

# An Evaluation of the Effectiveness of Adaptive Histogram Equalization for Contrast Enhancement

JOHN B. ZIMMERMAN, MEMBER, IEEE, STEPHEN M. PIZER, EDWARD V. STAAB, J. RANDOLPH PERRY, WILLIAM MCCARTNEY, AND BRADLEY C. BRENTON

**Abstract**—Adaptive histogram equalization (AHE), a method of contrast enhancement which is sensitive to local spatial information in an image, has been proposed as a solution to the problem of the inability of ordinary display devices to depict the full dynamic intensity range in some medical images. This method is automatic, reproducible, and simultaneously displays most of the information contained in the gray-scale contrast of the image. However, it has not been known whether the use of AHE causes the loss of diagnostic information relative to the commonly-used method of intensity windowing. In the current work, AHE and intensity windowing are compared using psychophysical observer studies.

In studies performed at North Carolina Memorial Hospital, experienced radiologists were shown clinical CT images of the chest. Into some of the images, appropriate artificial lesions were introduced; the physicians were then shown the images processed with both AHE and intensity windowing. They were asked to assess the probability that a given image contained the artificial lesion, and their accuracy was measured. The results of these experiments showed that for this particular diagnostic task, there was no significant difference in the ability of the two methods to depict luminance contrast; thus, further evaluation of AHE using controlled clinical trials is indicated.

## I. INTRODUCTION

MODERN imaging devices such as computed tomography (CT) and magnetic resonance (MR) scanners provide images with a large range of contrast information. Unfortunately, the range of displayable intensities available on common display devices such as CRT's or film is insufficient to display all this information simultaneously. The usual solution to this problem is to apply a contrast enhancement mapping such as intensity windowing (or more precisely, global linear min-max windowing, to be defined hereafter) to the image to display the information of interest. This approach has the difficulty that it throws away much of the available information; with intensity windowing it is often necessary to apply several windows

to the data in order to see the information available in different intensity ranges of the image. Further, intensity windowing requires manual intervention by the physician or technician to produce satisfactory results; often this manipulation must be repeated when the original film set is inadequate for complete interpretation.

An alternative method for contrast enhancement has been proposed by Pizer *et al.* [1] which has the property that it is adaptive to local information in the image. This method, adaptive histogram equalization (AHE), allows information in all intensity ranges of the image to be viewed simultaneously. It also has the advantages that it is automatic (i.e., no manual intervention is required) and is thus reproducible from study to study. Preliminary evaluation of AHE has shown that it has promise for routine use in a clinical setting; however, before this can be realized, it must be demonstrated that AHE does not cause a loss of diagnostic information relative to intensity windowing in most circumstances.

As a first step in this demonstration, this paper describes psychophysical observer experiments which assess the relative effectiveness of AHE and global linear min-max windowing in allowing the detection of contrast while using gray-scale display for a particular diagnostic task.

## II. DESCRIPTION OF THE TWO CONTRAST ENHANCEMENT MODALITIES

In the current studies, adaptive histogram equalization was compared to the most common clinical method of contrast enhancement, global linear min-max windowing. In this section these two contrast enhancement methods are briefly described.

Due to the coherence of objects within an image, the information about an object is generally contained in some subrange of the intensity values in the image. Global linear min-max windowing attempts to exploit this object coherence by enhancing the contrast within a subrange of the intensity values at the expense of the remaining intensity values. As shown in Fig. 1, a subrange or window of the data is chosen and these intensity values are linearly remapped to occupy the entire range of the display device. A single such mapping is applied to the entire image. Intensity values lying above or below the window are mapped to the maximum and minimum possible intensity values, respectively. The result is that the contrast of pixel

Manuscript received November 20, 1987; revised August 2, 1988.

J. B. Zimmerman is with the Department of Computer Science and Mallinckrodt Institute of Radiology, Washington University, St. Louis, MO 63130.

S. M. Pizer is with the Departments of Computer Science and Radiology, University of North Carolina, Chapel Hill, NC 27514.

E. V. Staab was with the Department of Radiology, University of North Carolina, Chapel Hill, NC 27514. He is now with the Department of Radiology, University of Florida, Gainesville, FL 32601.

J. R. Perry and W. McCartney are with the Department of Radiology, University of North Carolina, Chapel Hill, NC 27514.

B. C. Brenton was with the Department of Radiology, University of North Carolina, Chapel Hill, NC 27514. He is now with the Department of Radiology, Wayne Memorial Hospital, Goldsboro, NC 27530.

IEEE Log Number 8823777.

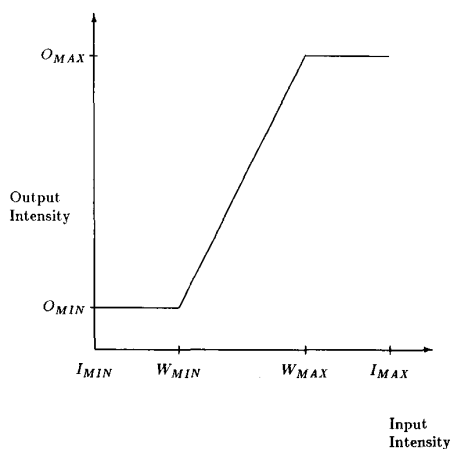


Fig. 1. Global linear min-max windowing.

values within the window is enhanced while information in other parts of the image is discarded. If other areas of the image are to be visualized, a new window must be chosen. Frequently, multiple windows are necessary to see all the information of interest in a clinical image. While standard windows can be selected to show particular areas of interest for a given anatomical region and imaging modality, manual intervention is usually necessary for the best results.

An alternative technique for contrast enhancement which has been widely used is global histogram equalization [2]. In this method, the intensity values in the image are altered such that the resulting image has a constant intensity histogram. This transformation may be accomplished by the use of the cumulative distribution function of the pixel intensities as the intensity remapping function. Such images utilize the available display levels well, but because the contrast enhancement is based on the statistics of the entire image, some levels will be used for the depiction of parts of the image which are diagnostically unimportant, such as the background.

