not followed by an unvoiced stop consonant, and



Fig. XI-25. Raw vowel (æ) duration taken from 50 single-syllable words and plotted as a function of following consonant for Subject 1 of Experiment 2.

Group	Subjects		
	1 (msec)	2 (msec)	3 (msec)
Voiced-stop consonant	290	280	350
Unvoiced-stop consonant	170	220	230
Unvoiced fricative	240	250	290
Voiced fricative	360	360	380
Nasal	190	240	330

Fig. XI-26. Comparative vowel durations (av.) for 3 subjects as a function of the following consonant.

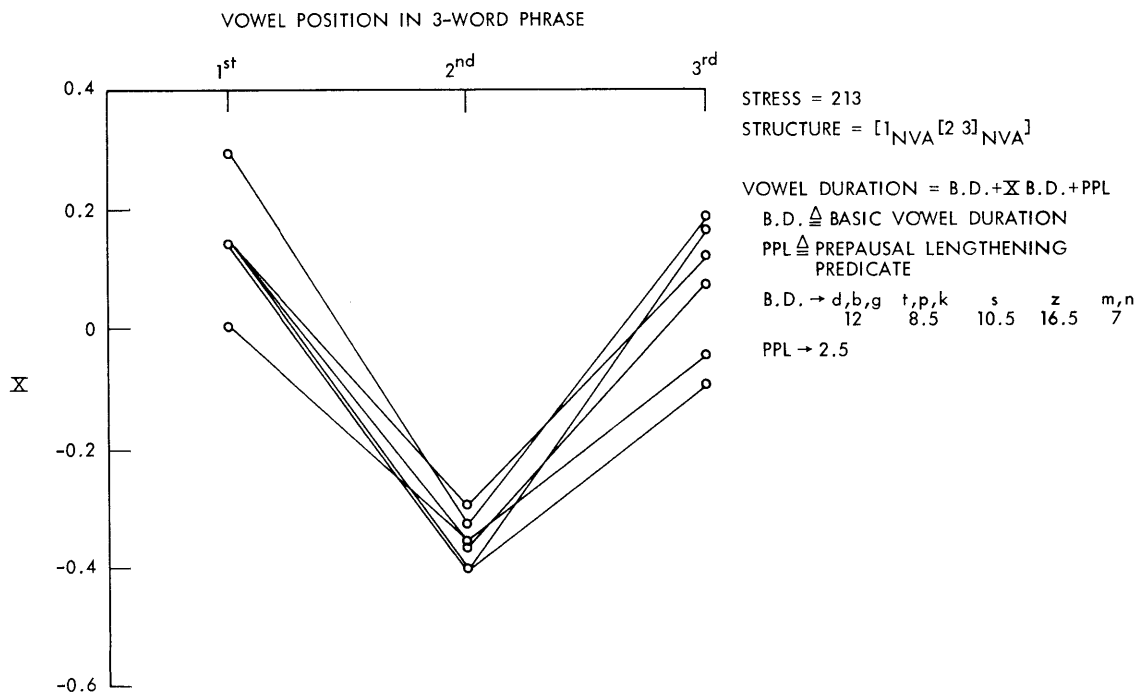


Fig. XI-27. X vs vowel position for 213 stress.  
 (Subject 1, Experiment 2.)

