

ENGINEERING

P
NW



FILM RECORDING OF DIGITAL COLOR IMAGES

by

Robert Hall Wallis

May 1975

Image Processing Institute
University of Southern California
University Park
Los Angeles, California 90007

Sponsored by
Advanced Research Projects Agency
Contract No. F08606-72-C-0008
ARPA Order No. 1706

Ad

DISCONTINUED SUPPLEMENT A
Approved by 2015 Fall 1991
Dunbar Unlimited



IMAGE PROCESSING INSTITUTE

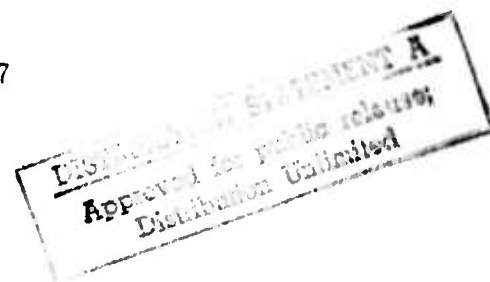
FILM RECORDING OF DIGITAL COLOR IMAGES

by

Robert Hall Wallis

May 1975

Image Processing Institute
University of Southern California
University Park
Los Angeles, California 90007



This research was supported by the Advanced Research Projects Agency of the Department of Defense and was monitored by the Air Force Eastern Test Range under Contract No. F03606-72-C-0008, ARPA Order No. 1706.

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

DOCUMENT CONTROL DATA - R & D

(Security classification of title, body of abstract and indexing annotation must be entered when the overall report is classified)

| | | | |
|--|--|---|----------------------|
| 1. ORIGINATING ACTIVITY (Corporate author) Image Processing Institute University of Southern California, University Park Los Angeles, California 90007 | | 2a. REPORT SECURITY CLASSIFICATION UNCLASSIFIED | |
| | | 2b. GROUP | |
| 3. REPORT TITLE FILM RECORDING OF DIGITAL COLOR IMAGES. | | | |
| 4. DESCRIPTIVE NOTES (Type of report and inclusive dates) Technical Report, May 1975 | | | |
| 5. AUTHOR(S) (First name, middle initial, last name) Robert H. Wallis | | | |
| 6. REPORT DATE 28 May 1975 | | 7. TOTAL NO. OF PAGES 196 | 8. NO. OF REFS 37 |
| 9. CONTRACT OR GRANT NO. F08606-72-C-0008 | | 9a. ORIGINATOR'S REPORT NUMBER(S) USCIP 14-570 | |
| 10. PROJECT TITLE ARPA Order 1706 | | 9b. OTHER REPORT NO(S) (Any other numbers that may be assigned this report) | |
| 10. DISTRIBUTION STATEMENT Approved for release: distribution unlimited | | | |
| 11. SUPPLEMENTARY NOTES | | 12. SPONSORING MILITARY ACTIVITY Advanced Research Projects Agency 1400 Wilson Boulevard Arlington, Virginia 22209 | |
| 13. ABSTRACT The goal of this study is the formulation of a technique for the quantitative control of color film. The specific application is the recording of digital images from a color monitor (or scanner) onto color film with maximum colorimetric fidelity. Inasmuch as human color perception is three-dimensional in nature, the recording process can be modeled as the passage of a three-dimensional signal through a trinary channel. Consequently, the distortions imposed upon the signal by the channel can be neutralized by a suitable pre-distortion of the input three-vector. Specifically, the pre-distortion should be the mathematical inverse of the channel mapping. A hierarchy of methods for the inversion of the channel mapping are developed, and tested. Although the numerical errors associated with the inversion process are shown to be insignificant, some subjective errors are apparent in the photographic reproduction if the predistortion is based on traditional fidelity criterion for color matching. A fidelity criterion based on a more relevant model of color vision, with particular attention to chromatic adaption, is developed. Using this model, photographic reproductions of color monitor images can be produced with excellent fidelity. ***** 14. Key words: Image processing, Image Restoration, Color Vision, Colorimetry. | | | |

391 141

Jones

| 14 | KEY WORDS | | LINK A | | LINK B | | LINK C | |
|----|-----------|--|--------|----|--------|----|--------|----|
| | | | ROLE | WT | ROLE | WT | ROLE | WT |
| | | | | | | | | |

ACKNOWLEDGEMENTS

The author wishes to express his appreciation to the members of his guidance committee: Dr. W. K. Pratt, Dr. H. C. Andrews, and Dr. J. Pierce, and to the staff and students of the Image Processing Institute, too numerous to mention, whose helpful suggestions and discussions aided in the preparation of this dissertation. In particular the guidance and advice of Dr. William K. Pratt, throughout the author's graduate career, is gratefully acknowledged.

This research was supported in part by the Advanced Research Projects Agency of the Department of Defense, and was monitored by the Air Force Test Range under Contract No. F-08.606- 72-C-0008

A

ABSTRACT

The goal of this study is the formulation of a technique for the quantitative control of color film. The specific application is the recording of digital images from a color monitor (or scanner) onto color film with maximum colorimetric fidelity. Inasmuch as human color perception is three-dimensional in nature, the recording process can be modeled as the passage of a three-dimensional signal through a trinary channel. Consequently, the distortions imposed upon the signal by the channel can be neutralized by a suitable pre-distortion of the input three-vector. Specifically, the pre-distortion should be the mathematical inverse of the channel mapping.

A hierarchy of methods for the inversion of the channel mapping are developed, and tested. Although the numerical errors associated with the inversion process are shown to be insignificant, some subjective errors are apparent in the photographic reproduction if the predistortion is based on traditional fidelity criterion for color matching. A fidelity criterion based on a more relevant model of color vision, with particular attention to

chromatic adaption, is developed. Using this model, photographic reproductions of color monitor images can be produced with excellent fidelity.

TABLE OF CONTENTS

| | |
|--|-----|
| 1. INTRODUCTION | 1 |
| 2. THREE COLOR THEORY OF HUMAN VISION | 7 |
| 2.1 Eye model and the trichromatic specification | 7 |
| 2.2 Tristimulus values and color coordinate systems | 16 |
| 2.3 Color adaption | 33 |
| 3. ADDITIVE COLOR SYSTEMS | 60 |
| 3.1 The additive principle and color television monitors | 60 |
| 3.2 Gamut and luminance restrictions | 69 |
| 4. SUBTRACTIVE SYSTEMS | 83 |
| 4.1 Subtractive color systems | 83 |
| 4.2 Subtractive systems; quantitative description | 88 |
| 5. COLOR PHOTOGRAPHY | 101 |
| 5.1 Chemistry and physics of color photography | 101 |
| 5.2 Mathematical model of transparency reversal film | 109 |
| 5.3 A comparison of empirical findings and theory | 120 |

| | |
|--|-----|
| 6. THE DISPLAY-FILM SYSTEM AS A MAPPING OF TRISTIMULUS VALUES | 130 |
| 6.1 Direct inversion of the film model | 131 |
| 6.2 Indirect inversion and computational considerations | 154 |
| 6.3 Experimental results | 160 |
| 7. GENERALIZATIONS, AND TOPICS FOR FUTURE RESEARCH | 165 |
| APPENDIX A | |
| Transformations between sets of tristimulus values | 171 |
| APPENDIX B | |
| Multidimensional mappings, and least squares curve fitting of surfaces | 174 |
| APPENDIX C | |
| Spatial Filtering of color images | 179 |
| APPENDIX D | |
| Transformations used in the film model | 183 |

| Figure | | page |
|---------|--|------|
| (1.0-1) | Conceptual block diagram of film recording process | 3 |
| (1.0-2) | Conceptual diagram of the tristimulus correction process | 4 |
| (2.1-1) | A first-order model for color vision | 10 |
| (2.1-2) | Color Mach bands (color plate) | xii |
| (2.1-3) | A color vision model incorporating spatial frequency effects | 14 |
| (2.2-1) | Experimental color matching apparatus | 17 |
| (2.2-2) | Color matching functions for spectral primaries | 20 |
| (2.2-3) | XYZ color matching functions | 23 |
| (2.2-4) | Chromaticity diagram | 26 |
| (2.2-5) | Constant luminance planes in Lab space (color plate) | xii |
| (2.2-6) | Television phosphor gamuts in ab plane of Lab space | 29 |
| (2.2-7) | Lab cylindrical coordinate system | 31 |
| (2.3-1) | Color shifts due to chromatic adaption | 36 |
| (2.3-2) | Von Kries model for chromatic adaption | 37 |
| (2.3-3) | Konig's receptor sensitivities | 43 |
| (2.3-4) | Adaptive bias model for chromatic adaption | 46 |

| | | |
|----------|--|-----|
| (2.3-5) | Chromatic adaption interpreted as a self-centering coordinate system | 48 |
| (2.3-6) | Effect of pre-nonlinearity gains on adaption shifts for one-third power nonlinearity | 50 |
| (2.3-7) | Effect of post-nonlinearity biases at low lightnesses, for one-third power nonlinearity | 51 |
| (2.3-8) | Effect of post-nonlinearity biases at medium lightnesses, for one-third power nonlinearity | 52 |
| (2.3-9) | Effect of post-nonlinearity biases at high lightnesses, for one-third power nonlinearity | 53 |
| (2.3-10) | Color shifts predicted by combined Von Kries and subtractive bias effects | 56 |
| (2.3-11) | Block diagram of adaption transform | 57 |
| (3.1-1) | Nonlinearities of color monitor display | 62 |
| (3.1-2) | Effect of brightness and gain controls | 65 |
| (3.1-3) | Effect of gamma correction on color film recording (color plate) | xii |
| (3.1-4) | Typical television phosphor characteristics | 66 |
| (3.1-5) | Comparison of modern phosphors and N.T.S.C. phosphors | 67 |
| (3.2-1) | Gamut of additive system in RGB space | |
| (3.2-2) | Gamut of additive system in XYZ space | 70 |
| (3.2-3) | Gamut of additive system in Lab space | 74 |
| (3.2-4) | Histogram of image hues in ab plane for test image of fig. 6.3-1 | 75 |

| | | |
|---------|--|-----|
| (3.2-5) | Geometrical interpretation of the gamut correction problem | 79 |
| (4.1-1) | The subtractive primaries | 86 |
| (4.2-1) | Mathematical model for absorption | 89 |
| (4.2-2) | Measurement of transmissivity | 94 |
| (4.2-3) | Gamut of colors in subtractive system | 99 |
| (5.1-1) | Hurter-Driffeld curves | 104 |
| (5.1-2) | Typical construction of color film | 105 |
| (5.1-3) | Schematic representation of modern color film processing | 107 |
| (5.2-1) | Mathematical model of color transparency film | 110 |
| (5.2-2) | Layer sensitivities for EKTACHROME-X film | 113 |
| (5.2-3) | Dye spectral densities for EKTACHROME-X film | 116 |
| (5.3-1) | Comparison of XYZ color matching functions and colorimeter characteristics | 121 |
| (5.3-2) | Actual vs. estimated XYZ color matching curves, using colorimeter | 124 |
| (5.3-3) | Comparison of actual vs. predicted tristimulus values, using the film model | 127 |
| (6.1-1) | The display-film cascade as a mapping | 132 |
| (6.1-2) | Spectral distributions of colors used in comparison of tristimulus integration methods | 151 |
| (6.3-1) | Photograph of original test image (color plate) | xii |

| | | |
|---------|--|-----|
| (6.3-2) | Photograph of test image pre-distorted by inverse of film model (color plate) | xii |
| (C-1) | Frequency response of chromatic channels in visual system model | 181 |
| (C-2) | Color image processed by high emphasis filtering (color plate) | xii |

Figure (2.1-2) Color Mach bands

Figure (2.2-5) Constant luminance planes in Lab space

Figure (3.1-3) Effect of gamma correction on color film recording

Figure (6.3-1) Photograph of original test image

Figure (6.3-2) Photograph of test image pre-distorted by inverse of film model

Figure (C-2) Color image processed by high emphasis filtering

Description of photographs in color plate

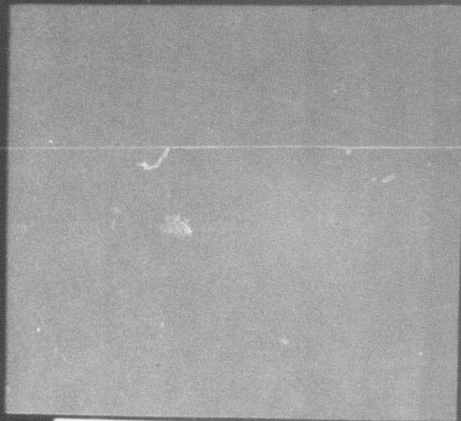


Figure (2.1-2)

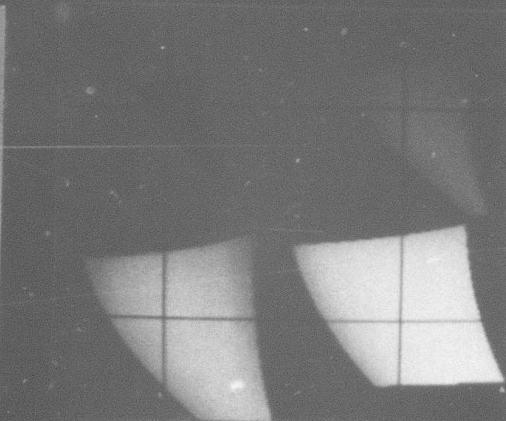


Figure (2.2-5)