Adaptive histogram equalization attempts to overcome the limitations of global linear min-max windowing and global histogram equalization by providing most of the desired information in a single image which can be produced without manual intervention. In this technique, developed independently by Ketcham *et al.*, Hummel, and Pizer *et al.* [3]–[5], the contrast enhancement mapping applied to a particular pixel is a function of the intensity values immediately surrounding the pixel. The basic method is illustrated in Fig. 2. For each pixel in the image, a region centered about the pixel, called its *contextual region*, is assigned. The intensity values in that region are used to calculate a histogram equalization mapping which is then applied to the pixel in question. The result is an image in which the mapping applied to each pixel is different and is adaptive to the local distribution of pixel intensities rather than the global information content of the image. In practice, this produces an image in which different objects whose intensity values lie in different

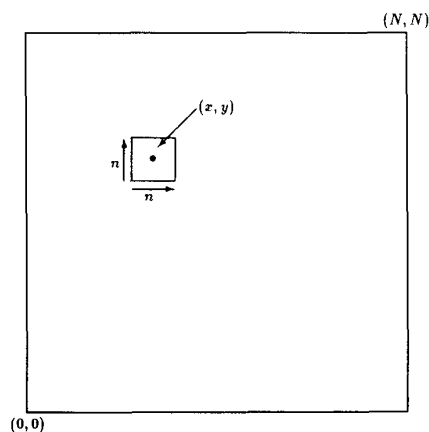
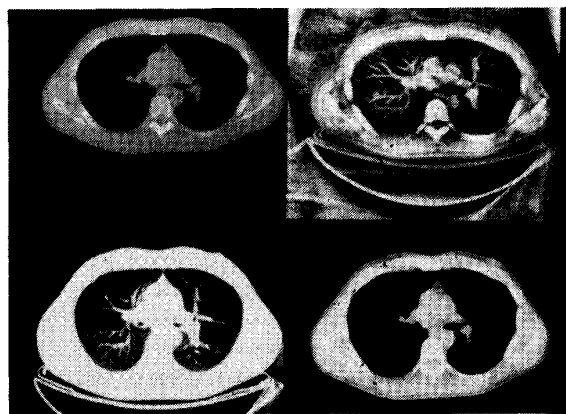
Fig. 2. Adaptive histogram equalization. The contextual region shown is a square of  $n \times n$  pixels about a pixel at location  $(x, y)$ .

Fig. 3. Adaptive histogram equalization. The original image is shown in the upper left. The upper right image has been processed with AHE. The lower left image is windowed to show the lungs; the lower right to show the mediastinum.

subranges of the intensity values are simultaneously visible.

In the method just described, the mapping for each pixel is determined by calculating a histogram equalization based on the pixel's contextual region. This method works well, but is too slow for implementation on a general-purpose computer, requiring about 2 h to process a  $512 \times 512$  image on a VAX11/780 computer. An alternative algorithm approximates the mapping required for each pixel by choosing a small number of contextual regions within the image and calculating the mapping for a given pixel as a bilinear interpolation of the mappings derived from nearby contextual regions. This interpolated version of AHE gives results which are very close to the uninterpolated version. A full description of this method may be found in Pizer [5]. Interpolated AHE allows  $512 \times 512$  medical images to be processed in a few tens of seconds on a general-purpose minicomputer or a few seconds on many special-purpose image processing systems.

An example of interpolated AHE is shown in Fig. 3.

Here, an original chest CT image is shown along with two windowed versions, one intended to depict the lung area and the other the mediastinum, and the same image processed with AHE. It is clear that AHE allows the simultaneous visualization of the major vessels in the mediastinum and the pulmonary vessels. It should be noted that unlike windowing and global histogram equalization, AHE does not preserve the rank ordering of the intensity values; that is, if two pixels  $a$  and  $b$  are such that  $I_a < I_b$  in the original image, this relationship may not be preserved in the AHE image, particularly if the pixels are spatially separated. Thus, quantitative comparisons of pixel intensities in different areas of the image are not possible using AHE images; however, such comparisons are not often needed in images used for diagnosis.

### III. EXPERIMENTAL METHODOLOGY

In order to evaluate the effectiveness of the two contrast enhancement methods described in the last section, formal observer studies were undertaken at North Carolina Memorial Hospital. The goal of these studies was to evaluate the methodologies using observer tasks which were as clinically realistic as possible while still allowing the results to be analyzed by well-understood statistical methods. In this section, the experimental methodology which was used is presented.

#### A. Overview of the Experiment

A receiver operating characteristic (ROC) rating experiment was performed to compare the ability of AHE and interactive global linear min-max windowing to depict luminance contrast. In this experiment, a set of test images was prepared from a number of normal CT images of the chest. In each normal image, four sites were chosen for the insertion of simulated lesions, two sites in the lungs and two in the mediastinum. For each field (lungs and mediastinum), three simulated lesions were prepared. The lesions varied in their linear size and gray-scale intensity. The linear sizes were chosen as appropriate for the given field; the intensities were chosen such that when a lesion was inserted into an image it was very subtle. From each normal image, twelve additional images were generated, each image having one of the lesions inserted at one of the selected sites. From each of these images, another image was prepared by processing the given image with AHE. The complete test set consisted of the normal images with no processing applied, the normal images with the lesions inserted and no processing applied, the normal images processed with AHE, and the normal images with the lesions inserted processed with AHE. In order that there be equal numbers of images with and without the artificial lesions, copies of the images without lesions inserted were included. The test set was then presented in random order to three trained radiologists.

For those images not processed with AHE, the observers were asked to perform interactive linear min-max windowing on the appropriate image field (lungs or mediastinum). After they had completed the windowing task,

they were presented with a replica of the lesion placed to the side of the image and the exact location of the prospective lesion was indicated with removable crosshairs. For those images processed with AHE, the observers were shown the lesion replica and the crosshairs immediately. In all cases, they were asked to rate their confidence that a lesion of the type depicted was present at the indicated site.