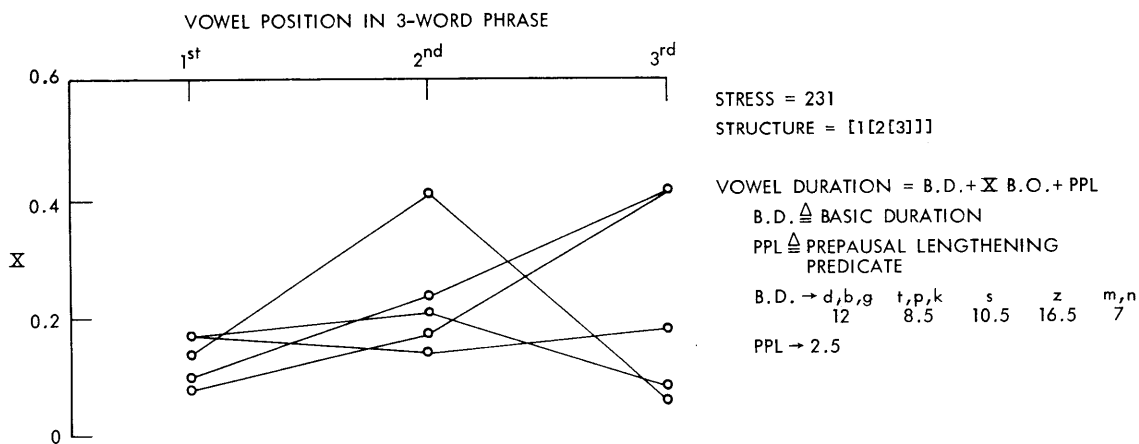


Fig. XI-28. X vs vowel position for 231 stress.  
 (Subject 1, Experiment 2.)

(XI. COGNITIVE INFORMATION PROCESSING)

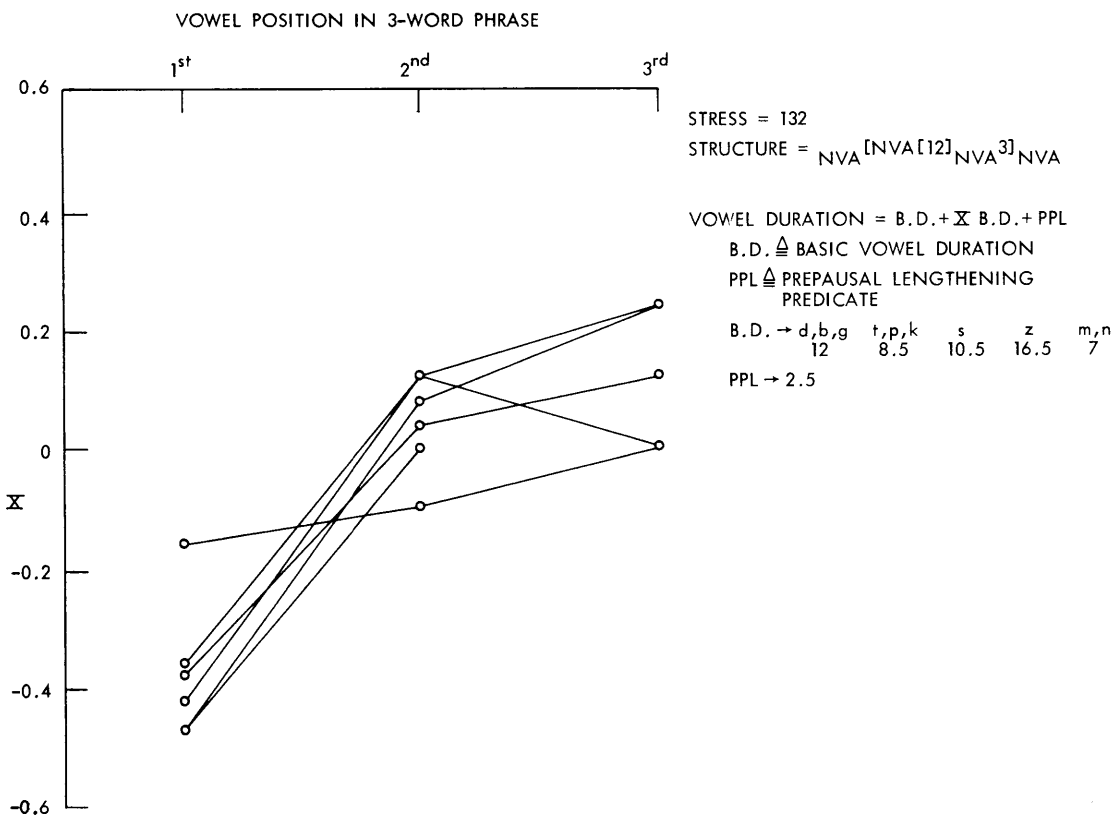


Fig. XI-29.  $\bar{X}$  vs vowel position for 132 stress.  
(Subject 1, Experiment 2.)