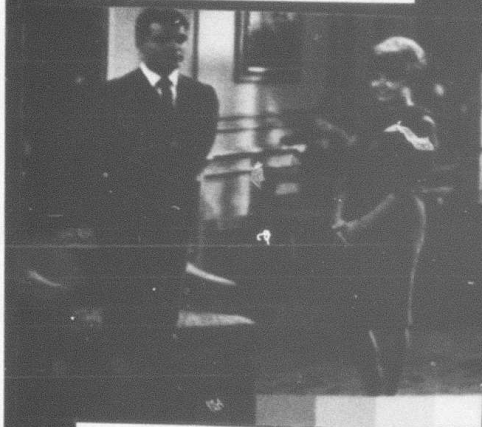


Figure (3.1-3)



Figure (6.3-1)

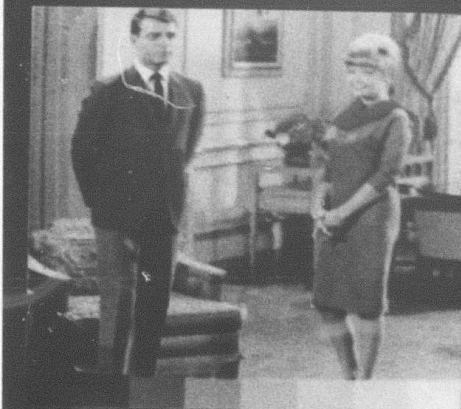


Figure (6.3-2)

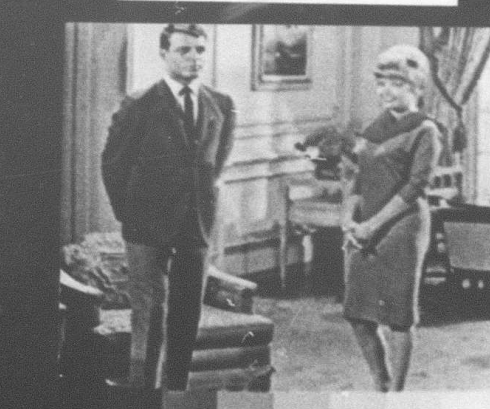


Figure (C-2)

NOTATION

| | |
|---------------------|---|
| λ | wavelength variable, expressed in nanometers |
| $s_k(\lambda)$ | spectral sensitivity of kth sensor (usually in reference to the receptors of the visual system) |
| S_k | output of kth sensor |
| P_k | weight on additive primary (usually TV phosphor) |
| $C(\lambda)$ | arbitrary spectral distribution |
| $R(\lambda)$ | reflectivity |
| $p_k(\lambda)$ | distribution of kth additive primary |
| $t_k(\lambda)$ | kth XYZ color matching function |
| T_k | kth XYZ tristimulus value |
| $W(\lambda)$ | distribution of reference white |
| $\tau(\lambda)$ | optical transmissivity |
| $D_k(\lambda)$ | optical density of kth dye |
| d_k | concentration of kth dye |
| $\langle * \rangle$ | footnote |
| [] | reference |
| \underline{T} | vector |
| A | matrix |

INTRODUCTION

The goal of the research to be described in this dissertation is the hopefully fruitful union of two diverse disciplines, color photography, and digital image processing. The first came into prominence in the 1930's and is today still unchallenged as a convenient means of recording color imagery. The second is of recent origin, and is rapidly gaining popularity with the increasing prevalence and economy of digital computing. The two fields are dissimilar in that the disadvantages of one often are the advantages of the other. Film recording is fast, very high resolution, relatively inexpensive, but can only be controlled with tedious and imprecise darkroom techniques. Digital image processing, on the other hand, offers perfect repeatability and control, but at the present time is slow and relatively expensive. One point at which the two disciplines meet is in the final recording of a digitally processed image. Unfortunately, this is often considered a mundane detail and little, if any, effort is made to understand or control the recording. Consequently, the result may be an image inadvertently processed to a greater degree by the photographic film than by the computer. In the case of

color images recorded from a television monitor, the distortions can be quite significant. However, if as much care is given to the modeling of the film recording system as is given to the processing of the image, quite faithful reproductions can be achieved. The philosophy to be used is that the film should be controlled not by chemical means, but by a pre-distortion of the image to be photographed.

The basis of color matching theory is the three dimensional nature of color vision. That is, the appearance of a color can be specified by only three numbers, known as tristimulus values. This is due to the fact that there exist only three types of color receptors in the human retina. Consider the schematic diagrams of figs. 1.0-1, and 1.0-2. The process of film recording is modeled as a mathematical correspondence (mapping) between the tristimulus values associated with the color being photographed (the input), and the tristimulus values associated with the color which results on the film (the output). Therefore, from a systems standpoint, a particular display-film combination can be represented as a trinary channel or "black box" with known input-output characteristics. By mathematically inverting the black

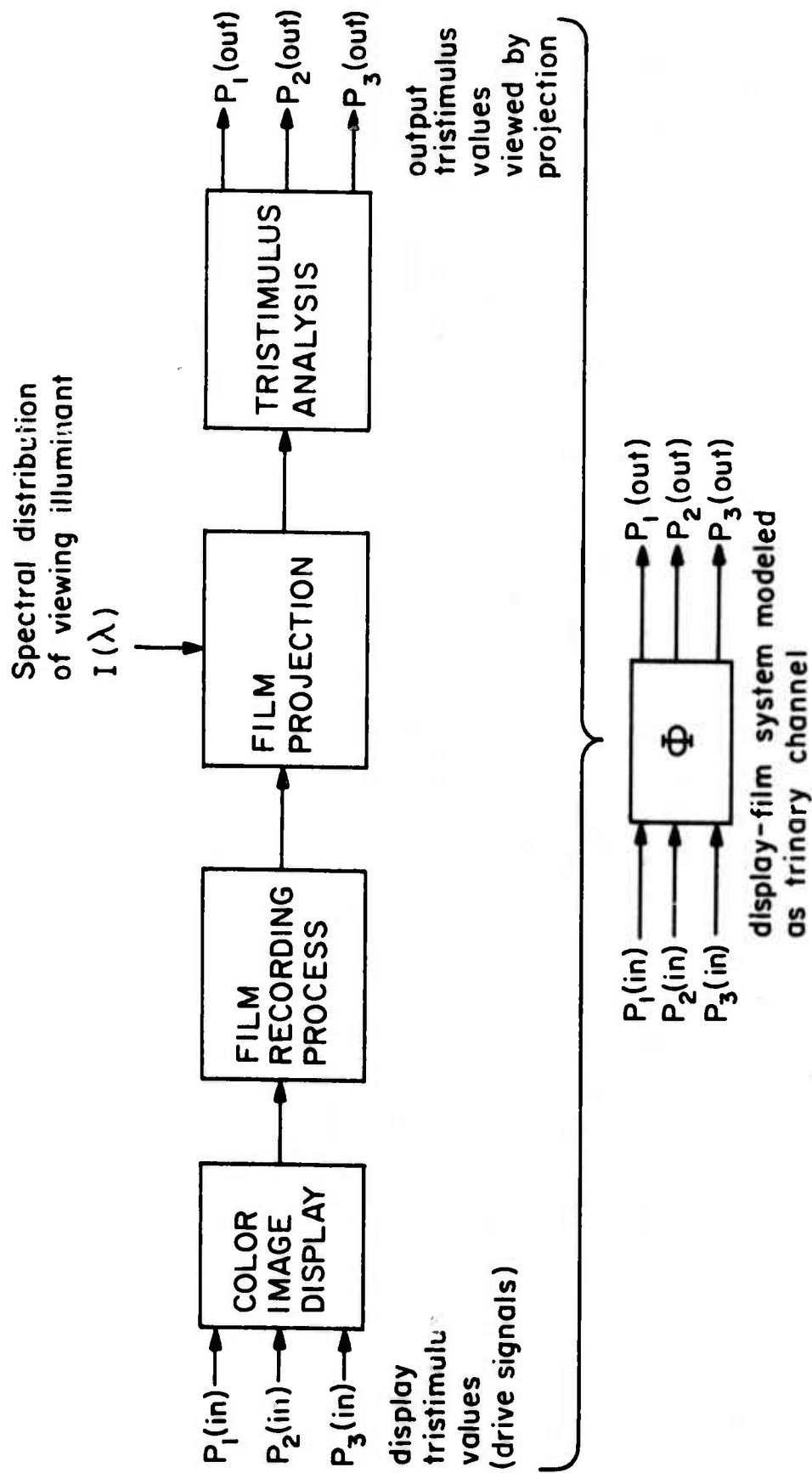


Figure (1.0-1) Conceptual block diagram of film recording process

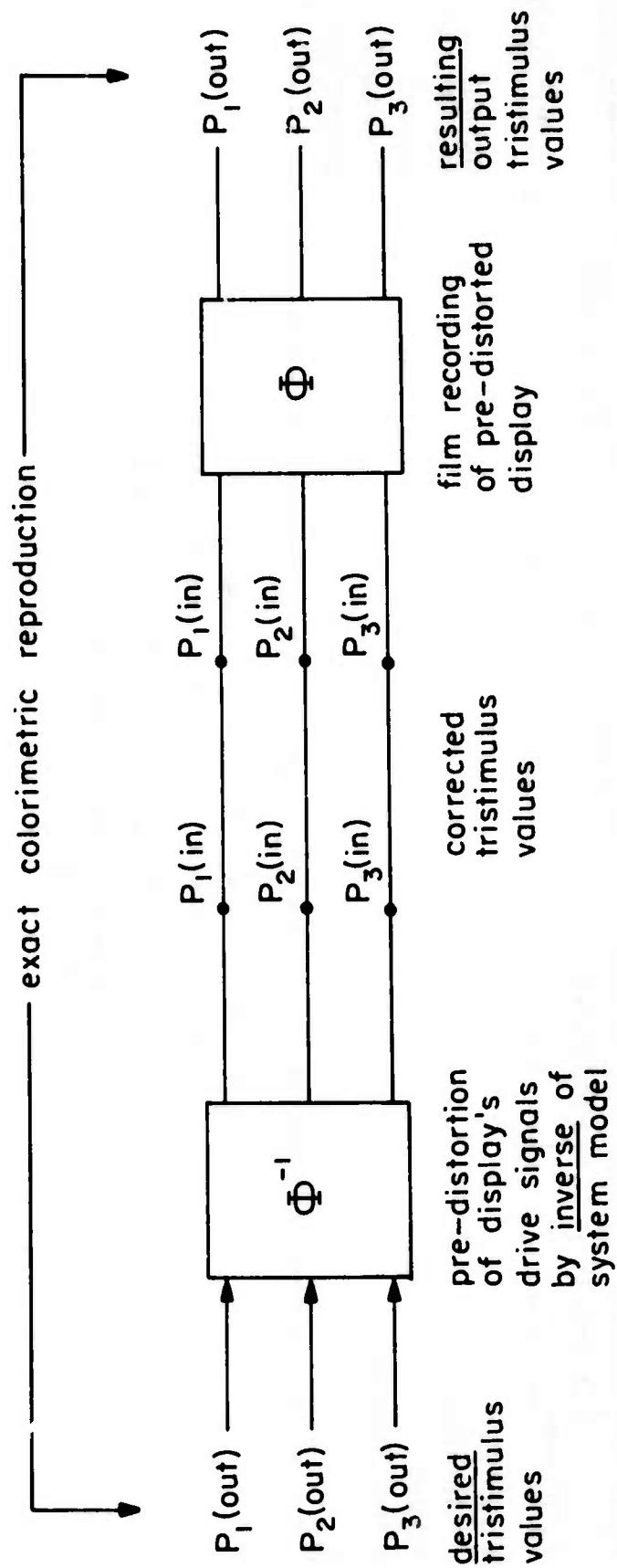


Figure (1.0-2) Conceptual diagram of the tristimulus correction process

box mapping, the input tristimulus vector necessary to produce a desired output tristimulus vector can be calculated. Consequently, this "inverse channel" can be used to "pre-distort" the tristimulus values of the image to be photographed, on a point-by-point basis. Thus, perfect colorimetric fidelity is possible in the recording process. The practical implementation of this approach is the subject of this dissertation.

Although digital techniques have been used for the analysis, coding, and transmission of color images [2], no studies on the colorimetrically exact film recording of digital images have been published. However, some previous work has been reported in which a color monitor image was pre-distorted to yield an "improved" photographic image [3,4]. However, the techniques were qualitative in nature and not concerned with an exact colorimetric match. While it can be argued that an exact duplication of tristimulus values is not necessary for a pleasing reproduction, the point is that qualitative techniques are rather empirical and lack precise control. However, once an accurate model of the film recording process has been assembled, the type of reproduction desired, whether exact or not, can be quantitatively

1.0

controlled.

REFERENCES

1. Cornsweet, T.N. Visual Perception, New York: Academic press. 1970
2. Pratt, W.K. Digital Image Coding and Transmission, USCEE Rept. 403, 1970
3. Lisk, K.G. and Evans, C.H. Color Television Film Recording from a Shadow Mask Picture Tube, Society of Mot. Pic. and Television Engrs., Oct., 1971
4. Robinson, S.P. and Horst, L.S. Color Film Recording from T.V., S.P.I.E. Journal Vol. 9 p. 196, 1971

THREE COLOR THEORY OF HUMAN VISION

The basis of color technology is the three dimensional nature of color vision. Consequently, color reproduction systems require only three degrees of freedom. Color television uses three guns and color film uses three layers. Even so called "four color" printing utilizes only three primary colors; the fourth color is black, and is used for reasons of economy (colored inks are expensive). The following chapter explores the assumptions of the three-color theory, its triumphs, and also its failures.

