

III. Information Theory and Image Coding

Information theory provides us with mathematical tools to determine the channel capacity required to transmit a certain image or class of images over the channel. Two results from this theory give us a basis for image coding, they are [3]

$$0 \leq H(U) \leq \log_2(A) \quad (1)$$

and

$$H(U_N) \leq NH(U) \quad (2)$$

where H is the entropy function,

$$H(U) = - \sum_i p_i \log_2(p_i) \quad (3)$$

and U is a scalar source, U_N is an N dimensional source, p_i is the probability of source symbol u_i , and A is a source with uniform distribution. The first equation states that compression of the image data can be achieved if the statistical distribution of the data is not uniform, while the second equation states that further compression can be achieved if the data is dependant or correlated.

The theory behind these two equations together with the rate distortion theory form the basis of all image coding algorithms. Basically the above two equations determine that the output of the coder should consist of an independent uniformly distributed sequence while the rate distortion function determines the minimum distortion

at which this can be achieved for a fixed rate. The grey scale histogram of an image is normally not uniform, i.e. figures 4 and 5 show the histograms of the test images. By applying equation (1), the average bit rate can be computed for this image. Most images also contain a fair amount of correlation between the image pels. If a method could be devised to remove this correlation, equation (2) predicts a lowering in the average bit rate. In transform image coding the correlation is reduced by using an energy compacting transform. This property of the transform will be discussed in more detail in Section VII.

Information theory also provides the means by which to analyze, in a mathematical sense, the transmission of sources where some distortion of the source is acceptable. This is done in the framework of the *rate distortion theory*. A brief overview of the origin of the rate distortion function is given in the rest of this section. Most of the overview is based on the references [3] and [18].

The problem addressed by the rate-distortion theory is the minimisation of the channel capacity requirement while holding the average distortion at or below an acceptable level. More specific, the rate distortion function $R(D)$ is the minimum value of the mutual information $I(U,V)$ for a given distortion level D [3]. By keeping the rate lower than the channel capacity C , i.e $R(D) < C$, the possibility of obtaining distortion D is ensured.

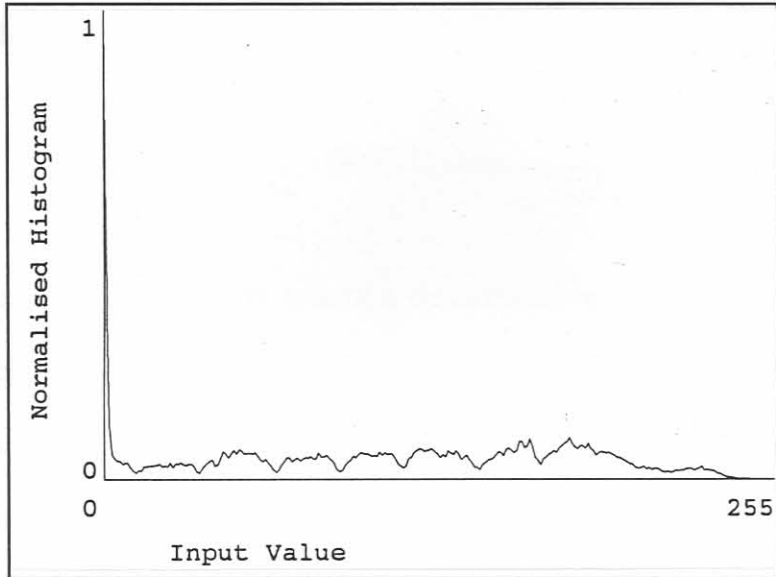


Figure 4 Histogram of the test image GIRL

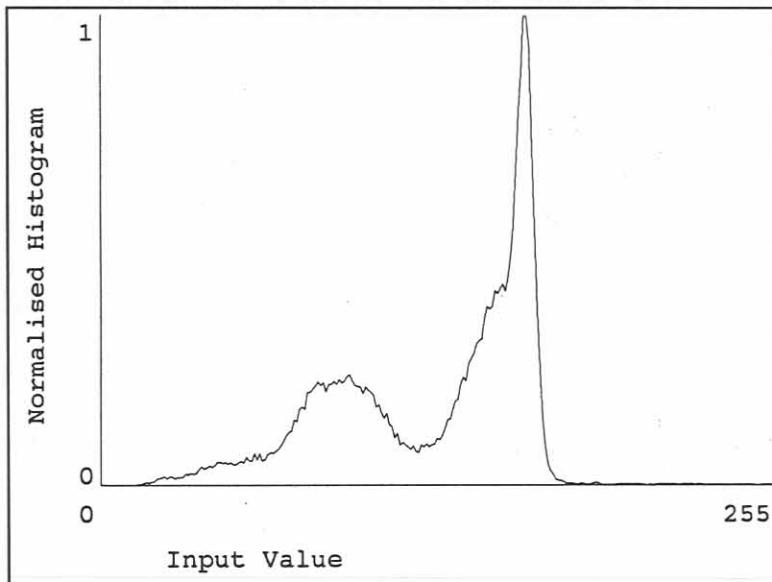


Figure 5 Histogram of the Image ROAD

The information theoretic measure of information transmitted is the average mutual information between U and V , and is defined for a block length N as:

$$I_N(U, V) = \sum_i \sum_j P(U_i) Q(V_j | U_i) \log \frac{Q(V_j | U_i)}{\sum_k P(U_k) Q(V_j | U_k)} \quad (4)$$