constant otherwise, and  $\bar{X}$  is the multiplicative effect of structure. The numbers were chosen for this model so that it described both the single-word data and the three word groups. The criterion for goodness was numbers giving the smallest spread of  $\bar{X}$  values for a particular structural configuration. The resulting  $\bar{X}$  values are plotted for Subject 1 for the three three-word groups in Figs. XI-27, XI-28, and XI-29. Two things are noteworthy about these data and the data for the two other subjects. First, for all subjects, the spread in  $\bar{X}$  values for a given structural configuration was quite small. Second, once more, individual differences were very much in evidence, and although the data of all three subjects conformed very well to the model, the B.D., PPL, and X's observed were quite different for different subjects.

A study was also made of the various fundamental frequency contours for the three word groups. The results of this study can be summarized as follows. For each structural group of 5 utterances, the fundamental frequency contours observed for each subject were very similar both in shape and frequency range. There were noticeable differences, however, between the different structural groups and the different subjects.

4. Discussion

A great deal may be said both for and against the results just presented. In their favor is the fact that they do suggest a model that relates numerically the observed vowel durations to the structure via the  $\underline{X}$  values. This was one of the purposes of the experiment. Likewise, further study shows that there is a great similarity among the fundamental frequency contours and among the  $\underline{X}$  values observed for the two-word compounds. This is also observable for several other structural units, in particular a one-stressed word not part of a compound, adjectives before other adjectives or one-stressed words, adjectives before a compound, and the third word in a three-word compound. These facts, together with the results of other experiments seem to suggest that a simple model may be used to describe both the fundamental-frequency results and the durational results.

Clearly, this experiment leaves a great many questions unanswered. First, there is no way to know from this single experiment whether the observed results hold for any vowel but /æ/. To test this, a second experiment was performed using the vowel /I/, a much shorter vowel than /æ/. The results of the second experiment were such as to strongly support the results of the first. The implication here is that the observed durational phenomena probably hold throughout the vowels.

Second, there are no data on multiple-syllable words. Clearly, a study must be made on multiple-syllable words to discover whether the X values or any of the other consistent results have any general relevance. Likewise, studies must be made which will expand the class of structures involved to cover all structures available through the parser. Studies of this kind are now in progress.

Another point of relevance concerns the observable variations among subjects. Throughout the experiment, we found that the subjects were internally very consistent, but showed great variation from subject to subject. These individual variations are generally observed phenomena in speech research, but they pose great problems in the writing of rules. It is probably best to follow the example of a single subject, but only extensive perceptual experiments will show whether this is correct.

One more area of intensive study associated with this experiment was an attempt to correlate the observed acoustic results with the stress numbers suggested by Halle and Chomsky.<sup>7</sup> A great many functional relationships were tried here, but none proved successful. This in no way detracts from current linguistic theory, since the stress numbers are really defined in a perceptual rather than physical sense. It does imply that whatever the relationship between the acoustic signal and the perceived stress may be, it is not simple. If one were to characterize the results on this point, one would have to say that it appears that the acoustic correlates observed relate directly to structure rather than to "stress numbers."



## 5. Summary

Experiments have been performed to study the acoustic correlates of prosodic features. From these experiments, a basis for a possible model has been observed, although the model is very incomplete. The model may and may not have general relevance to the field of speech communication, but it has the advantage of describing quite well a body of observed experimental results.