The task presented to the observers was thus one of detection of luminance contrast; the data constituted the results of a rating scale experiment from which a receiver operating characteristic curve could be derived, giving ordered pairs of true positive versus false positive responses for any desired level of observer confidence [6]. The data were analyzed using the CORROC program developed by Metz and his collaborators for correlated ROC data [7]. Observer performance was computed separately for each site in the image and for each simulated lesion. The areas  $A_z$  under the ROC curves and their standard deviations were calculated and used to compare the two modalities for a given site and lesion; the modalities were assumed to have no difference in their ability to depict luminance contrast if the difference in the areas of the two ROC curves was less than a preset criterion. The results were evaluated for two discrimination criteria, differences of 1.5 and 2.5 standard deviations in the areas under the ROC curves for the two contrast enhancement methods. In the following sections, more detail is given regarding the preparation of the test data and the experimental protocol.

#### B. Preparation of the Test Data

1) *Selection of Normal Images:* The set of normal chest images were chosen from CT scans of five separate patients who were classified as having no pathology in the areas of interest. The images were obtained on a Technicare 2060 CT scanner; the preliminary selection of images was done by an experienced radiologist. The images were taken from the scanner in digital form. The intensities in each image were calibrated in Hounsfield units, with the CT intensities approximately in the range  $-1000$  to  $+1200$ ; their spatial resolution was  $512 \times 512$  pixels. From about 100 slices in five patients, 32 slices were chosen. Adjacent slices were avoided to maximize anatomical differences between slices and reduce the possibility of memorization of normal variation by the observer. These 32 images were used to generate images for observer training and for the conduct of the actual experiment.

2) *Lesion Site Selection:* In each image, four sites, two sites in the lungs and two in the mediastinum, were chosen for the insertion of artificial lesions. The criteria for site selection were the presence of appropriate natural anatomy and the prevalence of real lesions at that site in clinical practice. Similar, but not identical, sites were chosen in each base image. This selection scheme ensured that the observer's task was approximately the same at each lesion site across the base set of images, a circumstance which enhances the probability that the rating scale

data can be pooled across observers and images for each site. Selection of the sites was done with the collaboration of an experienced radiologist. An example of an image with an artificial lesion inserted at one of the chosen sites is shown in Fig. 4; the intensity of the inserted lesion has been exaggerated for illustrative purposes.

The selection of multiple sites within the same image in two different image fields allows the comparison of a single AHE image with multiple different windowed images; this is important in that one of the apparent diagnostic advantages of AHE is that it allows the replacement of multiple windowed images with a single view which shows diagnostic information in many fields.

3) *Generation of Artificial Lesions*: The characteristics of the artificial lesions inserted into the chosen sites are given in Table I. A Gaussian intensity profile was chosen to approximate that which would be generated by a spherical tumor; the widths of the artificial lesions were chosen as appropriate to appear in the given field (lung or mediastinum).

The intensity profile for a given lesion is given by

$$I(x, y) = A(\Omega) \exp \left\{ -(x^2 + y^2)/2\sigma^2 \right\}. \quad (1)$$

The variance of the Gaussian  $\sigma^2$  determines the width of the lesion in pixels. There are two variances given in the table for each field (lung or mediastinum); the lesions corresponding to these two variances will be referred to as the narrow and wide lesions. These widths correspond visually to a noticeable difference in lesion size. The peak intensity of the lesion is given by  $A(\Omega)$ ; it depends on the local neighborhood  $\Omega$  of the lesion site, where  $\Omega$  is a square of side  $5\sigma$ , centered at  $(x, y)$ .

The function  $A(\Omega)$  must be chosen such that the lesion intensity is not too small, in which case insufficient information will be conveyed to the observer and the observer's decisions will proceed by guesswork, not too large, in which case the observer will make a correct choice on every trial and no discrimination of the contrast enhancement methods will be possible. A fixed lesion intensity is unacceptable, since it is well known that the ability of the eye to detect contrast is strongly dependent on both the mean intensity and the presence of structure in the near background of the lesion site. In addition, there is evidence [8] that the far background of the image also exerts influence on the detectability of the artificial lesion.

Preliminary tests performed on the images used in this experiment indicated that the appropriate lesion intensities must be chosen on an image-by-image basis. This choice was based on a simple measure of the structural complexity, the average absolute value of the Laplacian of the image intensity in a neighborhood  $\Omega$  of the lesion site. It was assumed that the desired peak intensity of the Gaussian  $A(\Omega)$  was given by

$$A(\Omega) = a \left| \nabla^2 I(x, y) \right| + b \quad (2)$$

a linear function of the average Laplacian value over the region  $\Omega$  which is a square  $5\sigma$  on a side centered at the position  $(x, y)$  of the lesion site.

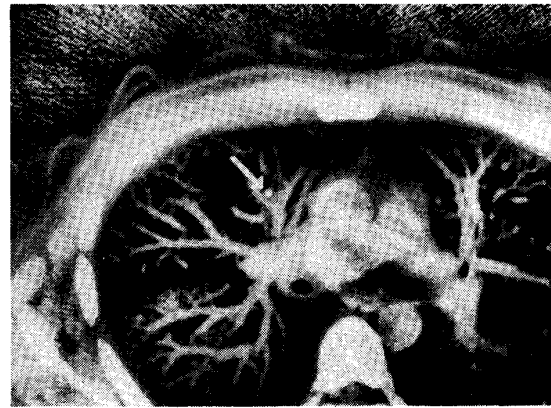


Fig. 4. Artificial lesion inserted into image of Fig. 3. The site shown is typical; the intensity of the lesion has been exaggerated for illustrative purposes.