Each block is described by one of a denumerable set of messages $\{U_i\}$ with probability $P(U_i)$. Any given system, i.e. channel, is described mathematically by the conditional probability $Q(V_j | U_i)$ of message V_j being output by the decoder for a given source output U_i . The mutual information may also be written as

$$I_N(U, V) = H_N(U) - H_N(U | V) \quad (5)$$

where $H_N(U | V)$ is the entropy of the source given the observed decoder output. In other words the mutual information is equal to the entropy of the source minus the entropy of the source given the decoder output V .

If distortion is introduced then the decoder output only contains statistical information about U and as a result thereof $I_N(U, V)$ decreases. In the worst case V contains no information about U , so that $H_N(U | V) = H_N(U)$ and $I_N(U, V) = 0$.

For the time discrete continuous-amplitude sources, the mutual information and average distortion functions are defined with integrals in place of the summations. Given the distortion $d(u,v)$ between u and v , the rate distortion results can be summarised as follows for block messages:

$$R(D^*) = \lim_{N \rightarrow \infty} \frac{1}{N} \inf_{q: D(q) \leq D^*} I_N(U, V) \quad (6)$$

$$I_N(U, V) = \int p(u) q(v|u) \log \frac{q(v|u)}{t(v)} dv du \quad (7)$$

$$D(q) = \int d(u, v) p(u) q(v|u) dv du \quad (8)$$

$$t(v) = \int p(u) q(v|u) du \quad (9)$$

It is very difficult to solve these equations and it is normally only done for homogeneous, isotropic, Gaussian sources with a mean square error criterion. The rate for this particular Gaussian source is given by [18]

$$R(D^*) = \frac{1}{2} \log \frac{\sigma^2}{D^*} \quad \text{for } \sigma^2 > D^* \quad (10)$$

The rate distortion theory is used in a later section to simulate the optimal encoding of images.

The next section looks at the statistical models used to facilitate a mathematical treatment of image coding.

IV. Image Statistical Models

Images are sometimes represented by simple stochastic models in order to develop useful algorithms or to compare the performance of various processes on an image mathematically. A stochastic process can be completely described by its joint probability density [1]. In general, high-order joint probability densities of images are usually not known, nor are they easily modeled [4]. For practical reasons the images are characterised by their mean and covariance functions.

A common model used, for natural images, is that of the two dimensional, stationary, first-order Markov process [1]. If f_{ij} represents the picture brightness at the point (i, j) , then for this process the autocorrelation function may be written

$$R(m, n) = E [f_{ij}, f_{i+m, j+n}] = \rho^{|m-n|} \quad (11)$$

where $0 < \rho < 1$

and zero mean is normally assumed $E[f_{ij}] = 0$. The assumption of zero mean is nonessential, since the mean can always be easily computed and subtracted if necessary to obtain a zero mean image.

To test the validity of the assumption the correlation of the two test images have been computed, horizontally and vertically, and is compared with the theoretical Markov model in figures 6 and 7. For large block sizes the images follow the model for small shifts from the origin.

V. The Human Visual System

The human visual system (HVS) is a very complex system. At this time

there is a lot of research going on in the area of the HVS. The HVS is a very complex system and it is not yet fully understood. However, there are some basic principles that can be used to design a visual coding system. One of the most important principles is the concept of spatial correlation. This is the degree to which the intensity of a pixel is related to the intensity of its neighbors. The higher the spatial correlation, the more similar the pixels are to each other. This is why images with high spatial correlation are easier to compress than images with low spatial correlation.

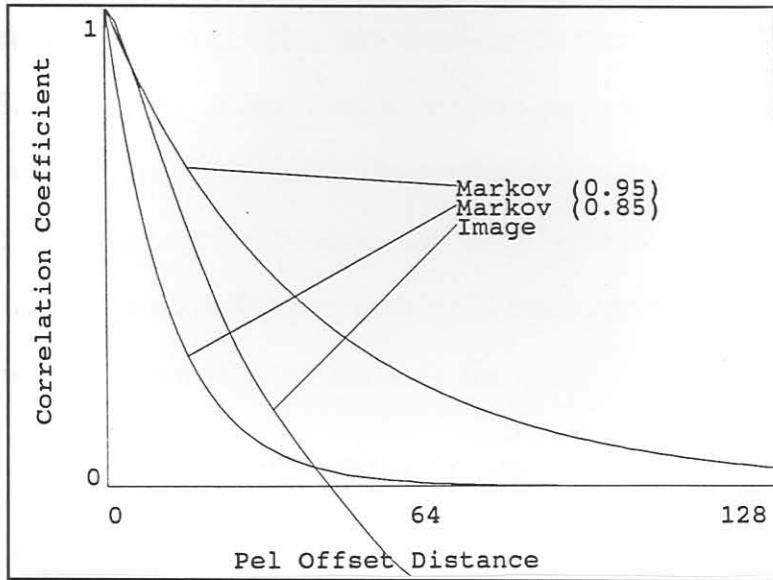


Figure 6 Spatial correlation for image GIRL

They also derived a nonlinear transfer function of the HVS given by

$$f(x) = \frac{1}{1 + e^{-x}}$$

This function is used to model the response of the HVS to different intensities of light.