T. P. Barnwell III

## References

1. J. K. Clemens, "Optical Character Recognition for Reading Machine Applications," Quarterly Progress Report No. 79, Research Laboratory of Electronics, M. I. T., October 15, 1965, pp. 219-227.
2. F. F. Lee, "A Study of Grapheme to Phoneme Translation of English," Ph. D. Thesis, M. I. T., 1965.
3. J. Allen, "A Study of the Specification of Prosodic Features of Speech from a Grammatical Analysis of Printed Text," Ph. D. Thesis, M. I. T., 1968.
4. T. P. Barnwell, "An Algorithm for the Transformation of a Sentence Representation into Control Parameters for a Speech Synthesizer," S. M. Thesis, M. I. T., 1967.
5. J. N. Holmes, I. G. Mattingly, and J. N. Shearme, "Speech Synthesis by Rule," *Language & Speech* 7, 127-143 (1964).
6. A. S. House, "On Vowel Duration in English," *J. Acoust. Soc. Am.* 33, 1174-1178 (1961).
7. N. Chomsky and M. Halle, Sound Patterns of English (Harper and Row Publishers, Inc., New York, 1968).

## G. ASYMPTOTIC RATE FOR FOURIER SERIES EXPANSION

If a sequence of Gaussian random variables  $\{x_i\}$ ,  $0 \leq i \leq N-1$  is encoded with respect to a mean-square-error criterion, the minimum rate per block for a mean-square error per sample  $d$  is given by

$$R_b(d) = \frac{1}{2} \log_2 \frac{|R_x|}{d^N} \quad \text{for } d < \min_i \lambda_i,$$

where the  $\lambda_i$  are the eigenvalues of  $R_x$ , the correlation matrix of the random variables. Thus the average rate per element  $R(d)$  is given by

$$R(d) = \frac{1}{N} R_b(d) = \frac{1}{2} \log_2 \frac{|R_x|^{1/N}}{d}.$$

Now  $|R_x|^{1/N}$ , the  $N^{\text{th}}$  root of the determinant of  $R_x$  is just the geometric mean of the eigenvalues of  $R_x$ . Letting  $P$  be the orthonormal matrix that diagonalizes  $R_x$ , we have

$$\begin{pmatrix} \lambda_0 & & 0 \\ & \ddots & \\ 0 & & \lambda_{N-1} \end{pmatrix} = \mathbf{P}\mathbf{R}_x\mathbf{P}^T.$$

Furthermore, setting  $y = \mathbf{P}x$ , where  $x^T = (x_0, \dots, x_{N-1})$ , we have  $E(y_i^2) = \lambda_i$ . Thus, the  $y_i$  can be coded independently, by using  $\frac{1}{2} \log_2 \frac{\lambda_i}{d}$  bits/ $y_i$  to yield the minimum average rate.

Next, let  $y = \frac{1}{\sqrt{N}} \mathbf{F}x$ , where  $\mathbf{F}$  is the Discrete Fourier Transform (DFT) matrix. Thus  $\mathbf{R}_y = \frac{1}{N} \mathbf{F}\mathbf{R}_x\mathbf{F}^T$ . If the  $y_i$  are encoded independently, the average rate is

$$R(d) = \frac{1}{2} \log_2 \frac{(\pi a_i)^{1/N}}{d},$$

where  $a_i \triangleq E(|y_i|^2) = (\mathbf{R}_y)_{ii}$ , the  $i^{\text{th}}$  diagonal element in  $\mathbf{R}_y$ . Clearly,  $R_F(d) \geq R(d)$ ; otherwise, there would be a method that is better than the optimum method. Thus

$$R_F(d) - R(d) = \frac{1}{2} \log_2 \left[ \frac{(\pi a_i)^{1/N}}{(\pi \lambda_i)^{1/N}} \right]$$

is the increase in rate when the Fourier transform is used instead of the optimum  $\mathbf{P}$ . The rest of this report will show that if the Gaussian sequence in question is stationary and Markov of degree  $K$ , then

$$\frac{(\pi a_i)^{1/N}}{(\pi \lambda_i)^{1/N}} \rightarrow 1 \quad \text{as } N \rightarrow \infty \quad (1)$$

so that the DFT is asymptotically optimum with respect to mean-square error for average rate of transmission.