EYE MODEL AND THE TRICHROMATIC SPECIFICATION

The perception of color is certainly a complex phenomenon, sharing some similarities with the sense of hearing. Both faculties exhibit wide dynamic ranges, and tend to adapt to ambient levels of stimulus. (i.e. the senses tend to perceive relative changes, rather than absolute levels). Nevertheless, hearing and color

2.1

perception differ in one fundamental respect, their intrinsic dimensionality. Consider, for example, duplicating a certain sound, such as a particular note on a violin. In order for the replica and original to be indistinguishable to the ear, it is necessary that their power spectra match over the audible range. This indicates that an auditory match involves a nearly infinite number of degrees of freedom. In the case of color vision, duplication of the spectrum over the visible range is sufficient for a match, but is by no means necessary. If this were the case, color reproduction systems such as film and television [1] would be substantially more complex, if not impossible to achieve. The relative simplicity of these systems is due to the three dimensional nature of color perception. Specifically, a given color can usually be matched perfectly by an additive mixture of only three "primary" colors, provided that none of the primaries is merely a sum of the other two. The range of colors which can be generated by this method is known as the gamut of reproducible colors (for that particular set of primaries). If a color lies outside the gamut, then a match can only be attained by diluting the color with one

2.1

of the primaries until it can be matched with the remaining two. Furthermore, color matching obeys the laws of additivity and proportionality, i.e. if colors A and B match, and colors C and D match, then colors (A plus C) and (B plus D) match. Also, a match remains undisturbed if the intensities of the two colors are varied equally. The simplest model of color vision compatible with the above laws, shown in fig. 2.1-1 [2], consists of three receptors with spectral sensitivities $s_k(\lambda)$. The output of the kth receptor is given by

$$S_k = \int_{\lambda} s_k(\lambda) C(\lambda) d\lambda \quad (2.1-1)$$

where $C(\lambda)$ is the spectral distribution of the color being observed. The validity of the model is also supported by physiological evidence indicating that the retina contains three kinds of color receptors. In fact, various kinds of color blindness have been explained as deficiencies in one or more of the receptors [3, p. 155]. As with any simple model, that of fig. 2.1-1 has its inadequacies, the greatest of which is its failure to explain chromatic adaption (A term used to describe the eye's ability to perceive a wide range of illuminants as white). This

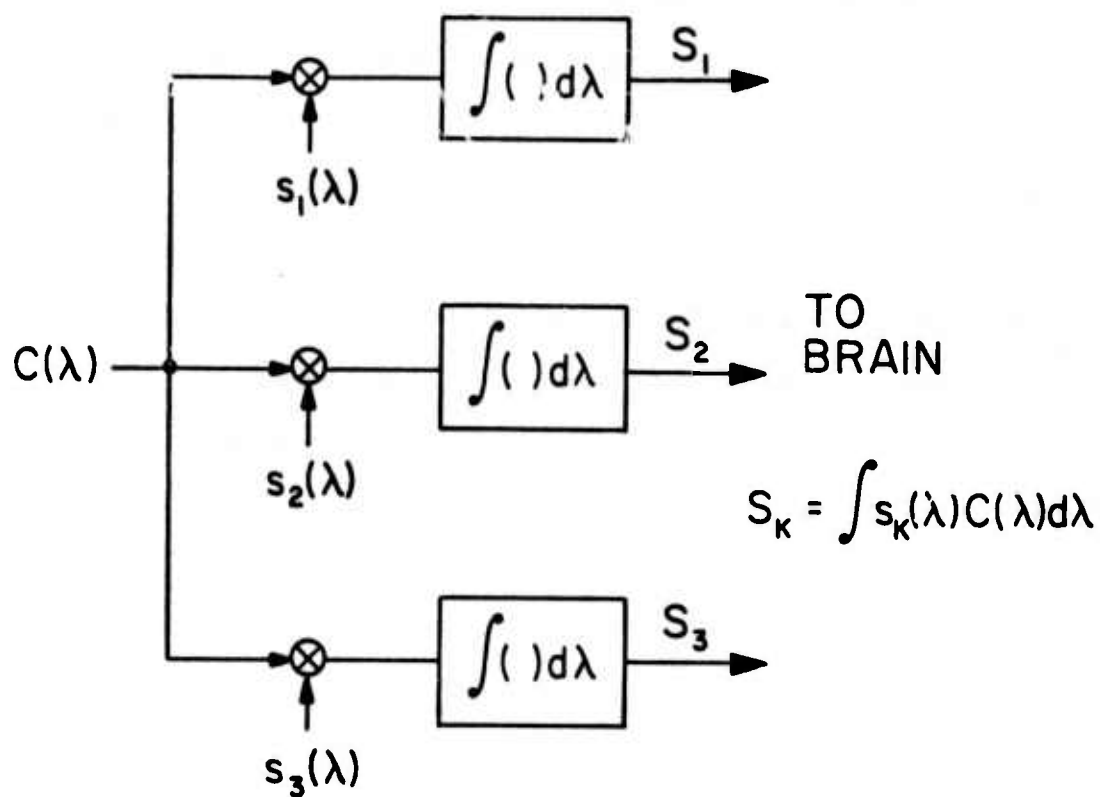


Figure (2.1-1) A first-order model for color vision

2.1

phenomenon is well known to photographers, since color film does not adapt to the ambient illumination, and hence must be balanced for one particular illuminant. For instance, if film balanced for daylight is used to photograph a scene illuminated with tungsten light, the resulting photograph will have a quite noticeable yellow bias. The question immediately arises, why doesn't the eye adapt to the yellow bias of the photograph? The answer is that the photograph usually subtends only a small portion of the field of view, and its surroundings serve as a reference which further exaggerates its poor color balance [4, p. 130]. This is the reason that a color print is less tolerant of errors in color balance than a transparency viewed in a darkened room. A projected transparency usually subtends a large field of view, and a darkened room deprives the eye of any reference by which to judge the color balance of the image. Chromatic adaption will be discussed more fully in Sec. 2.3.

A second but related phenomenon unexplained by the model of fig. 2.1-1 is known as simultaneous contrast. This is best illustrated by the color Mach bands of fig. 2.1-2 (color plate). The eight colors gradually shift from cyan to yellow, and are uniform across any

2.1

given bar. Perceptually, however, it appears that the cyan color becomes deeper just before the transition to yellow. Similarly, the yellow at the interface between the bars seems more intense. This is quite analogous to black and white Mach bands which appear as overshoots at the transitions between different grey levels [3, p. 278]. The phenomenon can best be appreciated by masking off the fields adjacent to a given bar with white cards. This will suppress the Mach band phenomenon, showing that each bar is indeed uniform. For unknown reasons, this effect appears to be strongest for transitions between cyan and yellow, and is substantially weaker or totally absent for other pairs of colors. In color images of very fine structure, another poorly understood phenomenon appears. This is the Van Bezold spreading effect [4, p. 181]. In this phenomenon, the color shifts are exactly opposite to those predicted by contrast effects. The color shifts appear to be an averaging or low-pass filtering effect in which neighboring colors are blended together by the eye. Furthermore, the model fails to explain why the eye responds to stimuli in a quasi-logarithmic manner [3, p. 373].

More complex visual models which attempt to explain the

2.1

above mentioned contrast and adaption effects have been proposed [5]. The block diagram of fig. 2.1-3 shows how the simple first order model can be extended to explain a number of color phenomenon. Fig. C-1 indicates the possible nature of the frequency response in the lightness and chromatic channels. Note that the cross-couplings and spatial filters in no way invalidate the color matching laws deduced from the simpler model. That is, if two spectral distributions generate the same receptor signals, S_k , then the two colors will still match, because no amount of subsequent processing will allow the visual system to distinguish identical signals. Later chapters will discuss how the model of fig. 2.1-3 can be used to analyze chromatic adaption phenomenon (Sec. 2.3), formulate "uniform" color coordinate systems (Sec. 2.2), and provide guidelines for color image enhancement (APPENDIX C).

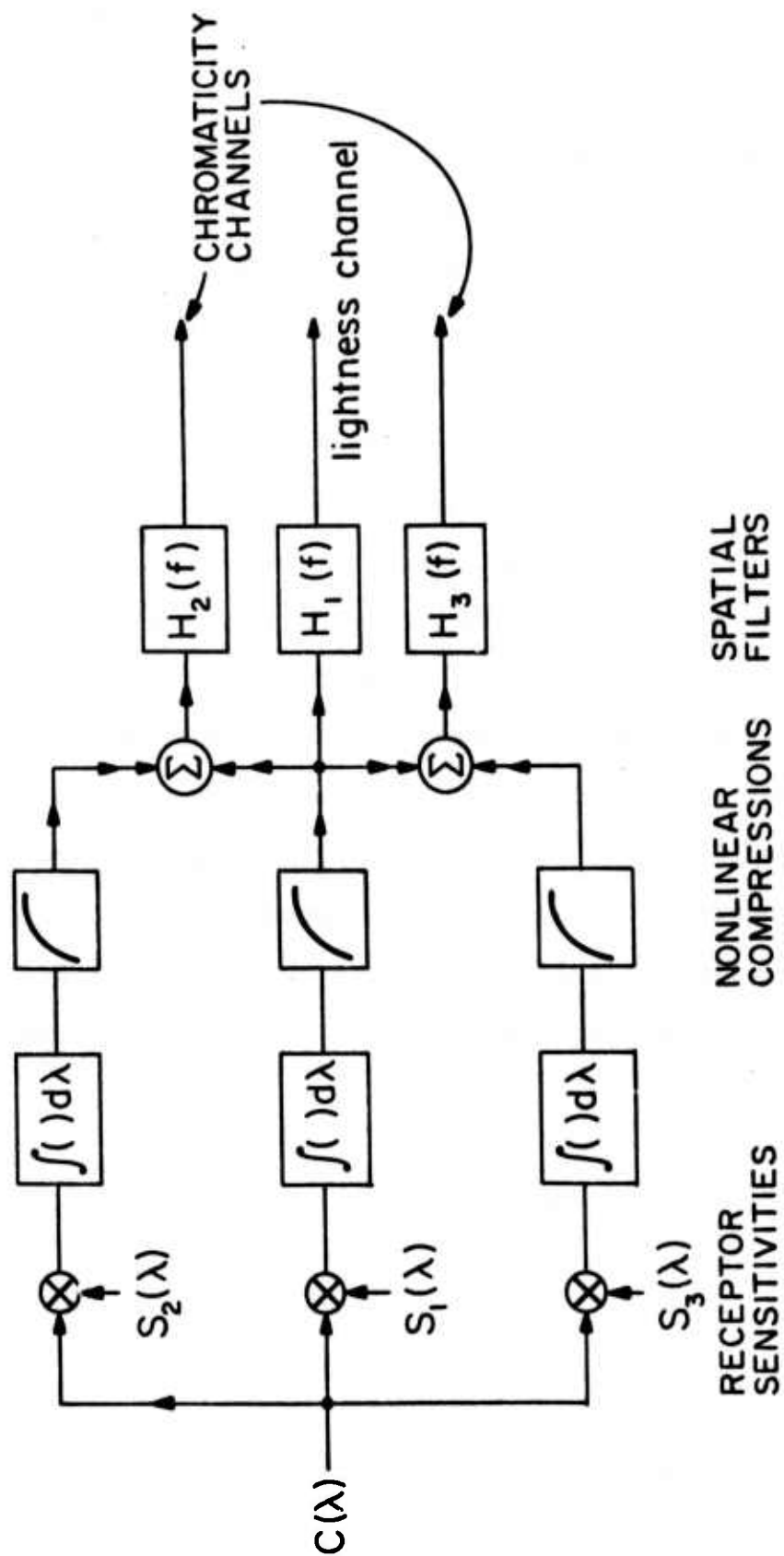


Figure (2.1-3) A color vision model incorporating spatial frequency effects

REFERENCES

1. Hunt, R.W.G. The Reproduction of Colour, London: Wiley and Sons. 1967
2. Wyszecki, G. and Stiles, W.S. Color Science, New York: Wiley and Sons. 1967
3. Cornsweet, T.N. Visual Perception, New York: Academic press. 1970
4. Evans, R.M. An Introduction to Color, New York: Wiley and Sons. 1948
5. Frei, W. Modeling Color Vision for Psychovisual Image Processing, USCEE Rept. 459 p. 112. 1973

TRISTIMULUS VALUES AND COLOR COORDINATE SYSTEMS

Having established the three dimensional nature of color vision, it would seem reasonable to specify a color's appearance by a triplet of numbers. Although the outputs of the three color receptors would be a logical choice for the specification, the spectral responses of the actual receptors in the eye are difficult to measure. However, since any linear combination of the receptor outputs can serve as an equally valid specification, it will be demonstrated that an indirect method such as color matching can be used. In color matching, a subject attempts to duplicate the appearance of a test color by an additive mixture of three primaries [1, p. 395]. A simplified schematic of the experimental setup is shown in fig. 2.2-1. Now consider the matching of a spectral line of unit power at wavelength $\lambda = \lambda_0$. Assuming the validity of the eye model (fig. 2.1-1 and eqn. 2.1-1), a match will be attained if and only if both colors evoke the same three "signals" from the eye's receptors, i.e.

$$\int_{\lambda} s_i(\lambda) \sum_{j=1}^3 t_j p_j(\lambda) d\lambda = \int_{\lambda} s_i(\lambda) \delta(\lambda - \lambda_0) d\lambda \quad (2.2-1)$$