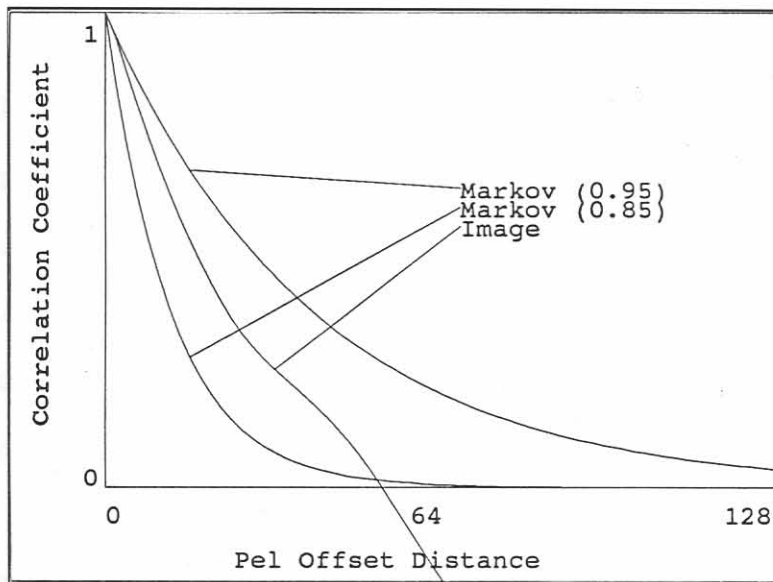


Figure 7 Spatial correlation for image ROAD

Correlation function for the image.

V. The Human Visual System

The human visual system (HVS) is a very complex system. At this time there is no accurate model that can simulate the visual interpretation of the HVS. However, a few simple models have been presented by various authors [16,21,22,23,27]. These empirical results were basically achieved by doing subjective threshold judging tests.

The results of the judging tests by Mannos and Sakrison [16] lead to a transfer function of approximately the form

$$A(f_r) \approx 2.6 (0.0192 + 0.114 f_r) \exp [-(0.114 f_r)^{1.1}] \quad (12)$$

where f_r is the radial frequency in cycles/degree. The equation has a peak value of one at 8 cycles/degree and diminishes at 64 cycles/degree. They also derived a nonlinear transfer function for the HVS given by

$$f(u) = u^{0.33} \quad (13)$$

Plots of the two equations are given in figure 8.

A. Rate Distortion Simulation

To determine the best picture that any coding scheme may produce, the same rate distortion simulation procedure was followed that led to the generation of the empirical equations (12) and (13) [16]. To calculate the rate distortion function one must specify both the *distortion measure* and the *probability distribution* of the source. At present we are unable to specify the probability distribution of an image source. The best we can do is to specify the mean and correlation function for the image.