There are essentially three parts to the proof. Part I establishes that the off-diagonal terms of  $\mathbf{R}_y$  vanish as  $1/N$ . Part II establishes that  $\|\mathbf{R}_{y_i}^{-1}\| \leq M < \infty$ , independent of  $N$  and  $i$ ,  $i = 0, \dots, N-1$ , where  $\mathbf{R}_{y_i}$  is the correlation matrix for  $y_i, y_{i+1}, \dots, y_{N-1}$ . Part III uses the result of Part II to establish Eq. 1.

### Part I

$$\text{Let } y = \frac{1}{\sqrt{N}} \mathbf{F}x, \text{ where } f_{k\ell} \triangleq \exp \frac{j2\pi k\ell}{N}.$$

Assume  $\sum \nu R_\nu \triangleq \frac{C}{2} < \infty$ . Then

$$\begin{aligned} a_{ij} &\triangleq E[y_i y_j^*] = \frac{1}{N} \sum_n \sum_m f_{in} f_{jm}^* R_{n-m} \\ &= \frac{1}{N} \sum_{n=0}^{N-1} f_{in} f_{jn}^* \left( \sum_{\nu=-n}^{N-1-n} f_{j\nu}^* R_\nu \right) \\ &= R_0 \delta_{ij} + \frac{2}{N} \operatorname{Re} \left[ \sum_{\nu=1}^{N-1} f_{j\nu}^* R_\nu \left( \sum_{n=0}^{N-1-\nu} f_{in} f_{jn}^* \right) \right]. \end{aligned}$$

Now

$$\sum_{n=0}^{N-1} f_{in} f_{jn}^* = 0 = \sum_{n=0}^{N-1-\nu} (f_{in} f_{jn}^*) + \sum_{n=N-\nu}^{N-1} (f_{in} f_{jn}^*) \quad i \neq j,$$

and so

$$\left| \sum_{n=0}^{N-1-\nu} f_{in} f_{jn}^* \right| \leq \sum_{n=N-\nu}^{N-1} 1 = \nu.$$

Thus

$$i \neq j \implies |a_{ij}| \leq \frac{2}{N} \sum \nu R_\nu \leq \frac{C}{N}.$$

The continuous-time form of this result is due to Root and Pritcher.<sup>1</sup>

## Part II

Lemma 1. For a stationary Markov K sequence  $\{x_i\}$ ,  $0 \leq i \leq N-1$   
 $\|R_x^{-1}\| \leq M < \infty$ , where M is independent of N.

Proof: Let

$$\begin{aligned} y_0 &= x_0 \\ y_1 &= x_1 - r_{10} x_0 \\ y_2 &= x_2 - r_{21} x_1 - r_{20} x_0 \\ &\vdots \\ y_k &= x_k - r_{k,k-1} x_{k-1} - \dots - r_{k0} x_0 \\ &\vdots \\ y_i &= x_i - r_{i,k-1} x_{i-k+1} - \dots - r_{i0} x_0, \quad N-1 \geq i \geq k, \end{aligned}$$

## (XI. COGNITIVE INFORMATION PROCESSING)

where the  $r_{ij}$  are chosen to make the  $y_i$  uncorrelated. Writing this in matrix form, we have

$$y = Px$$

or

$$R_x^{-1} = P^T R_y^{-1} P.$$

Thus

$$x^T R_x^{-1} x = (Px)^T R_y^{-1} Px.$$

Now

$$\|x\|^2 = \sum x_i^2 \leq 1 \implies \max_i y_i^2 \leq \max \left[ \left(1 + |r_{k,k-1}| + \dots + |r_{k0}|\right)^2, \right. \\ \left. (1 + |r_{k-1,k-2}| + \dots + |r_{k-1,0}|)^2, \dots, (1 + |r_{10}|)^2, 1 \right],$$

and so

$$\|R_x^{-1}\| = \sup_{\|x\|=1} (x^T R_x^{-1} x) \leq \frac{\left(1 + |r_{k,k-1}| + \dots + |r_{k0}|\right)^2}{\min \eta_i},$$

where  $\eta_i = E(y_i^2)$ . In the absence of a deterministic or singular component  $\min_{0 \leq i \leq k} \eta_i > 0$ .