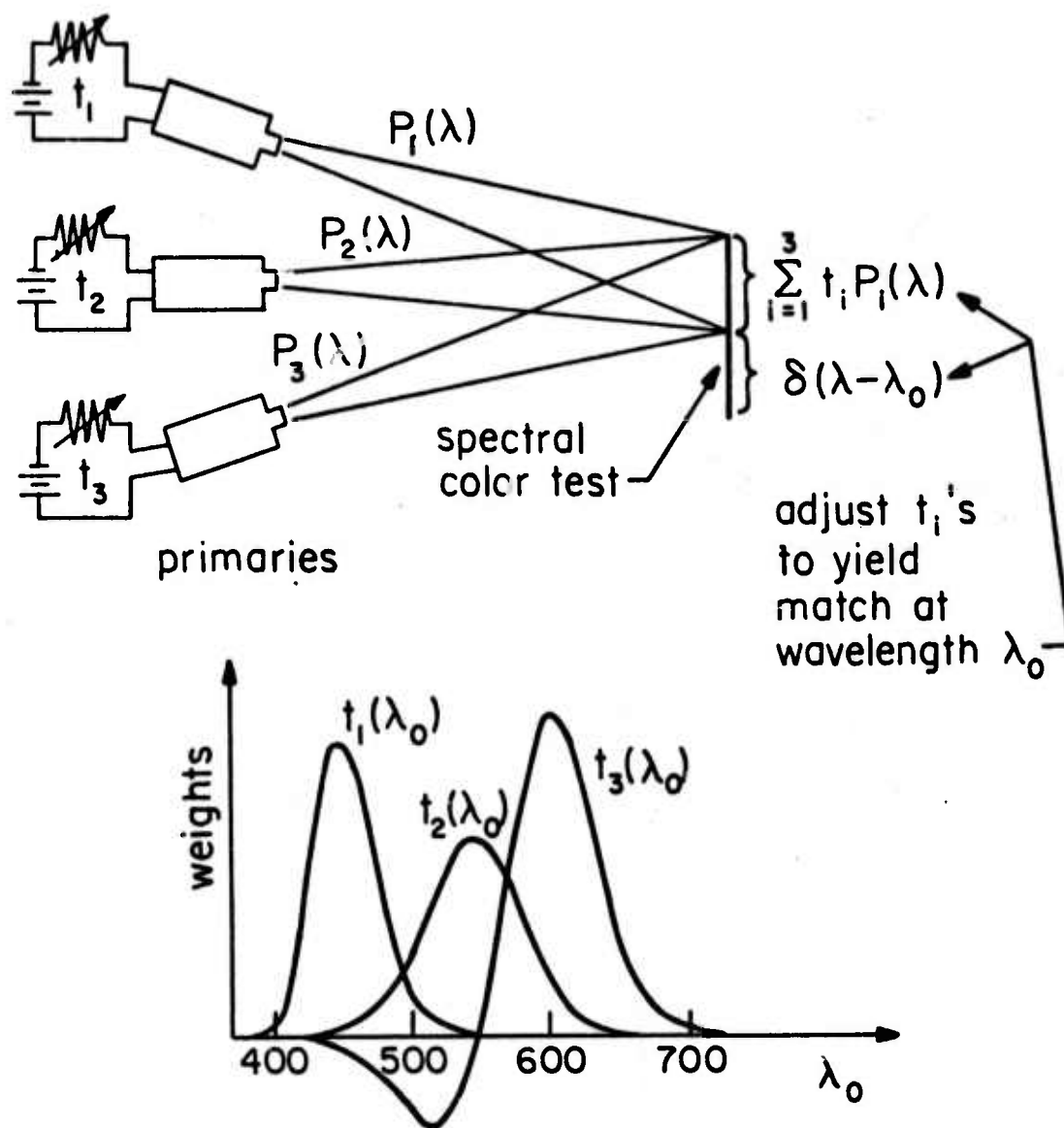


Figure (2.2-1) Experimental color matching apparatus

2.2

or

$$\sum_{j=1}^3 t_j \int s_i(\lambda) p_j(\lambda) d\lambda = s_i(\lambda_o) \quad (2.2-2)$$

for $i=1,2,3$, where the $p_j(\lambda)$ are the spectral distributions of the primaries, the t_j are weights to be determined, and the $s_i(\lambda)$ are the receptor sensitivities of the eye. Expressing 2.2-2 in matrix notation gives

$$\begin{bmatrix} a_{ij} \end{bmatrix} \begin{bmatrix} s_1(\lambda_o) \\ s_2(\lambda_o) \\ s_3(\lambda_o) \end{bmatrix} = \begin{bmatrix} t_1(\lambda_o) \\ t_2(\lambda_o) \\ t_3(\lambda_o) \end{bmatrix} \quad (2.2-3)$$

where

$$a_{ij} = \int s_i(\lambda) p_j(\lambda) d\lambda \quad (2.2-4)$$

Thus, the weights on the primaries required to match a spectral line are always linear combinations of the receptor sensitivities at that wavelength. If the wavelength of the spectral line to be matched is swept across the visible range, (380 to 775 nanometers), the weights, $t_j(\lambda_o)$, being functions of wavelength, form a

2.2

family of three curves, known as color matching functions [1, p. 447]. A group of such curves generated by using spectral primaries located at 700.0 nm. (red), 546.1 nm. (green), and 435.8 nm. (blue) is shown in fig. 2.2-2 [1, p. 223]. Now consider matching a continuous spectral distribution by a weighted sum of three primaries. This can be done by treating the continuous distribution as a "picket fence" of delta functions

$$C(\lambda) = \lim_{\Delta\lambda \rightarrow 0} \sum_k C(\lambda_k) \delta(\lambda - \lambda_k) \Delta\lambda \quad (2.2-5)$$

with convergence in the sense that

$$\lim_{\Delta\lambda \rightarrow 0} \int W(\lambda) \left[C(\lambda) - \left(\sum_k C(\lambda_k) \delta(\lambda - \lambda_k) \Delta\lambda \right) \right] d\lambda = 0 \quad (2.2-6)$$

or

$$\int W(\lambda) C(\lambda) d\lambda = \lim_{\Delta\lambda \rightarrow 0} \sum_k W(\lambda_k) C(\lambda_k) \Delta\lambda \quad (2.2-7)$$

where $W(\lambda)$ is any smooth function. Since the weights on the primaries needed to match a spectral line are merely the values of the color matching functions at that wavelength, the weights, T_i , required to match the

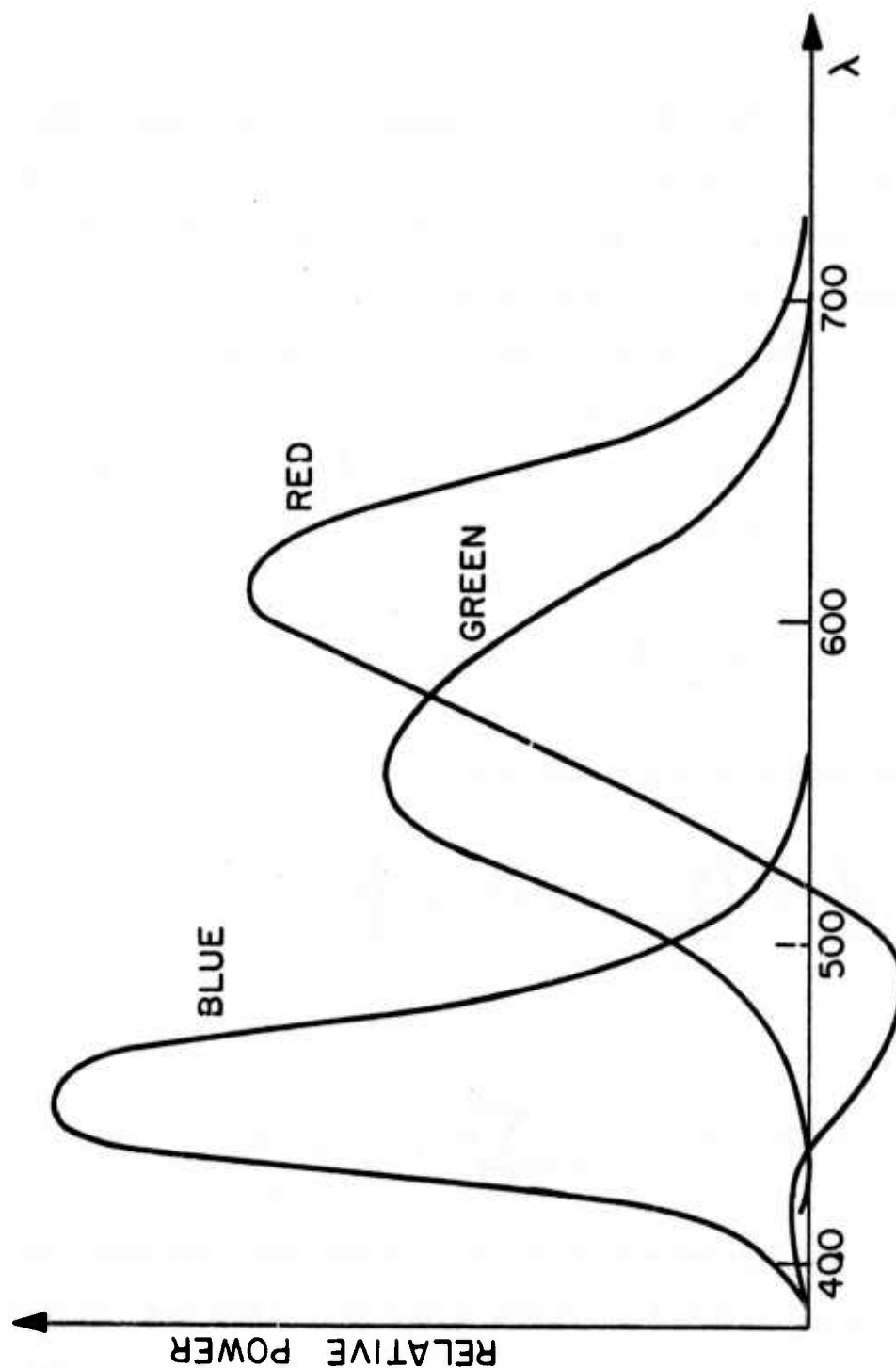


Figure (2.2-2) Color matching functions for spectral primaries

2.2

distribution, $C(\lambda)$, are given by the superposition

$$T_i = \lim_{\Delta\lambda \rightarrow 0} \sum_k C(\lambda_k) t_i(\lambda_k) \Delta\lambda \quad (2.2-8)$$

or

$$T_i = \int C(\lambda) t_i(\lambda) d\lambda \quad i=1,2,3 \quad (2.2-9)$$

Hence, a visual match to $C(\lambda)$ is given by the following weighted sum of primaries $\langle * \rangle$.

$$C(\lambda) \sim \sum_{i=1}^3 T_i p_i(\lambda) \quad (2.2-10)$$

If the color matching functions, $t_i(\lambda)$, are normalized such that the weights required to match a reference white, $W(\lambda)$, are all unity, i.e.

$$T_i = \int t_i(\lambda) W(\lambda) d\lambda = 1 \quad i=1,2,3 \quad (2.2-11)$$

then the weights, T_i , (defined in eqn. 2.2-9) are known as tristimulus values. Here, the notation convention will be

$\langle * \rangle$ The symbol \sim is used to denote a visual, or tristimulus equality, as opposed to a mathematical equality.

to denote scalar weights, such as tristimulus values, by uppercase letters, and the corresponding color matching functions by lowercase letters, (i.e. T_i and $t_i(\lambda)$). Clearly, the tristimulus values associated with a particular spectral distribution depend on the associated set of primaries. However, since all color matching functions are linear combinations of the receptor sensitivities, tristimulus values of different primary systems must also be linearly related. That is, tristimulus values can be transformed from one system to another by use a 3×3 matrix multiplication. Although there are an infinite number of possible color matching functions, some systems are more convenient to use than others. The most popular is the XYZ system $\langle * \rangle$, developed by the "Commission Internationale de L'Eclairage" (or CIE), in 1931 [1, p. 238]. The XYZ color matching functions are shown in fig. 2.2-3. A common way of specifying colors in this system (and other systems as well) is in the form of a chromaticity diagram. This avoids the problem of plotting vectors in a three

 $\langle * \rangle$ Unless stated otherwise, the vectors \underline{T} and $\underline{t}(\lambda)$ will refer to the XYZ coordinate system.