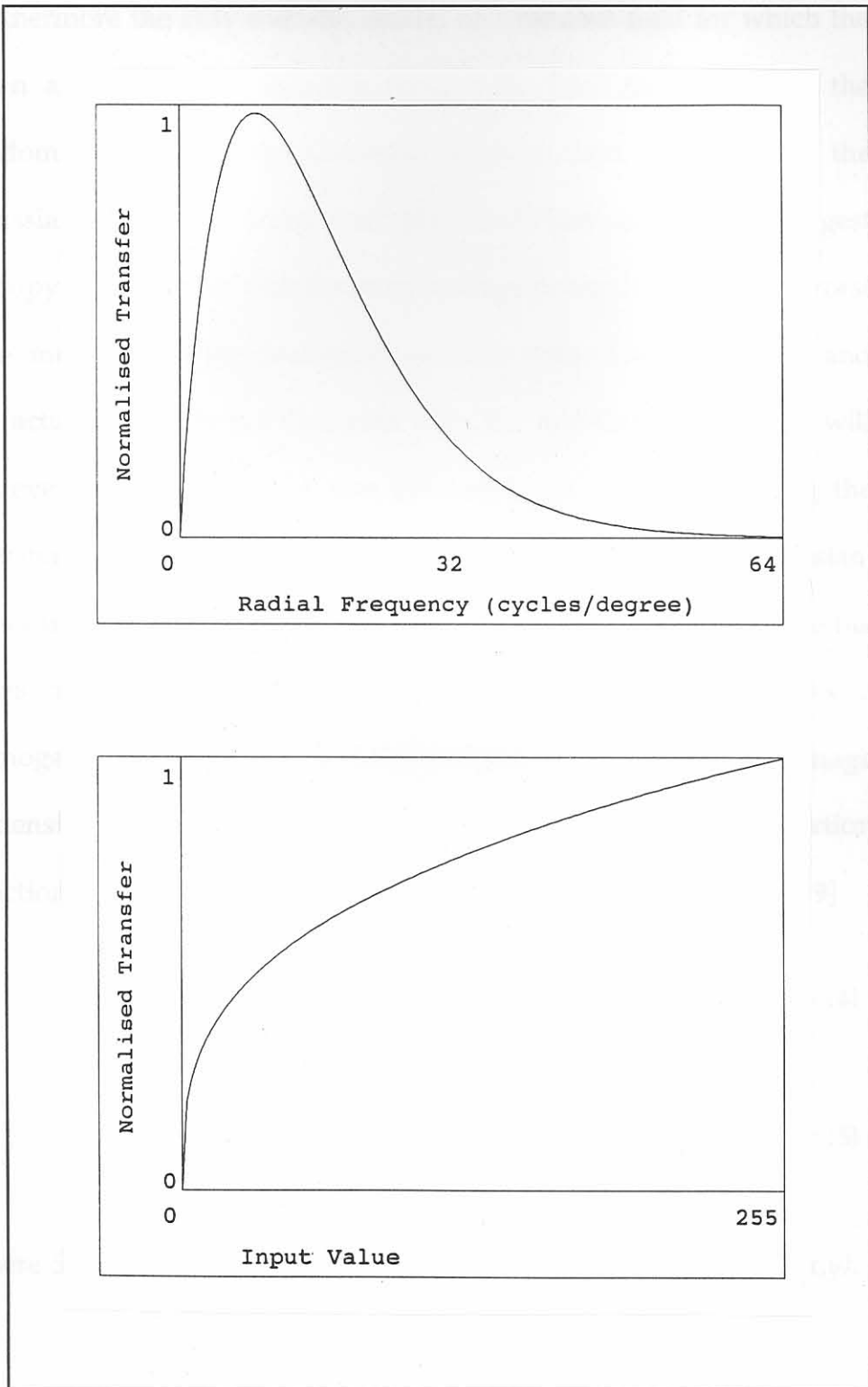


Figure 8 Transfer functions of the HVS

Furthermore the only tractable model of a random field for which the mean and correlation function specify the joint distribution of the random variables is the Gaussian random field. However, the Gaussian source is a "worst case" hypothesis because it has the largest entropy of all sources with the same average power [19,27]. This "worst case" means that if the simulations are done with a Gaussian model and the actual model is not Gaussian then the non-Gaussian source will achieve better results than was estimated. In transform coding the distribution of the errors in the spatial domain will be nearly Gaussian. This can be seen from the central limit theorem, i.e. the errors are the sums of independent errors in the transform coefficient. For a homogeneous (invariant to a shift in origin) Gaussian field with image dimension large compared to correlation distance, the rate distortion function is given in parametrical form by the pair of equations [19]

$$R(\mu) = \frac{1}{2} \iint_{S(f_x, f_y) > \mu} \log_2 [S(f_x, f_y) / \mu] df_x df_y \quad (14)$$

$$d^*(\mu) = \iint_{-\infty}^{\infty} \min[S(f_x, f_y), \mu] df_x df_y \quad (15)$$

where $S(f_x, f_y)$ is the power spectral density of the input image $V(x, y)$.

If $V_{k,j}$ is the Fourier coefficients of $V(x,y)$, then those coefficients for which

$$E[|V_{k,j}|^2] \doteq \lambda_{kj} < \mu \quad (16)$$

are not transmitted. The remaining coefficients are transmitted in such a way that the received coefficients $B_{k,j}$ have a distribution such that $V_{k,j} - B_{k,j}$ are Gaussian with zero mean, variance μ , and is independent of $B_{k,j}$. This can be simulated [16,19] by setting

$$B_{k,j} = \frac{\lambda_{k,j}^{-\mu}}{\lambda_{k,j}} [V_{k,j} + N_{k,j}] \quad (17)$$

in which $N_{k,j}$ are zero-mean complex random variables, independent of $V_{k,j}$ and whose real and imaginary parts are uncorrelated. The variance of $N_{k,j}$ is given by

$$E[V_{k,j}^2] = \frac{\mu \lambda_{kj}}{2} (\lambda_{kj}^{-\mu}) \quad (18)$$

The simulation then consists of the following steps:

1. Apply the nonlinear function (13) to the input image $u[x,y]$, i.e. $w[x,y]=f(u[x,y])$.
2. Compute the power spectral density $S_w[f_x, f_y]$ of $w[x,y]$ by using a smoothed periodogram technique, i.e. averaging the power spectral density computed from small overlapping sections of $w[x,y]$.

3. Weigh the power spectral density with the frequency sensitivity of the human visual system using (12), i.e.

$$S[f_r] = |A(f_r)|^2 S_w[f_r] \quad (19)$$

$$\text{where } f_r = \sqrt{f_x^2 + f_y^2}$$

4. Using (14) iteratively, compute the μ corresponding to the desired rate.
5. Compute $W[f_x, f_y] = \text{FFT} \{ w[x, y] \}$.
6. Weigh W with the HVS (12), i.e. $V[f_x, f_y] = A(f_r)W(f_x, f_y)$.
7. Compute the variance of the coefficients $V_{k, j}$, i.e. apply equation (16), setting all coefficients with variance lower than μ equal to zero and adding noise according to (17) & (18) to the rest.
8. Compute the inverse FFT of B .
9. Apply the inverse of the HVS function (12) to the result of step eight.
10. Apply the inverse of the nonlinear function (13) to the result of step nine, this gives the rate distortion simulated picture.