TABLE I  
PARAMETERS FOR THE ARTIFICIAL LESIONS. WIDTHS ARE GIVEN IN PIXELS. THE HEIGHT FACTOR MULTIPLIES THE REQUIRED INTENSITY PREDICTED BY THE AVERAGE VALUE OF THE LAPLACIAN (SEE TEXT)

Lungs		
Lesion Number	Width ( $\sigma$ )	Height factor
0	1.41	0.85
1	2.0	1.15
2	2.0	0.85

Mediastinum		
Lesion Number	Width ( $\sigma$ )	Height factor
0	2.0	0.85
1	2.82	1.15
2	2.82	0.85

To determine the appropriate parameters for this function, a small number of the images (4) were examined by hand; for each of these an informal method of limits experiment was performed to determine the minimum detectable lesion intensity at each site. Since the detectability of a given lesion will be affected by the processing (windowing or AHE) done on the image at the time of observation, it was necessary to perform the method of limits determination after processing of the image.

A plot was then made of the lesion intensities determined by this method against the average Laplacian. A straight line fit of the data was made and the resulting line equation was used to determine the parameters  $a$  and  $b$ . A representative plot of the data and fitted line is shown in Fig. 5 for the lungs. The predicted intensities were checked during the execution of training runs to verify their appropriateness. Separate equations were used for the lung and mediastinal fields; no significant difference between the wide and narrow lesions was discernible, and the same coefficients were used for both.

The use of  $A(\Omega)$  is an attempt to predict the intensity necessary for the Gaussian to appear approximately one just noticeable difference (JND) above the given background in an image. Since there was the possibility that these predictions might be systematically biased such that the Gaussian would be either too bright or too dim, it was

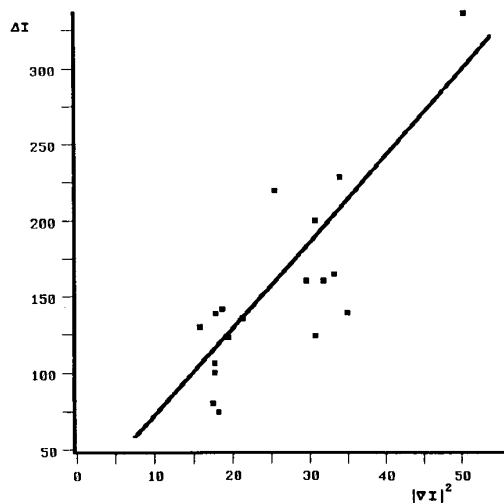


Fig. 5. A plot of the minimum detectable lesion intensity in Hounsfield units against the average absolute Laplacian value in the lung field, based upon observations conducted on four test images. The slope and intercept of the line shown was used to estimate the desired intensity level of the artificial lesion for the lung field in each base image. A similar determination was used for the artificial lesions in the mediastinum.

decided that two different possible lesion intensities, referred to as the dim and bright lesions, would be used at each site. These lesions were, respectively, 85 and 115 percent of the intensity predicted by the average Laplacian measure. These numbers were chosen experimentally to change the detectability parameter of the ROC curve to bracket the predicted value. Thus, there were four possible lesions which could be inserted at a given lesion site, corresponding to two choices for the width and two choices for the intensity. In this experiment, only three of these possible lesions were used; it can be seen in Table I that the bright, wide lesion was omitted.

4) *Preparation of the Trial Images:* The possible characteristics of a given trial image are shown in Table II. Each image is described by a 6-tuple of parameters that gives the base image, the lesion site, the type of lesion, lesion presence, etc. Of the 32 selected normal images, 24 were used as the base images for the generation of the trial set. Thus, there are 1152 possible trial images; all these were included in the trial set. In addition, approximately 10 percent (128) of the possible images were chosen at random and included twice, giving a total of 1280 images in the complete set. Hence, some images were seen twice by each observer, giving a check on the consistency of a given observer in rating the images. The 6-tuples for the 1280 images were generated and used to prepare the trial images in random order; this guaranteed no discernible pattern in the order of presentation. All calculations were made in the original range of the CT data and the final images were then scaled into the display range of the display device (0-255).

For each image, the inserted lesion was prepared by prescaling it to correspond to the local image structure as previously described. The prescaled lesion was then added (or not) to the image at the correct site. No attempt was

TABLE II  
PARAMETERS FOR THE TRIAL IMAGES

Lesion Parameter	Values	Comments
Base Image	24	Chest CT scans
Lesion Field	2	Lungs or Mediastinum
Lesion Type	3	As given in Table I
Lesion Site	2	Site within the field
Processing	2	AHE or windowing
Lesion Presence	2	Present or absent
Redundancy Factor	10%	Number of repeated images
Total	1280	All combinations

made to simulate the disarrangement of normal structure which would be present with real clinical lesions; the linear extent of the artificial lesions was sufficiently small that no significant effects would be expected. The same lesion was inset into the upper left corner of the image to serve as a reference for the observer. A typical site for an artificial lesion is shown in Fig. 4.

In addition to the trial set described above, the remaining 8 of the 32 normal images were used to generate 4 runs of 64 images each for use in training the observers. No image was used both for the trial set and the training set. One of the 4 runs was used exclusively to familiarize the observers with the equipment and experimental procedure; the other three formed an exhaustive set of three test images over all six parameters. These test runs were used to evaluate the choice of lesion parameters (shape and intensity) to ensure that they were reasonable. It was found that the scaling based on the average Laplacian produced acceptable results, in that the observer performances as measured by the areas under the ROC curve were about 70-80 percent correct and relatively consistent for a given lesion type. That is, the observers were performing neither with a high degree of certainty, nor by guesswork.

### C. Selection of Observers

Three observers were used in this study, all board-certified physicians and experienced in the reading of chest CT scans. The observers were guaranteed that no results of the experiment would be associated with a particular observer. The experiment was designed to occupy no more than 25 h of reading time, including the necessary training time. The observers were trained in the use of the experimental programs and allowed to do test runs to familiarize themselves with the appearance of AHE images.