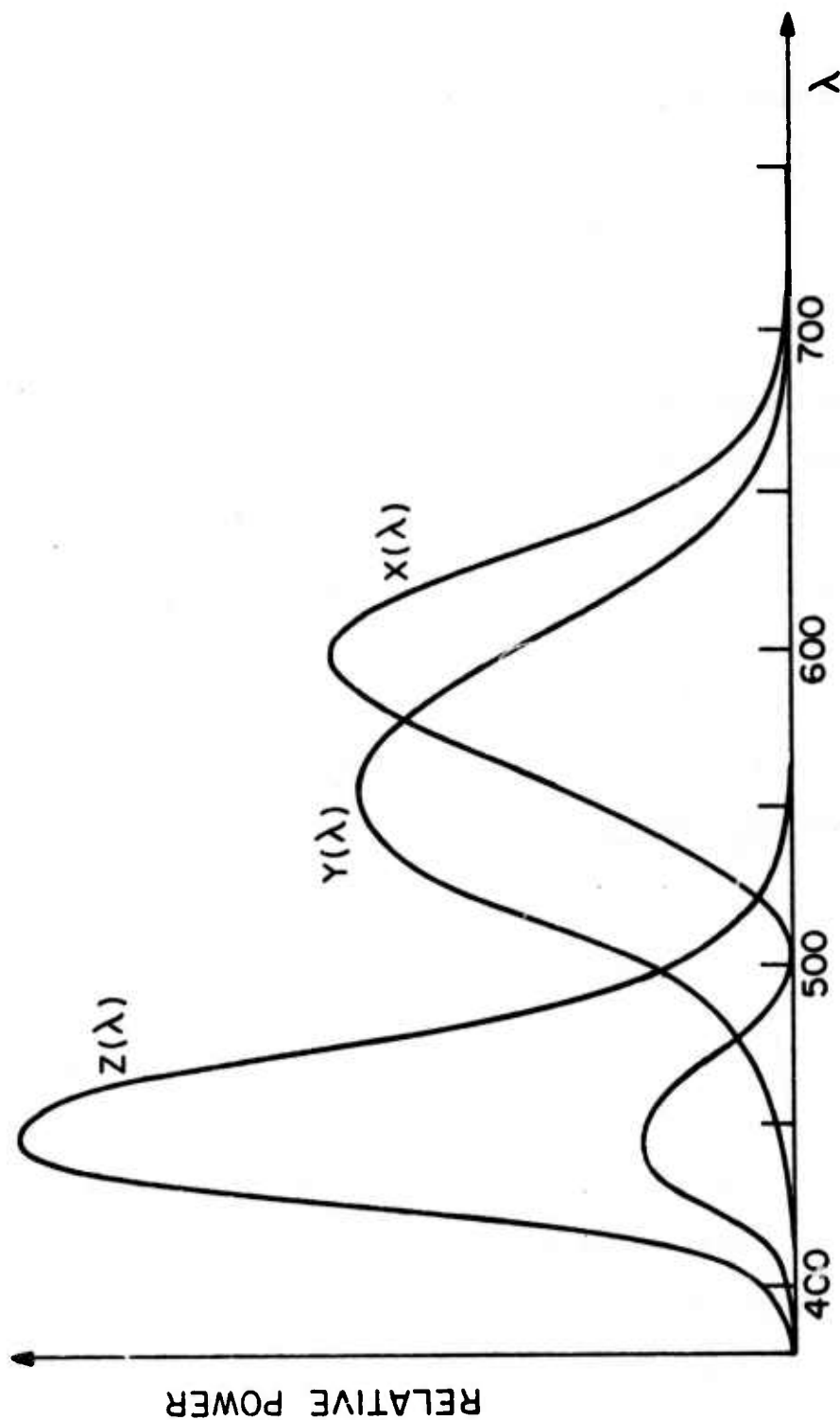


Figure (2.2-3) XYZ color matching functions

dimensional space by use of the transformation

$$x=X/(X+Y+Z) \quad (2.2-12a)$$

$$y=Y/(X+Y+Z) \quad (2.2-12b)$$

thus reducing color specification to two dimensions. This transformation retains hue and saturation information but suppresses lightness. A useful property of the XYZ system is that its color matching functions are nonnegative. This is accomplished by permitting the associated primaries to be non-physically realizable. A further advantage is that one of its color matching curves is the luminous efficiency function, $v(\lambda)$. This function is used to compute the luminance of a spectral distribution, defined as

$$Y = \int v(\lambda) C(\lambda) d\lambda \quad (2.2-13)$$

Its utility is that two colors matching in luminance appear to be of equal lightness.

A major shortcoming of the XYZ system is its unsuitability as a metric space. That is, if the distance between two color vectors, \underline{R} and \underline{S} , is defined by the

standard Euclidean metric

$$d = \sqrt{\sum_i (R_i - S_i)^2} \quad (2.2-14)$$

then unfortunately, there is little correspondence between the perceived disparity between the colors and their metric separation [1, p. 450]. As an illustration, consider the chromaticity diagram of fig. 2.2-4. The horseshoe-shaped line is the locus of spectral colors which are by definition maximally saturated for their particular hue. The ellipses are regions within which the chromaticity can be varied (at constant luminance) without any perceptible color shift <*>. Note that the ellipses in the green region (500-550 nm.) are considerably larger than those of the blue region (400-450 nm.) This means that although colorimetric errors in the blue and green portions of the color space might be equal from a Euclidean distance standpoint, the error in the blue color would appear significantly greater.

Many attempts at finding uniform coordinate systems

<*> The size of the ellipses has been exaggerated by a factor of 10.

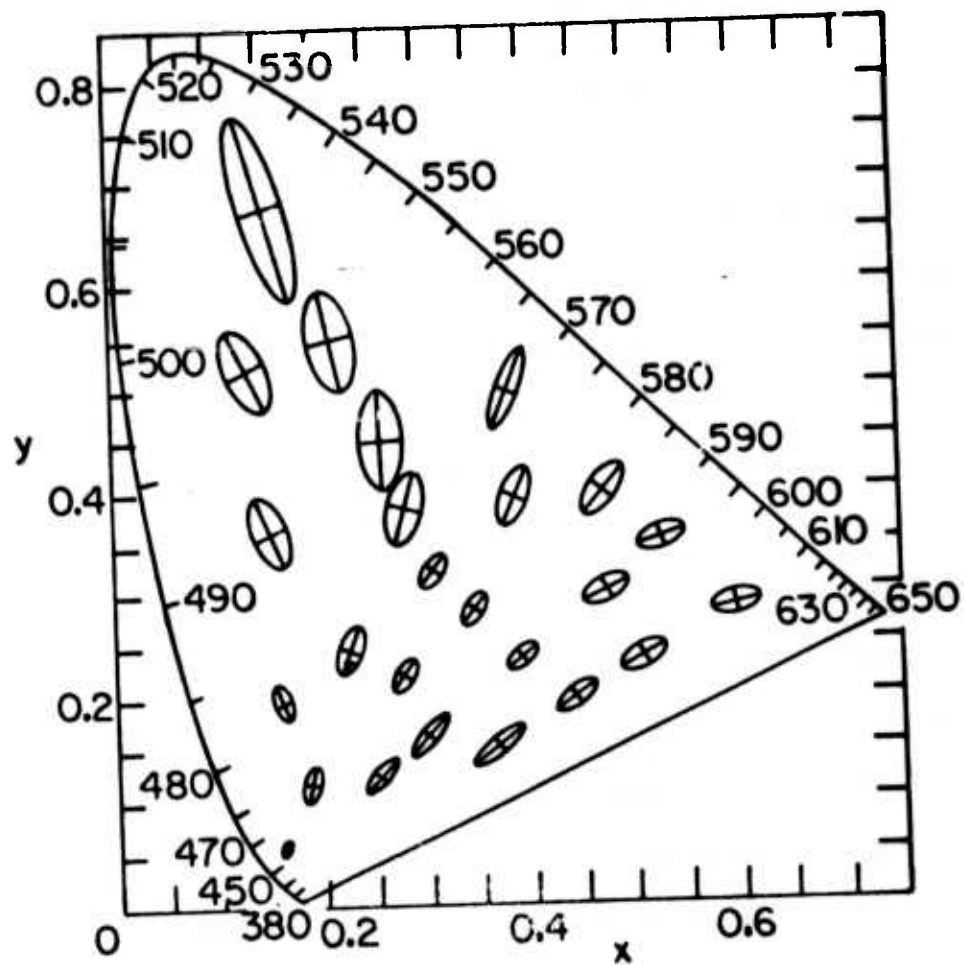


Figure (2.2-4) Chromaticity diagram

which would avoid this difficulty have been reported. These are usually nonlinear transformations of the XYZ system, and are useful in that an error between two colors can be expressed as the Euclidean distance between them [1, p. 450]. At the present time, two such systems are being considered for adoption by the CIE [2]. These are: (1) a modified version of the 1964 $U^*V^*W^*$ space, and (2) a simplified version of the Adams-Nickerson space. The first is defined as

$$L = 25(100Y/Y_0)^{1/3} - 16 \quad (2.2-15a)$$

$$u^* = 13L(u - u_0) \quad (2.2-15b)$$

$$v^* = 13L(v - v_0) \quad (2.2-15c)$$

with

$$u = 4X/(X + 15Y + 3Z) \quad (2.2-16a)$$

$$v = 9Y/(X + 15Y + 3Z) \quad (2.2-16b)$$

where u and v are the coordinates of the reference white (with u_0 and v_0 defined in the same manner as u and v). The second is defined as

$$L=25(100Y/Y_o)^{1/3}-16 \quad (2.2-17a)$$

$$a=500[(X/X_o)^{1/3}-(Y/Y_o)^{1/3}] \quad (2.2-17b)$$

$$b=200[(Y/Y_o)^{1/3}-(Z/Z_o)^{1/3}] \quad (2.2-17c)$$

where X_o , Y_o , and Z_o are the tristimulus values of the reference white to which the eye is adapted $\langle * \rangle$. This coordinate system is tentatively called the CIE 1976 Lab space, in anticipation of its formal adoption on that date. Although neither system is officially preferred by the CIE, it is felt by the author that the Lab space is superior. This conclusion has been reached by generating constant luminance planes in both systems on a calibrated color monitor, and comparing the two planes from the standpoint of uniformity. An example of four constant luminance slices in the Lab system is shown in fig. 2.2-5 (color plate) and fig. 2.2-6. The perimeters of each plane are determined by the constraint that the weights on the display primaries be nonnegative and less than a maximum value. The Lab system has several other

$\langle * \rangle$ A predecessor of the 1976 Lab system is the cube root coordinate system, [3].

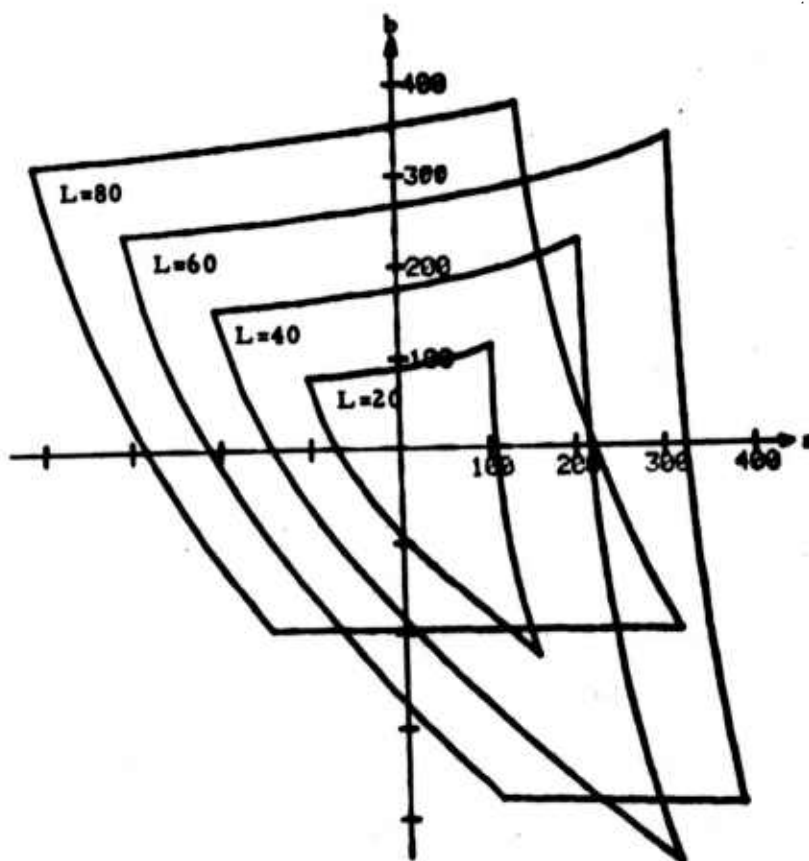


Figure (2.2-6) Television phosphor gamuts in ab plane of Lab space

advantages. First, if considered as a cylindrical coordinate system (see fig. 2.2-7), it can be used to decompose a color into its hue, saturation, and lightness. In this regard, it resembles the empirical "Munsell system" [1, p. 477], which attempts to arrange colors into planes of constant lightness, cylinders of constant saturation, and planes of constant hue. In addition, the Lab system concurs almost exactly with the accepted physiological model of color vision [2]. Recall that the visual model proposed in fig. 2.1-3 transformed the receptor outputs into differences of nonlinear functions. Note that the Lab system also forms differences of nonlinear functions. Therefore, assuming that the receptor sensitivities can be approximated by the XYZ color matching curves $\langle * \rangle$, the Lab system models the physiology of the eye as well. For the above reasons the 1976 Lab system will be used throughout this study to represent colors and colorimetric errors.

$\langle * \rangle$ The nature of the actual receptor sensitivities is discussed in Sec. 2.3.