It should be noted that no provision has been made to accommodate overheads in the average bit rate using this simulation, i.e. coding the position of coded coefficients. To avoid this problem and to keep the class of images as wide as possible, the simulations were done using the first order Markov model. The power spectral density for this model is given by (Appendix D),

$$S(\omega) = 2 \sigma^2 \frac{\alpha}{\alpha^2 + \omega^2} \quad (20)$$

where σ^2 is the image variance,

and where the correlation $\rho = e^{-\alpha}$

The simulation results are shown in figures 9 and 10. The results show that very good quality images are possible for bit rates of 1.0 and 0.5 bits/pel. These pictures can now be compared to those produced by the actual codec. Based on the comparison one should be able to decide how much more can be done to improve the codec performance. For example one would be able to determine if the added complexity is worth the gain in image quality. The next paragraph looks at the range of frequencies for which (12) is valid in normal viewing situations.



Figure 9 Rate Distortion Simulation: GIRL rate=1 bit/pel.

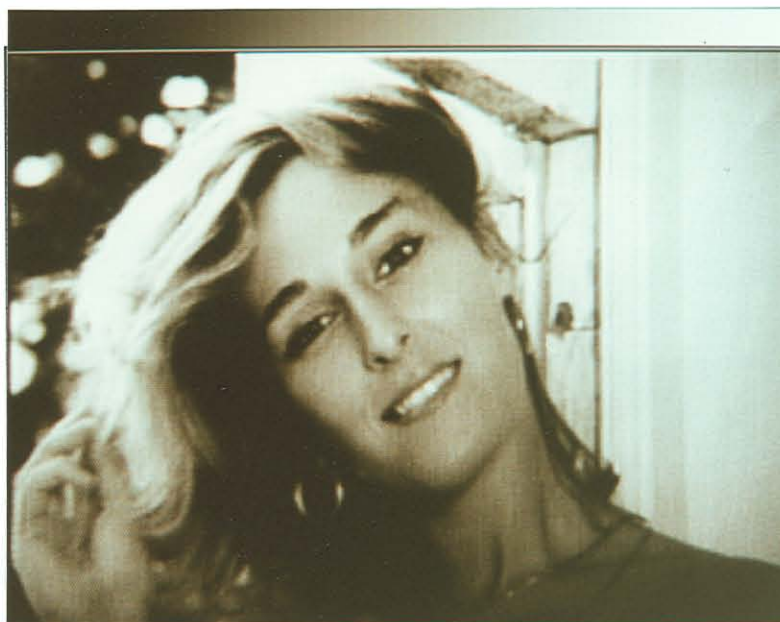


Figure 10 Rate Distortion Simulated Image: GIRL rate=0.5 bit/pel.

B. Parameters for the HVS

To apply (12) to the images used in this simulation it is necessary to compute the frequency range for which it is valid. This can be done by examining the system in figure 11. From this figure it is easy to determine that the angle spanned by the viewing device is given by the relationship

$$\theta = 2 \tan^{-1} \left(\frac{r}{d} \right) \quad (21)$$

where d is the distance from the screen, and $2r$ is the size of the screen. For a distance from the viewing device of 50cm and a *horizontal* screen size of 7.2" (9" diagonal) the horizontal viewing angle is 20.73 degrees. For an image resolution of 256 pels horizontally, the maximum number of cycles would be 128, which gives a frequency of $128/20.73 = 6.18$ cycles/degree.

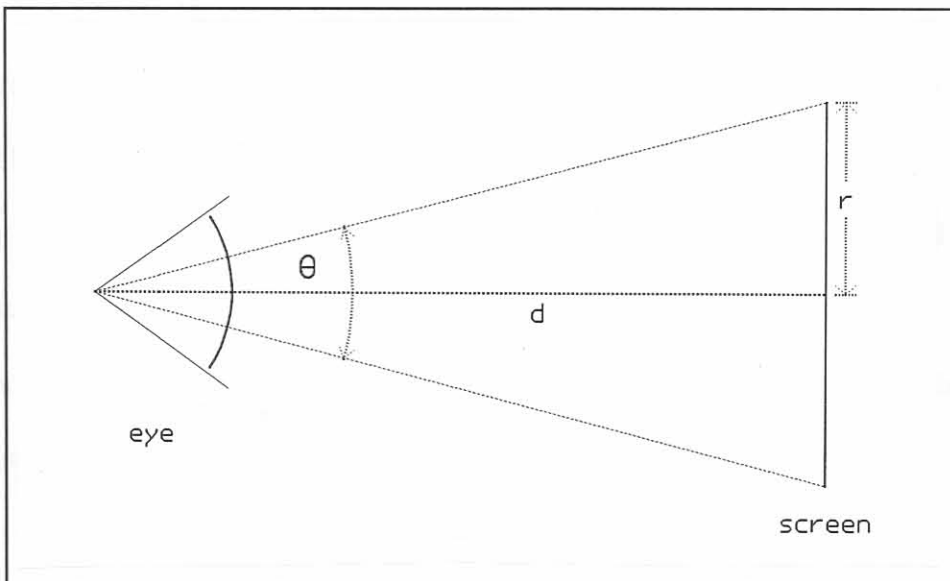


Figure 11 Determination of viewing angle from viewing distance and screen size.