(Note that  $\eta_i = \eta_k$  for  $i \geq k$  and furthermore,  $r_{ij}$  and  $\eta_i$  are not functions of  $N$ .) Thus

$$\|R_x^{-1}\| \leq M < \infty,$$

independent of  $N$ .

Lemma 2. Let  $R_{x_i}$  be the correlation matrix for  $x_1, \dots, x_{N-1}$ . Then

$$\|R_{x_i}^{-1}\| \leq \|R_x^{-1}\|.$$

Proof: Since  $\|R_x^{-1}\| = \frac{1}{\min_i \lambda_i}$ , where the  $\lambda_i$  are the eigenvalues of  $R_x$ , it is both necessary and sufficient that the minimum eigenvalue of  $R_{x_i}$  be greater than or equal to the minimum eigenvalue of  $R_x$ . Assume, therefore, that for some  $\lambda > 0$ ,  $R_{x_i} - \lambda I$  is singular and, furthermore, that (i) this is the smallest such  $\lambda$  for which the assumption of singularity is true, and (ii)  $\lambda < \min_i \lambda_i$ . Then we can write  $x_j = \hat{x}_j + n_j$  as an orthogonal

decomposition ( $j = i, \dots, N-1$ ), where  $R_{\hat{x}_i} = R_{x_i} - \lambda I$  is non-negative definite, and  $R_{n_i} = \lambda I$ . Since  $R_{\hat{x}_i}$  is singular, the  $\{\hat{x}_j\}$  have a deterministic component with probability 1; that is, they are linearly dependent random variables.

Next consider  $x_j = \hat{x}_j + n_j$   $j = 0, \dots, N-1$ , where the  $n_j$  have correlation matrix  $\lambda I$  and are independent of the  $\hat{x}_j$ , having correlation matrix  $R_x - \lambda I$ . Note that this is valid, since we are assuming that  $\lambda < \min_i \lambda_i$ . Thus  $R_{\hat{x}} = R_x - \lambda I$  is non-negative definite. In fact, it is positive definite and hence nonsingular. But it contains a deterministic component in the range  $j > i$ ; thus, it must be singular. Contradiction: Hence  $\lambda \geq \min_i \lambda_i$ , so  $\lambda_{\min} \geq \min_i \lambda_i$ . Q. E. D.

Corollary.  $\|R_{y_i}^{-1}\| \leq M < \infty$   $i = 0, \dots, N-1$

with  $M$  independent of  $N$ , where  $y$  is the DFT of  $x$ .

Proof:  $R_y = \frac{1}{N} F R_x F^\dagger$  (where the dagger indicates conjugate transpose).

Hence

$$R_y^{-1} = \frac{1}{N} F^\dagger R_x^{-1} F,$$

and so

$$y^\dagger R_y y = \frac{1}{N} (y^\dagger F^\dagger) R_x^{-1} (F y).$$

Now

$$\|x\| \leq 1 \iff \left\| \frac{1}{\sqrt{N}} F y \right\| \leq 1.$$

Hence

$$\begin{aligned} \sup_{\|y\| \leq 1} (y^\dagger R_y y) &= \sup_{\|x\| \leq 1} x^\dagger R_x^{-1} x = \|R_x^{-1}\| \\ &= \|R_y^{-1}\|. \end{aligned}$$

Thus

$$\|R_y^{-1}\| \leq M < \infty,$$

independent of  $N$ . So, by Lemma 2,

$$\|R_{y_i}^{-1}\| \leq \|R_y^{-1}\| \leq M < \infty \quad i = 0, \dots, N-1,$$

independent of  $N$ .

Q. E. D.

Part III

Let  $\sigma^2 \triangleq \lim_{N \rightarrow \infty} |R_x|^{1/N}$ . The limit exists, since we are assuming that  $\{x_i\}$  is stationary Markov. Let  $R_y \triangleq \frac{1}{N} FR_x F^\dagger$ , and  $a_i^{(N)} \triangleq (R_y)_{ii}$   $i = 0, \dots, N-1$

Theorem.  $\left( \prod_{i=0}^{N-1} a_i^{(N)} \right)^{1/N} \rightarrow \sigma^2$  as  $N \rightarrow \infty$ .