### D. Experimental Apparatus

The trial images were displayed in gray-scale on a Tektronix 690SR RGB monitor using a Comtal 10/24 frame buffer. The Tektronix monitor has very stable performance; our experience has shown that its intensity display characteristics are stable on a time scale of weeks. The monitor was calibrated and converged prior to the beginning of the experiment and checked periodically throughout the course of the experiment. The display scale was linearized using the procedures described by Johnston *et al.* [9]. The luminance range of the monitor was  $6 \times 10^{-4}$



foot-lamberts for a driving intensity of 0 display units and 26.2 foot-lamberts for a driving intensity of 255; the size of the displayable screen area was  $26 \times 26 \text{ cm}^2$ . Spatial and intensity resolution of the RGB monitor was somewhat degraded relative to an equivalent monochrome monitor; however, careful control of lighting conditions allowed for excellent image quality. The images were displayed at a resolution of  $512 \times 512$  pixels on the Comtal 10/24 frame buffer.

### E. Experimental Layout

The observers were seated before the Tektronix monitor with the data tablet on a table in front of them. A light box was placed on top of the monitor facing away from the observers to provide an ambient light intensity of 3.3 lux, similar to that used for the linearization of the monitor. The observers were allowed to position themselves comfortably; no attempt was made to constrain their movements.

The physical environment of the observers was controlled insofar as this was possible in the experimental area. The room lights were extinguished and all extraneous sources of light shielded except for the light box used for ambient illumination. The observers were seated at an average distance from the screen of 1 m, with the displayed image subtending an angle of  $15^\circ$ . The observers were required to acclimate themselves to the environment for a period of 1–2 min before beginning the experiment.

### F. Observer Procedure

The experiment was divided into 20 sessions of length approximately 1 h each. At each session, the observer was asked to rate 64 images. Before embarking upon these sessions, each observer performed four 1 h runs for training purposes. Prior to the beginning of their training, the observers were required to read a document explaining the purpose and methods of the experiment and detailing the criteria on which their responses were to be based. After familiarization with the interaction devices and the use of the software, the observers did one training run of 64 images in which they received feedback on the correctness of their answers. This allowed them to calibrate their perceptions against the appearance of the displayed images. Following this, they performed three training runs without feedback to provide them with further experience and to allow the calibration of the experimental procedure against their performance. The results of these training runs allowed minor adjustments in the experimental procedure and the techniques for preparation of the trial runs to be made before the beginning of actual data collection.

The observers were shown each image in the trial set and asked to provide a rating of their confidence that the simulated lesion was present. Table III shows the definition of the five confidence categories. It should be noted that the rating scale is not symmetric; thus, there is no response in the category "lesion definitely present," nor is there any equivocal category. This requires the ob-

TABLE III  
RATING CRITERIA FOR THE ROC RESPONSE CATEGORIES

Rating	Confidence
0	Definitely not present
1	Probably not present
2	Possibly not present
3	Possibly present
4	Probably present

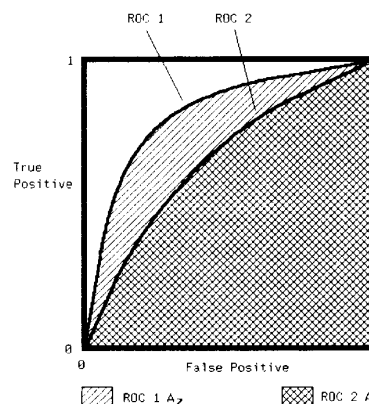


Fig. 6. Comparison of the areas beneath two ROC curves. These areas correspond to the parameter  $A_z$ , the percentage of correct answers. In this example, the method represented by ROC 1 is superior, since its area is greater.

server to make a decision as to which belief was stronger, the presence or absence of the lesion. The observers were urged to employ all five rating categories, using a 0 response on the occasions when they were most certain the lesion was not present and a 4 response they were most certain that the lesion was present.

### G. Analysis of Results

The result of the current experiments was a set of 1280 ratings for each observer; these can be considered as the outcome of 24 simultaneous experiments ( $2 \text{ fields} \times 2 \text{ sites per field} \times 3 \text{ lesions per site} \times 2 \text{ enhancement modalities}$ ). Each such experiment was analyzed separately using a program written by Metz *et al.* [7], CORROC. CORROC uses the maximum likelihood estimation technique of Dorfman and Alf [10] to estimate the ROC curve due to a given set of data in the case where the images are correlated. In the present experiments, the images are correlated, since each image is shown at least twice, once as a windowed image and once as an AHE'd image. For experiments which produced similar statistical results, such as trials which used the same field and site but different lesions, the results were pooled and reanalyzed. Goodness of fit estimations were also produced, allowing an analysis of statistical confidence.

The principal statistic used in comparing two similar experiments using different contrast enhancement modalities was the area under the ROC curve  $A_z$ ; this is equal to the percentage of correct responses one would obtain from a two-alternative forced-choice experiment and gives an overall measure of the performance of a given methodology. Fig. 6 shows two typical ROC curves and the

areas  $A_z$  beneath them. This approach precludes the comparison of the two enhancement methodologies at different levels of observer confidence; however, the integrated area is more stable against statistical fluctuations in the data. With the relatively small sample size used in the current experiment (approximately 50 data points per experiment), this is an important criterion.

#### IV. RESULTS

A summary of the results for each observer is shown in Table IV; the full results are given in Tables V–VII. In these tables, the headings of each column describe the parameters of a given experiment (lesion type, site, and processing modality), the integrated area under the resulting ROC curve, and the standard deviation of the area. The standard deviations shown have been corrected for the correlation of the data by the CORROC program. There is no data pooling in these results. The sixth column shows the number of standard deviations  $n_\sigma$  by which the areas of the two ROC curves differ. The seventh column is the two-tailed  $p$ -value; it represents the confidence with which the null hypothesis, in this case that the two ROC curves have the same area, can be rejected. A small value of  $p$  indicates that the two areas being compared are unlikely to have arisen from the same underlying distribution. The final column indicates which, if either, method was found preferable. Those studies marked with a “W” showed a preference for windowing with a difference in standard deviation of at least 1.5; those marked with an “A” showed a preference for AHE at the same level. The results were evaluated at two levels; if the entry has no asterisk, the areas under the ROC curve differ by more than 1.5 standard deviations but less than 2.5 standard deviations; if the asterisk is present, the areas differ by more than 2.5 standard deviations. To allow ready comparison, the data are arranged so that experiments with similar parameters (field, lesion type, lesion site, *etc.*) are adjacent with only the contrast enhancement modality different. The values of 1.5 and 2.5 standard deviations correspond to two-tailed  $p$  values of 0.1336 and 0.0124, respectively.