A few different viewing angles have been computed for changing viewing distances and screen sizes, and are given in table I. From this table it is clear that the frequency response of an image, with resolution of 256x256, displayed on a device with the given physical parameters is lacking in comparison with the resolution ability of the human eye (64 cycles/ degree).

Viewing Distance	Horizontal Screen Size		
	4" Photo (5" diag.)	7.2" (9" diag.)	11.2" (14" diag.)
30cm	6.66	3.78	2.52
50cm	11.03	6.18	4.03
1m	22.01	12.25	7.91
2m	43.99	24.45	15.73

Table I Maximum normalised frequency available from different viewing devices for an image resolution of 256x256

For images with a resolution of 512x512 pels the frequency response will be exactly double that shown in the table.

The next section introduces the distortion measures that will be used to determine image quality in a more mathematical sense. It should be kept in mind that a subjective evaluation of the image coding results should be used as final evaluation measure.

VI. Image Quality

A very important part of the design of a high quality image coder is the evaluation of the performance of different algorithms inside the image codec itself. To do this it is necessary to specify a fidelity criterion, i.e. some measure by which we can determine the quality of the reconstructed image. The mean square error is frequently used in image processing applications. It is defined as follows for images

$$e_{ms}^2 = \frac{1}{N^2} \sum_{x=0}^{N-1} \sum_{y=0}^{N-1} (i(x, y) - o(x, y))^2 \quad (22)$$

where the input pels are designated by i , and the output pels are designated by o . There are two definitions of signal to noise ratios (SNR) which use the definition of the above error. The first is the normalised signal to noise ratio

$$NSNR = 10 \text{Log}_{10} \frac{\sigma^2}{e_{ms}^2} \quad (23)$$

where σ^2 is the variance of the original image. The second definition is that of the peak signal to noise ratio (PSNR), defined as

$$PSNR = 10 \text{Log}_{10} \frac{(255)^2}{e_{ms}^2} \quad (24)$$

Although the mse measure does not agree with subjective evaluations of coded images, the mse gives an indication of the physical accuracy of reconstruction, and as such it is a useful measure.

The mse measure indicates the accuracy of both subjectively-important and subjectively-redundant image reconstruction. It has been found from experience that the mse only starts to fail as a good measure when the signal to noise ratio is relatively low, i.e. when coding of subjectively redundant information cannot be tolerated. Since this is normally the case for bit rates below one bit per pixel, a HVS frequency response weighted mse measure will also be used. This measure was also used by Davisson [18], and is defined as follows: If I_f is the spectrum of the input image and O_f is the spectrum of the reconstructed image then the weighted mse (wmse) is given by

$$e_{wms}^2 = \frac{1}{N^2} \sum_{x=0}^{N-1} \sum_{y=0}^{N-1} |A_f(x,y)|^2 (I_f(x,y) - O_f(x,y))^2 \quad (25)$$

where A_f is the frequency response of the HVS given by (12). The HVS was implemented for a nine inch diagonal viewing device viewed from half a metre. From table I the maximum observable frequency would then be 6.18 cycles/degree. This would mean that the HVS frequency response would be that of a high pass filter.

The next section examines the basic transform image codec structure. It gives a description of all the different sections, gives results of images coded with this codec and then proceeds to an analysis of the origins of the errors in transform coding.

VII. Basic Transform Image Coding Structure

For the high quality coding and transmission of still pictures one needs to exploit all structure or redundancy present within the image. Transform image coding has been found to be a robust and efficient way to achieve this [1,4,5,6,7,61,62]. This section starts with a brief overview of the structure of a typical intraframe codec, which is followed by a more detailed discussion of the different aspects involved.

In transform coding the image is first decorrelated by using a suitable transform. This step is normally reversible and contributes little to the overall image degradation, except for round-off errors in integer implementations. With a suitable transform the decorrelation step can achieve efficient energy compacting, as will be discussed in more detail in a following section.

The next step involves the quantisation of the coefficients, and it is in this step that the actual image compression takes place. This step is *irreversible* and is responsible for the majority of degradation in the reconstructed image. The statistics of the coefficients are normally exploited in the quantisation step to minimize the distortion.

The last step is channel encoding, and involves adding some redundancy to minimise the effect of channel errors. This step is important because of the efficient representation of the image data, which implies that a small error might have a significant influence on

the output image. The decoder section at the receiver starts with channel decoding followed by a reconstruction of the coefficients and an inverse transform.

The next section will look at the selection and desirable properties of the transform to make the transform coding as efficient as possible.