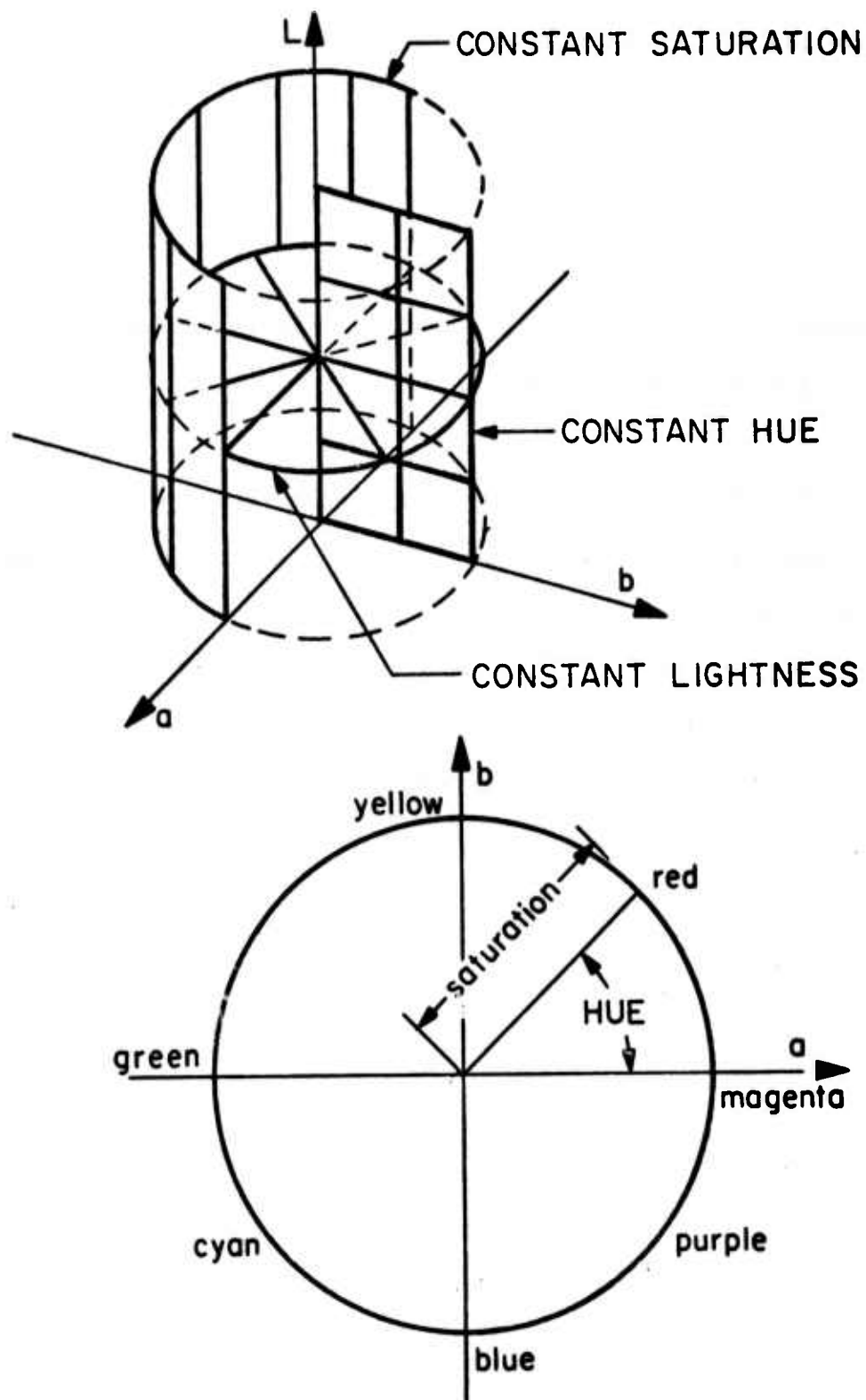


Figure (2.2-7) Lab cylindrical coordinate system

REFERENCES

1. Wyszecki, G. and Stiles, W.S. Color Science, New York: Wiley and Sons. 1967
2. C.I.E. Colorimetry Committee. Proposal for Study of Color Spaces, Tech. Notes in J. Opt. Soc. Am. Vol.64, p. 896, 1974
3. Glasser, L.G., McKinney, A.H., Reilley, C.D., and Schnelle, P.D. Cube Root Color Coordinate System, J. Opt. Soc. Am. Vol.48, p.736, 1958

CHROMATIC ADAPTION

All human perception is to some extent relative. When it is concluded that something is hot, cold, loud or quiet, it is usually with respect to a reference of some sort. The evaluation of color is no exception; a given spectral distribution can be perceived as almost any color, depending on the surround, adaption, etc. In most work involving colorimetry, this difficulty is conveniently evaded by stipulating identical viewing conditions for the comparison of colors. Unfortunately, for most practical situations such as photography, the original scene and the reproduction are invariably viewed under different conditions. The preferred white point for color television is about 6500°K , (i.e., the tristimulus values of the white point are the same as those of a black body radiator at a temperature of 6500 degrees Kelvin). On the other hand, the white point of most slide projectors is closer to 2800°K . If these two "whites" are viewed side by side, the TV will appear blue and the slide projector yellow. However, if viewed individually, with dark surrounds, the eye will quickly adjust to the balance of either color, and accept it as white. The phenomenon

is enormously complicated and exhibits temporal, spatial, and even psychological aspects.

This remarkable ability of the human visual system to discount the effects of the viewing illuminant is known as chromatic adaption. Qualitatively, the appearance of a color is determined by the direction and magnitude of its deviation from the reference white, where "reference white" describes either the viewing illuminant, or the prevailing color balance of the scene being viewed. Quantitatively, a description is much more difficult. Although some very precise experimental studies have been undertaken [1, p. 435], no unified theory explaining all the experimental results has emerged. The inconsistencies in the published data are doubtless due to differences in the experimental technique used, and individual variations among observers. Nonetheless, the colorimetric shifts due to adaption appear to be quite systematic. For instance, in the recent study by Sobagaki, et al. [2], the influence on color perception by a change in the illuminant from daylight to tungsten was investigated using 14 observers and 95 sample colors. Although there was a noticeable variance in the color shifts perceived by different individuals, there was a definite trend. The

plot of fig. 2.3-1 illustrates their results for eight colors <*>. Each of the individual arrows in a given cluster corresponds to the chromatic shift perceived by a particular observer, for daylight to tungsten adaption. For instance, color number 8 (Munsell notation 10 Purple) tends to shift towards red. Note that the color shifts appear as a counterclockwise rotation in color space.

Although the investigators in the Sobagaki study restricted attempts to fit their data to linear transformations, it is usually conceded that a nonlinear transformation is generally required to model the phenomenon with acceptable accuracy [3]. The traditional linear hypothesis purporting to explain adaption, due to Von Kries [3], is represented by the "gain controls" or variable scale factors in each of the color receptors in fig. 2.3-2. This hypothesis asserts that the gain controls are varied in such a way as to insure that the integrated responses from each of the three receptors are scaled to unity when the reference white is being viewed

<*> The original Sobagaki data was plotted in the uv coordinate system, and has been converted to the Lab system used throughout this dissertation.

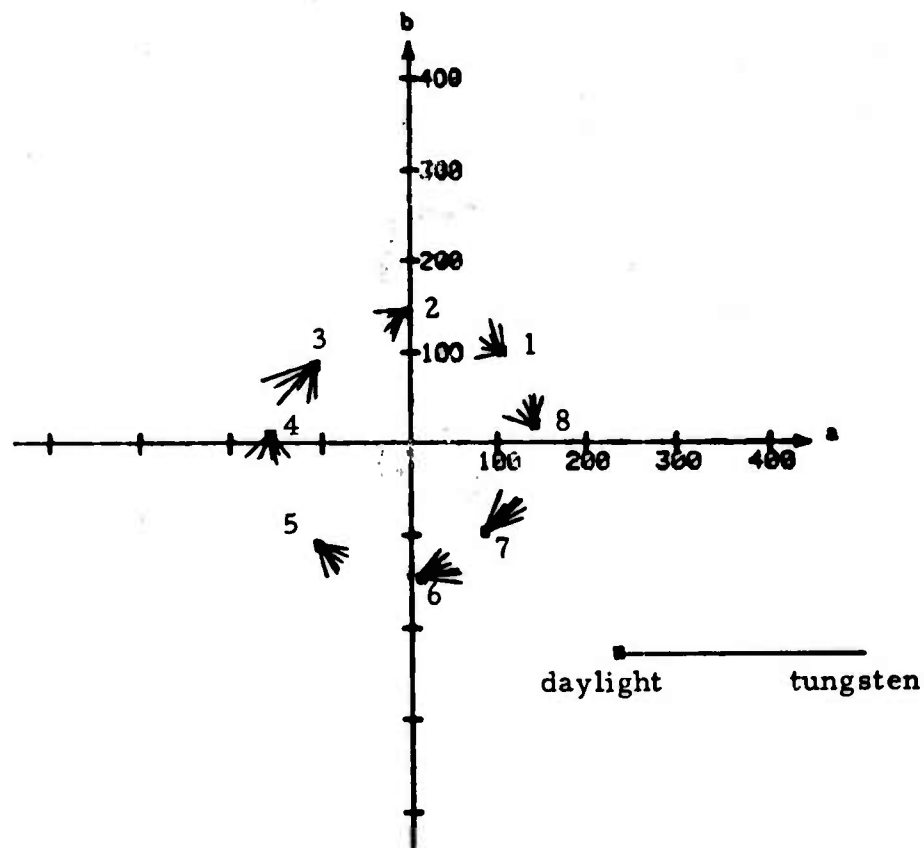


Figure (2.3-1) Color shifts due to chromatic adaption

i.e.,

$$k_j = \frac{1}{\int s_j(\lambda) W(\lambda) d\lambda} \quad (2.3-1)$$

where $W(\lambda)$ is the spectral distribution of the reference white, and $s_j(\lambda)$ is the spectral sensitivity of the j th visual receptor. Equivalently, the response of each receptor can be formulated as <*>

$$S_j = \frac{\int s_j(\lambda) C(\lambda) d\lambda}{\int s_j(\lambda) W(\lambda) d\lambda} \quad (2.3-2)$$

The validity of this attractive linear hypothesis has been attacked on a number of grounds. Historically, interest in the Von Kries hypothesis was motivated by the need to calculate the actual visual receptor sensitivities. Theoretically, this can be accomplished by performing an eigenvector analysis on color matching data taken under

<*> Note that if the receptor sensitivities, $s_j(\lambda)$, are assumed to be equivalent to the XYZ color matching functions, this type of scaling is identical to that used in the Lab coordinate system defined in eqn. 2.2-17.

two different states of adaption [1, p.438]. Denoting the actual (physiological) receptor signals evoked by the same color viewed under two different states of adaption by the vectors \underline{S} and \underline{S}' , the Von Kries hypothesis asserts that

$$\underline{S}' = D \underline{S} \quad (2.3-3)$$

where D is a diagonal matrix. Denoting the experimentally determined tristimulus vectors corresponding to the two different adaption states as \underline{T} and \underline{T}' , the physiological tristimulus vectors must be related as

$$\underline{S} = A \underline{T} \quad (2.3-4)$$

and

$$\underline{S}' = A \underline{T}' \quad (2.3-5)$$

where A is an unknown 3×3 matrix. The relationship between the \underline{T} and \underline{T}' vectors should also be linear, and thus, a 3×3 matrix B can be fitted to the experimental data such that

2.3

$$\underline{T}' = B \underline{T} \quad (2.3-6)$$

Substituting equations 2.3-4, 2.3-5, and 2.3-3, yields

$$\underline{T}' = A^{-1} D A \underline{T} \quad (2.3-7)$$

Therefore

$$B = A^{-1} D A \quad (2.3-8)$$

or

$$D = A B A^{-1} \quad (2.3-9)$$

This means that the matrix A , which diagonalizes the experimentally determined matrix B , can be used to determine the actual receptor sensitivities. Eqn. 2.3-4 implies

$$\underline{s}(\lambda) = A \underline{t}(\lambda) \quad (2.3-10)$$

where the $t_j(\lambda)$ are the color matching functions associated with the primaries used in the matching

experiments. Although the diagonalization of 2.3-9 can be accomplished by conventional eigenvector techniques, there is no guarantee that the resulting matrix, A , will be real. Indeed, this technique for finding the fundamental sensitivities, $\underline{s}(\lambda)$, is plagued with complex eigenvector solutions. The conclusion seems to be that the Von Kries hypothesis itself must be somewhat simplistic.

In order to investigate alternative color adaption mechanisms, it is necessary to have some data on the fundamental sensitivities of the eye. Fortunately, estimates of the $s_j(\lambda)$ can be obtained by alternative methods, such as color blindness studies [4]. A number of fundamental response curves have been proposed. One set which was deduced from color blindness data, and is in good agreement with more recent curves obtained directly by Wald, is due to Konig <*>. These fundamentals are defined in terms of the XYZ color matching functions by the transformation

<*> A good summary of this work is given in [4], and [5].

$$\begin{bmatrix} S_1 \\ S_2 \\ S_3 \end{bmatrix} = \begin{bmatrix} .070 & .940 & -.010 \\ -.477 & 1.409 & .068 \\ .000 & .000 & 1.000 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} \quad (2.3-11)$$

The corresponding sensitivities are plotted in fig. 2.3-3. The König fundamentals will be used to simulate various chromatic adaption models, in order to compare the results with experimental data.

An alternative to the Von Kries approach is the hypothesis that some form of adaption takes place subsequent to the compression-type nonlinearities depicted in fig. 2.3-2. For the case of a variable gain taking place immediately following a logarithmic nonlinearity, the net result is effectively a power law scaling

$$\begin{aligned} S' &= a \log(S) \\ &= \log(S^a) \end{aligned} \quad (2.3-12)$$

This type of structure is roughly equivalent to the nonlinear mechanism proposed by MacAdam to model chromatic adaption [3].