For three cases, the CORROC program failed to converge. In these cases, the data were analyzed as if uncorrelated using the ROCFIT program of Metz; this program performs a maximum likelihood estimation on uncorrelated data. These results are shown with an “I” in the Results column. Experience showed that the values for  $\sigma$  derived using this method did not differ substantially from those derived with CORROC; thus it can be asserted that none of these three results shows any significant difference between the two contrast enhancement modalities.

##### A. Unpooled Results

The results of the full experiment show that in most cases there is no significant difference in the diagnostic performance between the two modalities. For a difference of 1.5 standard deviations in the lung field, windowing was superior in two experiments, AHE in two, and no

TABLE IV

SUMMARY OF RESULTS FOR THE UNPOOLED RATING SCALE DATA. EACH COLUMN GIVES THE PREFERRED ENHANCEMENT METHOD, WINDOWING (W) OR AHE (A), AT THE SHOWN DISCRIMINATION LEVEL. FOR EACH FIELD, LUNGS OR MEDIASTINUM, THERE WERE SIX RESULTS, CORRESPONDING TO TWO SITES WITHIN THE FIELD  $\times$  3 POSSIBLE ARTIFICIAL LESIONS

Observer 1					
Field	W (1.5 $\sigma$ )	W (2.5 $\sigma$ )	A (1.5 $\sigma$ )	A (2.5 $\sigma$ )	None
Lungs	0	0	1	0	5
Mediastinum	0	2	0	0	4

Observer 2					
Field	W (1.5 $\sigma$ )	W (2.5 $\sigma$ )	A (1.5 $\sigma$ )	A (2.5 $\sigma$ )	None
Lungs	0	1	0	0	5
Mediastinum	1	1	0	0	4

Observer 3					
Field	W (1.5 $\sigma$ )	W (2.5 $\sigma$ )	A (1.5 $\sigma$ )	A (2.5 $\sigma$ )	None
Lungs	0	1	1	0	4
Mediastinum	2	0	0	0	4

TABLE V

RESULTS FOR OBSERVER 1. THE COLUMNS INDICATE THE LESION TYPE, SITE, AND PROCESSING METHOD; THE AREA UNDER THE ROC CURVE AND ITS STANDARD DEVIATION; THE NUMBER OF CORRECTED STANDARD DEVIATIONS BY WHICH THE AREAS DIFFER; THE TWO-TAILED  $p$ -VALUE; AND THE RESULT. IN THE RESULTS COLUMN, THE LETTER INDICATES THE PREFERRED METHOD (W = WINDOWING, A = AHE) AT THE 1.5 $\sigma$  LEVEL; THE ADDITION OF AN ASTERISK INDICATES THAT THE RESULTS WERE SIGNIFICANT AT THE 2.5 $\sigma$  LEVEL. AN I IN THE RESULTS COLUMN INDICATES THAT THE DATA WERE ANALYZED AS IF UNCORRELATED. THE LESION TYPES ARE AS GIVEN IN TABLE I. NO DATA POOLING WAS PERFORMED

Observer 1							
Lungs							
Type	Site	Proc	Area	$\sigma$	$n_\sigma$	pval	Result
0	0	W	0.7352	0.0694	0.6380	0.5235	
0	0	A	0.6808	0.0764	0.6380	0.5235	
1	0	W	0.7877	0.0628	-0.1635	0.8701	
1	0	A	0.7980	0.0650	-0.1635	0.8701	
2	0	W	0.6478	0.0747	-0.5953	0.5516	
2	0	A	0.7038	0.0741	-0.5953	0.5516	
0	1	W	0.7152	0.0680	-1.5024	0.1330	A
0	1	A	0.8303	0.0633	-1.5024	0.1330	
1	1	W	0.8866	0.0472			I
1	1	A	0.8836	0.0512			
2	1	W	0.7478	0.0682	-0.2562	0.7978	
2	1	A	0.7694	0.0672	-0.2562	0.7978	

Mediastinum							
Type	Site	Proc	Area	$\sigma$	$n_\sigma$	pval	Result
0	0	W	0.8567	0.0513	2.5621	0.0104	W*
0	0	A	0.6619	0.0761	2.5621	0.0104	
1	0	W	0.8205	0.0652			I
1	0	A	0.7499	0.0735			
2	0	W	0.6616	0.0748	0.5819	0.5606	
2	0	A	0.6031	0.0893	0.5819	0.5606	
0	1	W	0.7293	0.0705	0.6991	0.4845	
0	1	A	0.6720	0.0748	0.6991	0.4845	
1	1	W	0.7526	0.0693	2.5052	0.0122	W*
1	1	A	0.5795	0.0797	2.5052	0.0122	
2	1	W	0.6703	0.0769	1.0135	0.3108	
2	1	A	0.5993	0.0791	1.0135	0.3108	

significant difference was seen in 14 experiments. For 2.5 standard deviations in the lung field, two experiments showed windowing preferable; the rest showed no signif-

TABLE VI  
RESULTS FOR OBSERVER 2. NO DATA POOLING

Observer 2							
Lungs							
Type	Site	Proc	Area	$\sigma$	$n_\sigma$	pval	Result
0	0	W	0.6865	0.0736	-0.4104	0.6815	
0	0	A	0.7237	0.0720	-0.4104	0.6815	
1	0	W	0.8679	0.0465	0.1426	0.8866	
1	0	A	0.8609	0.0496	0.1426	0.8866	
2	0	W	0.8506	0.0524	2.6395	0.0083	W*
2	0	A	0.6378	0.0728	2.6395	0.0083	
0	1	W	0.7973	0.0573	-0.7875	0.4310	
0	1	A	0.8440	0.0516	-0.7875	0.4310	
1	1	W	0.9371	0.0318	1.3777	0.1683	
1	1	A	0.8779	0.0456	1.3777	0.1683	
2	1	W	0.7558	0.0708			I
2	1	A	0.8268	0.0624			