A. Transforms

A truly optimum transform would result in the best picture quality using the least number of bits. This is a criterion which is difficult to specify quantitatively. A simpler criterion is to require that the transform coefficients be statistically independent, but this requires knowledge of the probability density function of images which we do not yet have. Using second order statistics we can find a transform that results in uncorrelated coefficients. From the view of the information theory (see equation 2), the transform attempts to decrease the entropy of the image by taking advantage of dependencies in the source.

A significant unitary transform is the Karhunen-Loeve transform (KLT) for random fields. It is the complete orthonormal set of basis images $b(\cdot)$ determined from the eigenvalue equation

$$\sum_m \sum_n r(k, l; m, n) b(i, j; m, n) = \lambda_{i,j} b(i, j; k, l) \quad (26)$$

where $r(\cdot)$ is the image covariance function.

The optimality of the KLT for image processing stems from the following two properties [51]:

1. It completely decorrelates the transform coefficients, i.e.

$$\text{Cov}[v_{k,l}(T), v_{m,n}(T)] = \sigma_{k,l}^2(T) \delta_{k,m} \delta_{l,n} \quad \text{for } T=\Phi \quad (27)$$

where T denotes an arbitrary $N^2 \times N^2$ unitary transform, Φ is the KLT and $\sigma^2(T)$ is the variances of the T -transform coefficients $v(T)$.

2. Compared to all other unitary transforms, the KLT packs the maximum expected energy in a given number of samples, say M , i.e.

$$\sum_{k,l \in S(\Phi)} \sigma_{k,l}^2(\Phi) \geq \sum_{k,l \in S(T)} \sigma_{k,l}^2(T) \quad \forall 1 \leq M \leq N^2 \quad (28)$$

where $S(T)$ is the set containing M index pairs (k,l) corresponding to the largest M variances in the T -transform domain. This property serves as a basis for transform data compression techniques.

Although the optimal transform is explicitly known, its use in practice results in problems such as:

- * Different basis functions for every class of images as a result of the non-stationarity of images,
- * Singularities may exist in the covariance matrixes which means that all the basis functions cannot be computed,
- * No fast transform exist for the KLT.

For the Markov image model as described in section IV the discrete cosine transform (DCT) has been found to be very similar to the KLT, and can be derived as the limiting case of the KLT as the correlation coefficient approaches one, i.e. $\rho \rightarrow 1$. The DCT has also been found to perform better than other transforms in many image coding applications [36,40,41,45]. The DCT does not perform well for negative correlations or for correlations below 0.5.

The forward DCT is defined by Ahmed, Natrajan and Rao [31]:

$$F(u) = \frac{2 c(u)}{N} \sum_{j=0}^{N-1} f(j) \cos \left[\frac{(2j+1) u \pi}{2N} \right];$$

$$u = 0,1,\dots,N-1$$

where

(29)

$$c(u) = \frac{1}{\sqrt{2}} \quad \text{for } u = 0$$

$$c(u) = 1 \quad \text{for } u = 1,2,\dots,N-1$$

$$c(u) = 0 \quad \text{elsewhere.}$$

and the inverse transform is

$$f(j) = \sum_{u=0}^{N-1} c(u) F(u) \cos \left[\frac{(2j+1) u \pi}{2N} \right]$$

$$j = 0,1,\dots,N-1$$
(30)

Since the DCT is separable, the two dimensional transform can be computed by transforming the rows of the image followed by transforming the columns using the one dimensional transform. Many fast algorithms has been proposed for the DCT [31,32,33,34,35,37].

Normally the whole image is not transformed in a single transform since the statistics in an image is highly spatially variant, and from the information theory it would be better to group areas of similar statistics together. Coding images using a full-image transform normally results in a loss of detail, it is also time consuming, requires a large amount of memory, and require higher precision mathematical

processors. This leads to the *spatially adaptive* quantisation that will be discussed in the next section.

1. Transform Block Size

The question that remains is that of block size. In real images the assumption of a Markov model are not valid especially if the block size is small [7]. Computer simulations on real pictures show [7] that the mean square error (mse) produced by transform coding improve with the size of the sub-picture, but it does not improve meaningfully beyond 16x16. Natrevali [7] also showed that the subjective image quality does not appear to improve much with block size beyond 4x4.

To determine the influence of block size on the image quality, the test images were coded to 1.0 bits/pel for block sizes 4x4, 8x8, and 16x16. The coefficients were coded using a Laplacian pdf optimised Lloyd-Max quantiser with optimal bit assignment for mean square error minimisation. A graph of peak signal to noise ratio versus block size is given on figure 12. Figures 13 to 15 contains the coded image results for the different block sizes. The signal to noise ratio (S/N) increased with an increase in block size. The difference in S/N decreases between successive block sizes as the block size increases. The 4x4 block size showed a noticeable vector quantiser type noise along edges in the image. There were little difference, in subjective image quality, between transform block sizes of 8x8 and 16x16. Most of the simulations done further on in the thesis will use a block size of 8x8.

The reasons being that it allows for better adaption to local statistics, and also because DCT-chips exist that can do 8×8 transform in real time. The existence of the DCT-chips probably means that most of the current real world systems will use this block size.