A hypothesis which has received little attention is a

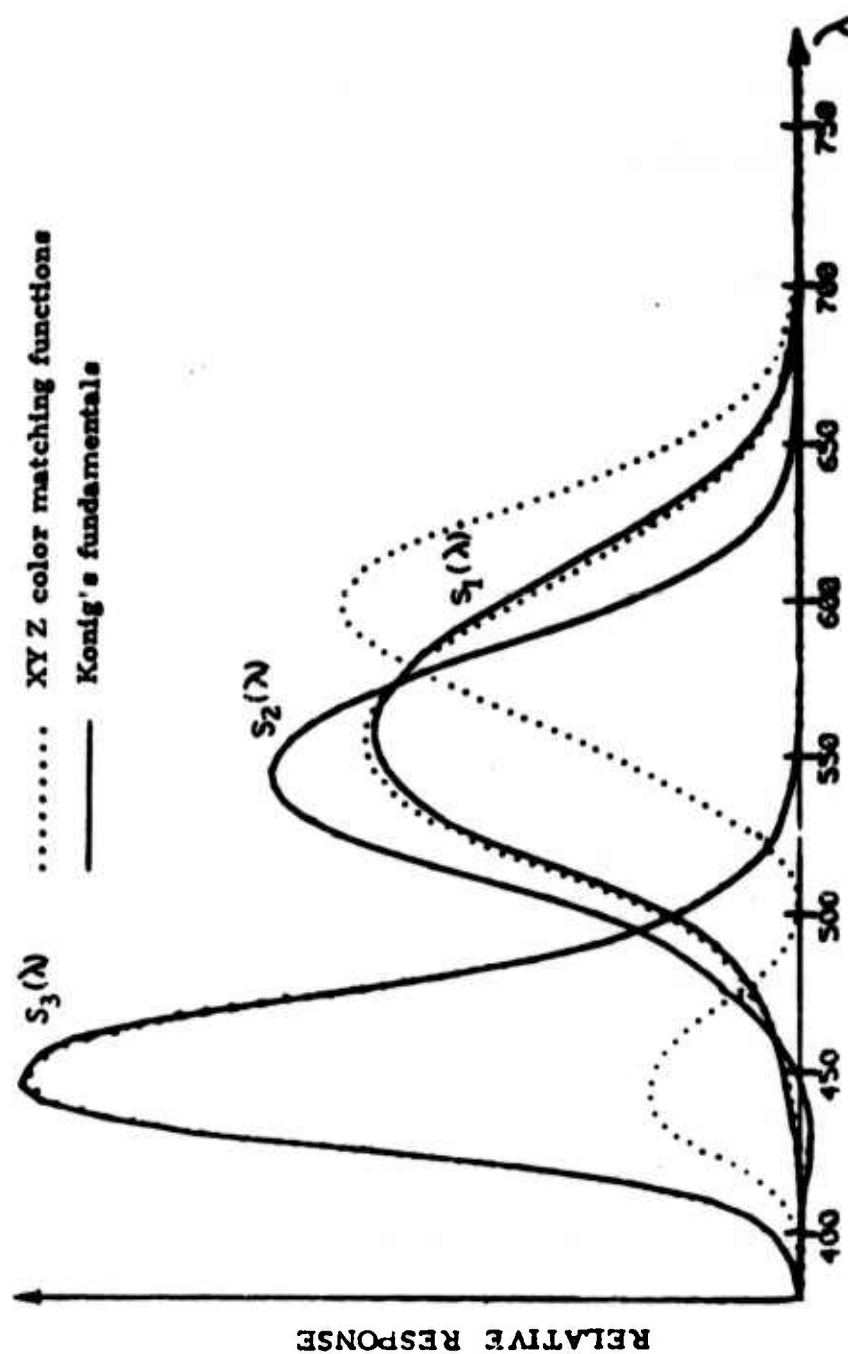


Figure (2.3-3) Konig's receptor sensitivities

2.3

process based on spatial frequency considerations and homomorphic filtering [6]. The rationalization for such an approach is the multiplicative nature of image formation. That is, the perceived spectral distribution is generated by multiplying the reflectance function by the illuminant.

$$C(\lambda) = I(\lambda) R(\lambda) \quad (2.3-13)$$

where $I(\lambda)$ is the illuminant and $R(\lambda)$ is the spectral reflectance. Altering the illuminant is approximately equivalent to imposing multiplicative biases on the "signals" generated by the visual receptors. Furthermore, a logarithmic nonlinearity transforms products into sums, changing the biases to additive form. Therefore, the eye's ability to "discount the illuminant" could be due to a neural inhibition mechanism which linearly filters out slowly varying components in the visual field (or equivalently, exaggerates sudden changes).

Consider the structure of the visual model in fig. 2.1-3. There are three separate spatial filters, one

for each channel. The high-emphasis filter operating on the luminance channel is responsible for monochrome spatial effects, such as black and white Mach bands. The fact that chromatic Mach bands are weaker indicates that any spatial filters operating on the chromatic channels must be of a different nature. However, if these filters in the chromatic channels basically reject zero spatial frequency, and then "roll off" at the middle frequencies (see fig. C-1), a definite adaption mechanism would take place. A crude, first order approximation to such an effect is illustrated by the subtractive biases in fig. 2.3-4. This corresponds to an unrealistic spatial filter which completely rejects zero spatial frequency and passes all other spatial frequencies with unity gain $\langle * \rangle$. For the case in which all three compressive nonlinearities are perfect logarithms, this formulation is exactly equivalent to the Von Kries hypothesis. However, a perfect logarithmic response is clearly unrealistic [7], and therefore the effect of the homomorphic process

$\langle * \rangle$ Clearly, zero spatial frequencies cannot be completely rejected, because this would imply that all completely uniform fields would be colorimetrically neutral.

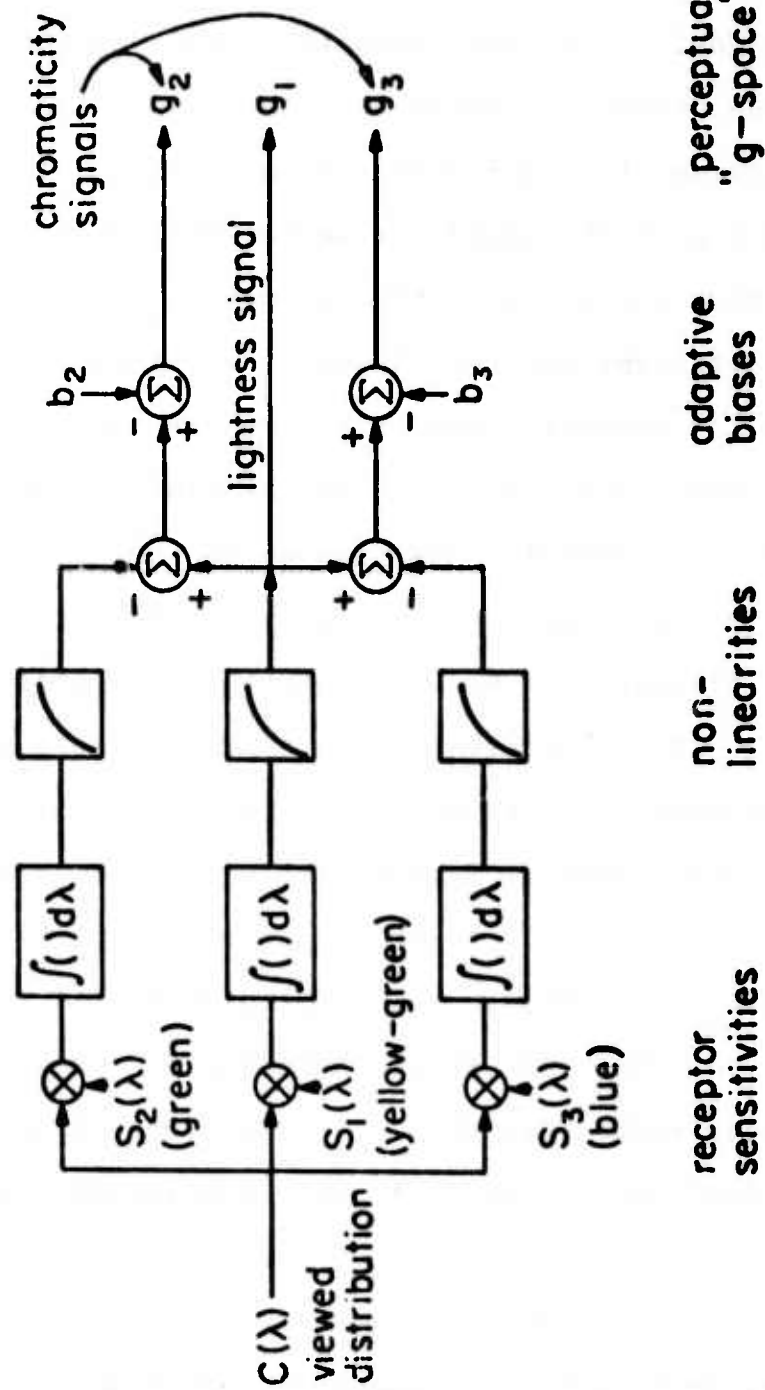


Figure (2.3-4) Adaptive bias model for chromatic adaption

2.3

described above will differ somewhat from the simple scaling postulated by Von Kries, especially at low brightness levels. A graphical interpretation of this type of normalization is shown in fig. 2.3-4. When plotted in the perceptual "g-space" <*> (see figs. 2.3-4, and 2.3-5), the biases resulting from the zero frequency rejection cause a translational shift of the coordinate axes. Thus, this formulation would predict that the eye would tend to "center its coordinate system" at the white point to which it is adapted. For the case in which an array of colors is being viewed, the effective white point would be the statistical average of the colors in the scene. For colors which deviate markedly from the "preferred white" (approximately 6500° K), it is likely that the eye never completely adapts, but reaches a sort of intermediate compromise. This effect can be seen by viewing a low color temperature tungsten bulb; at the usual viewing lightnesses it will always retain a yellowish tinge.

<*> The perceptual "g-space" differs from Lab space in that the actual receptor sensitivities, $s_j(\lambda)$, are used, as opposed to the CIE color matching functions, $t_j(\lambda)$. Otherwise, the mathematical structure is the same.

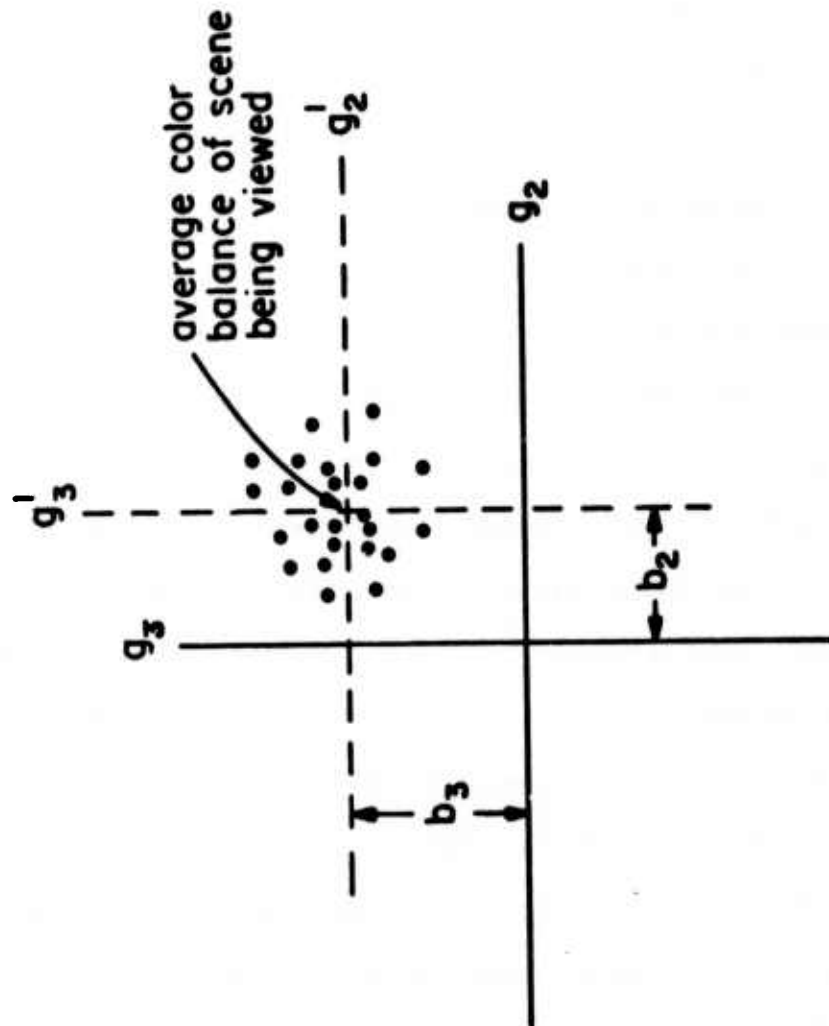


Figure (2.3-5) Chromatic adaption interpreted as a self-centering coordinate system

2.3

Using the Lab coordinate system, plots of the colorimetric shifts predicted by the mechanisms described above were generated. The presumed form of the nonlinearity was a cube root, and the assumed change in illuminant was from daylight to tungsten (2854°K). The effect of placing the gains before the nonlinearities is shown in fig. 2.3-6. The shifts predicted by the subtractive bias hypothesis, for lightnesses of 40, 60, and 80 are shown in figs. 2.3-7, 2.3-8, and 2.3-9. Note that the predicted shifts are a strong function of lightness. This phenomenon is of sufficient interest to merit further discussion.