Mediastinum							
Type	Site	Proc	Area	$\sigma$	$n_\sigma$	pval	Result
0	0	W	0.9061	0.0503	1.9669	0.0492	W
0	0	A	0.7840	0.0609	1.9669	0.0492	
1	0	W	0.8929	0.0430	2.5087	0.0121	W*
1	0	A	0.7323	0.0693	2.5087	0.0121	
2	0	W	0.6924	0.0713	-0.1845	0.8536	
2	0	A	0.7067	0.0709	-0.1845	0.8536	
0	1	W	0.7587	0.0650	0.4053	0.6852	
0	1	A	0.7249	0.0704	0.4053	0.6852	
1	1	W	0.7698	0.0632	0.3883	0.6978	
1	1	A	0.7321	0.0704	0.3883	0.6978	
2	1	W	0.6363	0.0754	-0.6317	0.5276	
2	1	A	0.6881	0.0744	-0.6317	0.5276	

TABLE VII  
RESULTS FOR OBSERVER 3. NO DATA POOLING

Observer 3							
Lungs							
Type	Site	Proc	Area	$\sigma$	$n_\sigma$	pval	Result
0	0	W	0.7827	0.0640	2.5205	0.0117	W*
0	0	A	0.5963	0.0789	2.5205	0.0117	
1	0	W	0.6363	0.0960	-0.7300	0.4654	
1	0	A	0.7088	0.0683	-0.7300	0.4654	
2	0	W	0.5517	0.0837	-1.0492	0.2941	
2	0	A	0.6380	0.0737	-1.0492	0.2941	
0	1	W	0.7760	0.0897	-0.3544	0.7230	
0	1	A	0.8074	0.0577	-0.3544	0.7230	
1	1	W	0.6269	0.0955	-1.7126	0.0868	A
1	1	A	0.7915	0.0597	-1.7126	0.0868	
2	1	W	0.8184	0.0782	0.3866	0.6990	
2	1	A	0.7851	0.0619	0.3866	0.6990	
Mediastinum							
Type	Site	Proc	Area	$\sigma$	$n_\sigma$	pval	Result
0	0	W	0.7544	0.0635	0.9862	0.3241	
0	0	A	0.6807	0.0705	0.9862	0.3241	
1	0	W	0.8094	0.0618	1.2161	0.2239	
1	0	A	0.7258	0.0685	1.2161	0.2239	
2	0	W	0.6595	0.0733	-0.1898	0.8495	
2	0	A	0.6768	0.0812	-0.1898	0.8495	
0	1	W	0.8078	0.0611	2.2891	0.0221	W
0	1	A	0.6362	0.0766	2.2891	0.0221	
1	1	W	0.6350	0.0771	1.6268	0.1038	W
1	1	A	0.5249	0.0830	1.6268	0.1038	
2	1	W	0.6073	0.0767	-0.3827	0.7020	
2	1	A	0.6405	0.0771	-0.3827	0.7020	

icant difference. For 1.5 standard deviations in the mediastinum, windowing was preferred in six experiments, AHE in none, and no difference was seen in 11 cases. For 2.5 standard deviations in the mediastinum, three cases

TABLE VIII  
SUMMARY OF RESULTS AFTER POOLING OF DATA ACROSS LESIONS SITES. THE COLUMNS ARE AS GIVEN IN TABLE IV. THREE RESULTS ARE GIVEN, CORRESPONDING TO THE THREE DIFFERENT LESION TYPES AS GIVEN IN TABLE I

Observer 1					
Field	W (1.5 $\sigma$ )	W (2.5 $\sigma$ )	A (1.5 $\sigma$ )	A (2.5 $\sigma$ )	None
Lungs	0	0	0	0	3
Mediastinum	1	1	0	0	1
Observer 2					
Field	W (1.5 $\sigma$ )	W (2.5 $\sigma$ )	A (1.5 $\sigma$ )	A (2.5 $\sigma$ )	None
Lungs	0	0	0	0	3
Mediastinum	1	0	0	0	2
Observer 3					
Field	W (1.5 $\sigma$ )	W (2.5 $\sigma$ )	A (1.5 $\sigma$ )	A (2.5 $\sigma$ )	None
Lungs	1	0	1	0	1
Mediastinum	2	0	0	0	1

using windowing showed a significant difference and the remaining 13 showed no difference. Windowing did better in more cases in the mediastinum than in the lungs. This may be attributable to two factors: first, the overall data range in the mediastinum is usually less than in the lungs; hence, the amount of contrast enhancement which is achievable by windowing is greater. Second, the narrowness of the window width in the mediastinum is limited by CT noise, whereas in the lung it is limited by the presence of normal structure. Little difference is seen between observers compared to the intraobserver variation.

### B. Results Pooled Across Lesion Sites

Given the full results, it appeared to be valid to pool the data for the same lesion when shown in the same field (lungs or mediastinum) but at different sites in the field. These results are summarized in Table VIII and bear out the results from the full data in that in the lung field for 1.5 standard deviations, no methodology seemed preferable (one experiment favored windowing, one AHE, and seven showed no significant difference) while in the mediastinum for 1.5 standard deviations, five experiments showed windowing to be preferable while four showed no significant difference. For 2.5 standard deviations, only one experiment in the mediastinum showed windowing to be superior.

One difference which was easier to see in the pooled data is the difference in detectability of the different lesions. Lesions 0 and 2 appeared to be about equally detectable regardless of image field, while lesion 1 seemed to be more obvious than the other two. Lesion 1, as described in Table I, had the greatest intensity. The difference in width of lesions 0 and 2 seemed to make no significant difference.