Since the quantisation forms such an important section of the coder structure, it will be discussed in a separate section. The following section will investigate the different approaches that is used in the quantisation of the transform coefficients.

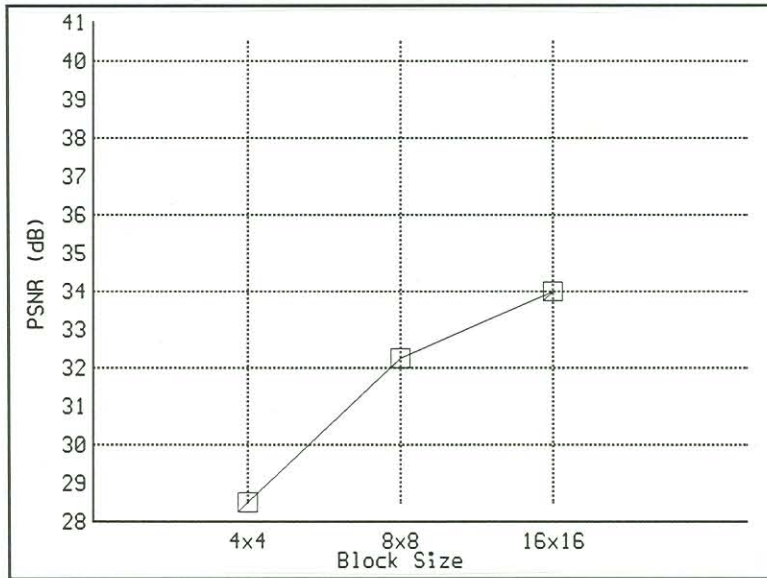


Figure 12 Signal to Noise Ratio versus Block Size for the image GIRL coded to 1.0 bits/pel.



Figure 13 GIRL coded using a DCT of size 4x4, 1.0 bits/pel.

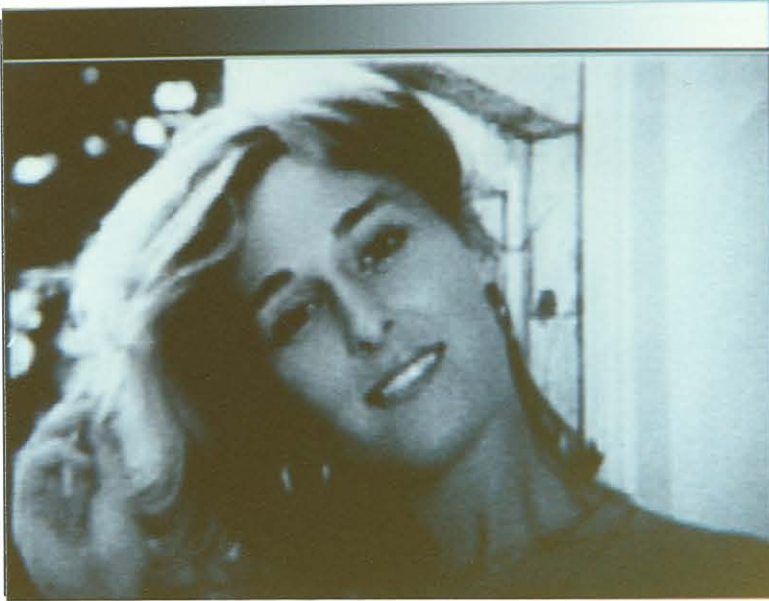


Figure 14 GIRL coded using an 8x8 DCT, 1.0 bits/pel.

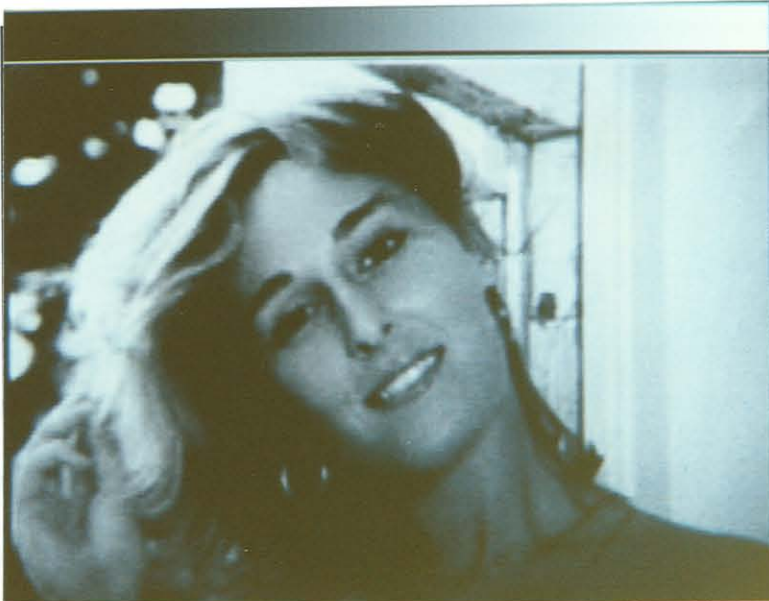


Figure 15 GIRL coded using a 16x16 DCT, 1.0 bits/pel.