Proof: We write  $R_y^{(N)}$  instead of  $R_y$  to denote the dependence of  $R_y$  on  $N$ . Consider

$$R_y^{(N)} \xi^{(N)} = C^{(N)} \triangleq \begin{pmatrix} a_0^{(N)} \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad (2)$$

and

$$R_y^{(N)} \xi = C_1^{(N)},$$

where

$$\xi = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}. \quad (3)$$

Then

$$C_1^{(N)} = \begin{pmatrix} a_0^{(N)} \\ 0\left(\frac{1}{N}\right) \\ \vdots \\ 0\left(\frac{1}{N}\right) \end{pmatrix},$$

where  $0\left(\frac{1}{N}\right)$  denotes terms that tend to zero as  $1/N$ . This follows from the results of Part I.

Combining (2) and (3), we have

$$R_y^{(N)} [\xi^{(N)} - \xi] = C^{(N)} - C_1^{(N)}.$$

Now  $C^{(N)} - C_1^{(N)} \rightarrow 0$  in  $\ell^2$ . From the results of Part II, we know that  $\| [R_y^{(N)}]^{-1} \| \leq M < \infty$ ; hence,  $\xi^{(N)} \rightarrow \xi$  in  $\ell^2$ . Thus  $\xi_o^{(N)} \rightarrow \xi_o = 1$  as  $N \rightarrow \infty$ .

By Cramér's rule,

$$1 \rightarrow \xi_o^{(N)} = \frac{\begin{vmatrix} a_o^{(N)} & | & \\ \hline 0 & | & R_{y_1}^{(N)} \\ \vdots & | & \\ 0 & | & \end{vmatrix}}{|R_y^{(N)}|} = a_o^{(N)} \frac{|R_{y_1}^{(N)}|}{|R_y^{(N)}|}.$$

Now repeat the procedure, with  $R_{y_1}^{(N)}$  instead of  $R_y^{(N)}$  to yield

$$a_1^{(N)} \frac{|R_{y_2}^{(N)}|}{|R_{y_1}^{(N)}|} \rightarrow 1.$$

Continue this procedure until  $N-1$  and multiply all terms. Cancellation of  $|R_{y_i}^{(N)}|$  occurs for  $i = 1, \dots, N-2$ , and we get

$$\frac{\prod_{i=0}^{N-1} a_i^{(N)}}{|R_x^{(N)}|} \leq \prod_{i=0}^{N-1} \left( 1 + \frac{MC}{\sqrt{N}} \right),$$

where  $M$  and  $C$  are obtained from

$$\| R_{y_i}^{-1} \| \leq M,$$

independent of  $i$  and  $N$ , and

$$|(R_y^{(N)})_{ij}| \leq \frac{C}{N} \quad \text{for } i \neq j,$$

independent of  $i, j, N$ .

Taking  $N^{\text{th}}$  roots, we get

$$\begin{aligned} \frac{\left[ \prod_{i=0}^{N-1} a_i^{(N)} \right]^{1/N}}{|R_x|^{1/N}} &\leq \left( \prod_{i=0}^{N-1} \left( 1 + \frac{MC}{\sqrt{N}} \right) \right)^{1/N} \\ &= 1 + \frac{MC}{\sqrt{N}}. \end{aligned}$$

(XI. COGNITIVE INFORMATION PROCESSING)

Thus

$$\left[ \prod_{i=0}^{N-1} a_i^{(N)} \right]^{1/N} \rightarrow \lim_{N \rightarrow \infty} |R_x|^{1/N},$$

with an error term of the order of (or less than)  $\sigma^2 MC/\sqrt{N}$ . In terms of rate, the error would be  $\log(1+MC/\sqrt{N}) \leq MC/\sqrt{N}$ .

Q. E. D.

J. W. Woods

References

1. W. L. Root and T. S. Pritcher, "On the Fourier-Series Expansion of Random Functions," *Ann. Math. Statist.* 26, 313-318 (June 1955).