### C. Time Series Analysis

One question of interest was whether the performance of the observers improved over time while conducting the



experiment. It might be expected that as observer familiarity with the AHE methodology increased, performance might improve. In order to investigate this possibility, the data pooled across the sites were arranged in the order in which the observer performed the runs; the rearranged data were divided into quarters and reanalyzed. The ROCFIT program was used since the division of the data precluded correlated analysis; this was not a drawback, since the purpose was to show trends in the data rather than absolute performance. There was no apparent increase in performance for the later runs. That is, the 4 h of training seems to have been adequate for the observers to learn fully to judge AHE'd images. No decline in performance was seen as a result of observer fatigue, either within a single trial run or over the course of the entire experiment.

## V. DISCUSSION AND CONCLUSIONS

In interpreting the results above, two circumstances must be kept in mind. First, the task which the observers were asked to perform was extremely difficult. The artificial lesions used had to be sufficiently subtle to allow some ambiguity in the interpretation of the image even when the exact location of the potential lesion was known exactly. Such lesions would be essentially undetectable in normal clinical practice. Given the preceding results, there is little question that both AHE and windowing are able to adequately depict supra-threshold contrast.

Second, in evaluating the levels of significant difference shown in Tables V–VIII, it must be realized that the criterion of a difference at the  $1.5\sigma$  level is very weak, with the corresponding  $p$ -value indicating a better than 1 in 8 chance of such a result arising from chance. A difference at the  $2.5\sigma$  level with  $p$ -value of 0.0124 is a much stronger criterion; however, the pooled data of Table VIII show that only one of the 18 results is significant at this level. Further, an examination of the actual percent correct (column 4 of Tables V–VIII) shows that in only one of the unpooled results and none of the pooled results does the probability of a correct answer differ by more than 0.2 between AHE and windowing.

Thus, the results shown indicate that there is very little difference in the ability of global linear min-max windowing and AHE to depict gray-scale contrast in an image. It seems probable that windowing holds some advantage in areas of the body (such as the mediastinum or liver) where the amount of contrast enhancement is limited by the image noise rather than the presence of normal structure; this advantage is unlikely to be realized in clinical practice, where the windows are usually not chosen as narrowly as in the current experiments. Given the fact that AHE is automatic and reproducible, further investigation through the use of controlled clinical trials seems desirable.

Recently, Pizer *et al.* [11], [12] have proposed a variant

of AHE which limits the amount of contrast enhancement in areas of the image with low variability. This contrast-limited adaptive histogram equalization (CLAHE) produces images in which the noise content of an image (such as is seen in the background of Fig. 3) is not excessively enhanced, but in which sufficient contrast enhancement is provided for the visualization of structures within the image. Images processed with CLAHE have a more natural appearance and facilitate the comparison of different areas of the image. However, the reduced contrast enhancement of CLAHE may hinder the ability of an observer to detect the presence of some significant gray-scale contrast. Psychophysical investigations similar to those described in this paper are currently underway to determine the relative ability of AHE and CLAHE to depict gray-scale contrast.

## ACKNOWLEDGMENT

The authors gratefully acknowledge the following for useful discussions and critiques of an earlier version of this paper: R. E. Johnston, F. P. Brooks, Jr., and J. J. Koenderink. They thank D. P. Schreiter for photographic assistance.

## REFERENCES

- [1] S. M. Pizer, J. B. Zimmerman, and E. V. Staab, "Adaptive grey level assignment in CT scan display," *J. Comput. Assist. Tomogr.*, vol. 8, pp. 300–308, 1984.
- [2] K. R. Castleman, *Digital Image Processing*. Englewood Cliffs, NJ: Prentice-Hall, 1979.
- [3] D. J. Ketcham, R. Lowe, and W. Weber, "Real-time enhancement techniques," *Seminar on Image Processing*, Hughes Aircraft, 1976, pp. 1–6.
- [4] R. Hummel, "Image enhancement by histogram transformation," *Comput. Graph. Image Processing*, vol. 6, pp. 184–195, 1977.
- [5] S. M. Pizer, "An automatic intensity mapping for the display of CT scans and other images," in *Proc. VIIIth Int. Meet. Inform. Processing in Med. Imaging*, 1983, pp. 276–309.
- [6] D. M. Green and J. A. Swets, *Signal Detection Theory and Psychophysics*. New York: Krieger, 1974.
- [7] C. E. Metz, P. Wang, and H. B. Kronman, "A new approach for testing the significance of differences between ROC curves measured from correlated data," in *VIIIth Conference on Information Processing in Medical Imaging*, F. Deconinck, Ed. The Hague: Martinus Nijhoff, 1984.
- [8] G. J. Van der Wildt and R. G. Waarts, "Contrast detection and its dependence on the presence of edges and lines in the stimulus field," *Vision Res.*, vol. 23, pp. 821–830, 1983.
- [9] R. E. Johnston, J. B. Zimmerman, D. C. Rogers, and S. M. Pizer, "Perceptual standardization," in *Proc. 3rd Int. Conf. Picture Archiving and Commun. Syst. (PACS III) for Med. Appl.*, SPIE, 1985, vol. 536.
- [10] D. D. Dorfman and E. Alf, Jr., "Maximum-likelihood estimation of parameters of signal-detection theory and determination of confidence intervals—Rating method data," *J. Math. Psychol.*, vol. 6, pp. 487–496, 1969.
- [11] S. M. Pizer, J. D. Austin, J. R. Perry, and J. B. Zimmerman, "Adaptive histogram equalization for automatic contrast enhancement of medical images," in *Proc. 4th Int. Conf. Picture Archiving and Commun. Syst. (PACS IV) for Med. Appl.*, SPIE, 1986, vol. 626.
- [12] S. M. Pizer, E. P. Amburn, J. D. Austin, R. Cromartie, A. Geselowitz, T. Greer, B. ter Haar Romeny, J. B. Zimmerman, and K. Zuiderveld, "Adaptive histogram equalization and its variations," *Comput. Vision, Graph. Image Processing*, vol. 39, pp. 355–368, 1987.