If the nonlinearities in the subtractive bias model were perfect logarithms, then the predicted color shifts would not be functions of the absolute level of the adapting illumination, since a logarithm transforms any gain into an additive term. However, for a more realistic nonlinearity, such as a cube root, some interesting level-dependent phenomenon take place. For instance, assume that the subtractive biases in the coordinate centering mechanism described above are adjusted so that a particular illuminant (at a given lightness level) is perceived as white. If the lightness level of the

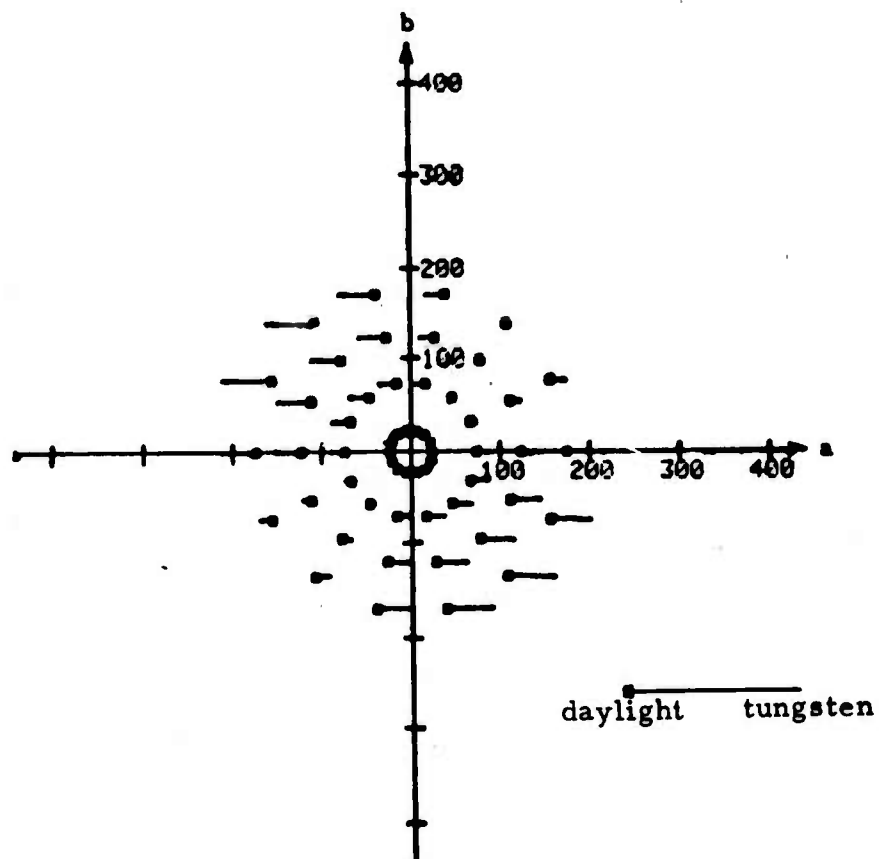


Figure (2.3-6) Effect of pre-nonlinearity gains on adaption shifts for one-third power nonlinearity

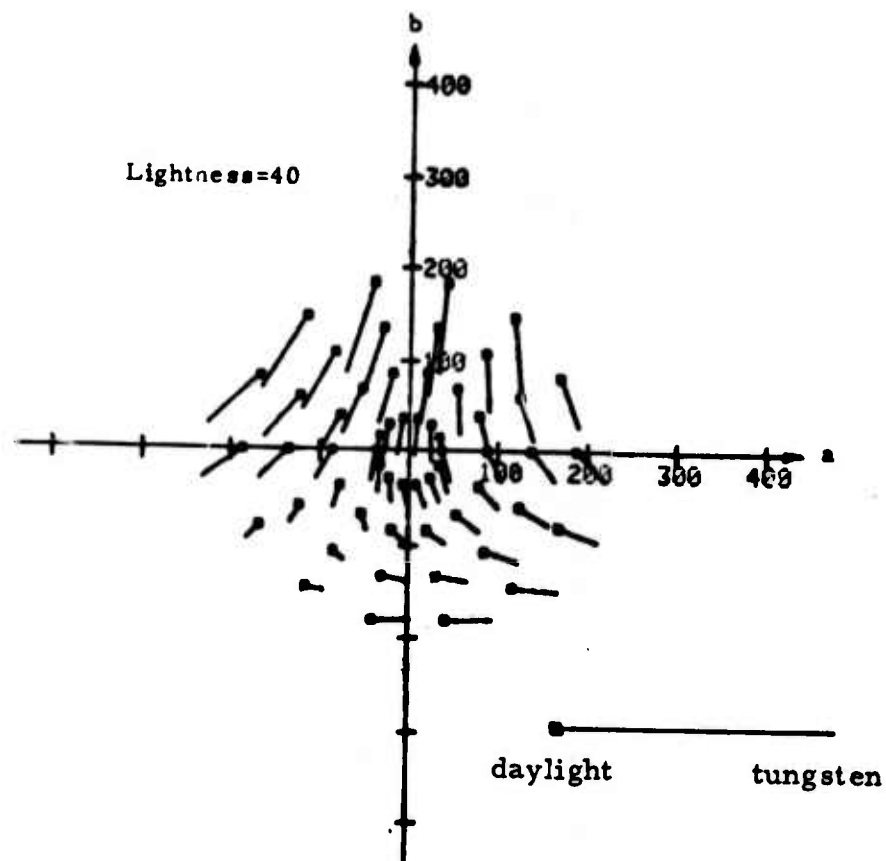


Figure (2.3-7) Effect of post-nonlinearity biases at low lightnesses, for one-third power nonlinearity

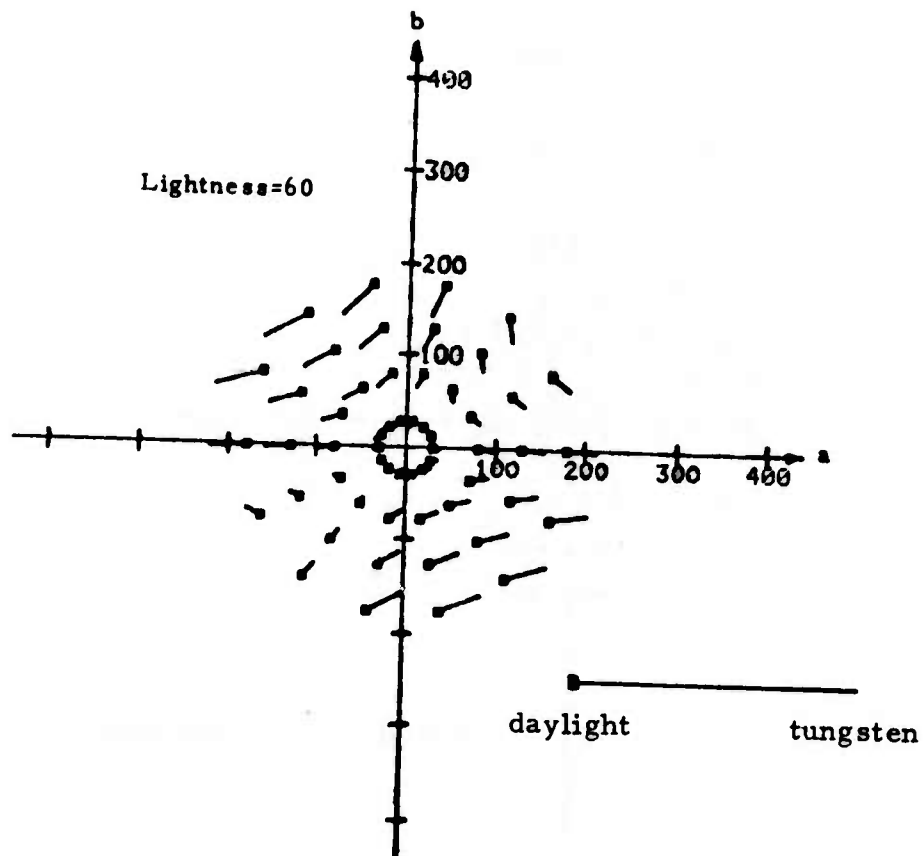


Figure (2.3-8) Effect of post-nonlinearity biases at medium lightnesses, for one-third power nonlinearity

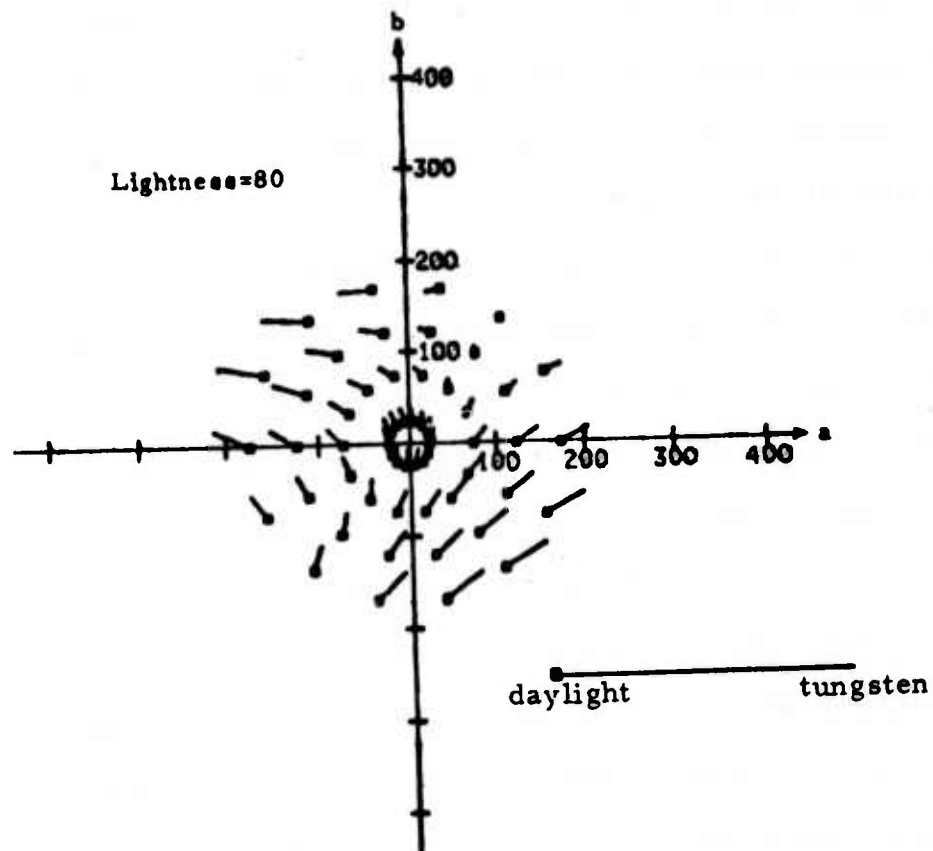


Figure (2.3-9) Effect of post-nonlinearity biases at high lightnesses, for one-third power nonlinearity

reference white is changed, slightly different subtractive biases would be required to adapt to it (unless the nonlinearities are perfect logarithms). This implies that color shifts should depend on whether the lightness of the color being viewed is greater than or less than the lightness of the adapting field. That is, if a complex field is being viewed, the adapting white can be considered the statistical average of the distributions being viewed (the zero spatial frequency component). Therefore, lighter colors should undergo different shifts than their darker counterparts. In the plots of figs. 2.3-7, 2.3-8, and 2.3-9, it was assumed that the nonlinearity was a cube root, and that the lightness of the adapting white was 60 (on a scale of 100). Note that the lighter hues ($L=80$) tend to shift towards the illuminant (yellow), while the darker hues ($L=40$) tend to shift towards its complement (blue). Significantly, these are precisely the types of shifts which are perceived in complex visual fields. By comparing these predictions with the empirical results illustrated in fig. 2.3-1, it is evident that a combination of the pre-nonlinearity scaling and the subtractive bias transform can be used to approximate chromatic adaption effects. This is

2.3

reasonable from a physiological standpoint, in that it is commonly thought that adaption proceeds with both a slow and fast component. A scaling in the early stages of the visual system might be due to chemical phenomenon, (such as bleaching), with long time constants, while the neural inhibition responsible for the subtractive bias component would be essentially instantaneous. If such multiple adaption mechanisms are indeed present in the visual system, the variance in results reported by different experimenters attempting to quantify the phenomenon might be explained by the dominance of one mechanism over the other.

In this dissertation the approach taken will be to divide the adaption effects between a pre-nonlinearity scaling (Von Kries hypothesis), and a subtractive biasing. One method which can be used to formulate this, is to modify the coefficients k_j of eqn. 2.3-1 by an inverse power law compression, so that their effect on the adaption is incomplete.

$$k'_j = k_j^{1/n} \quad (2.3-14)$$

The remainder of the adaption can be performed by the

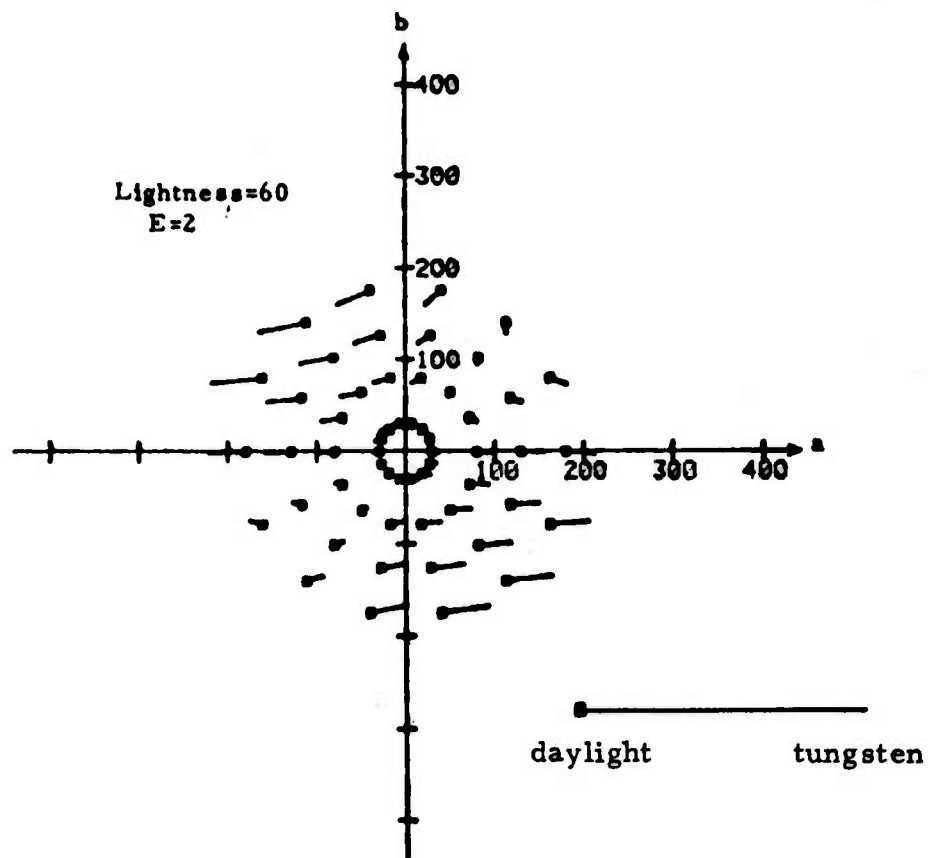


Figure (2.3-10) Color shifts predicted by combined Von Kries and subtractive bias effects

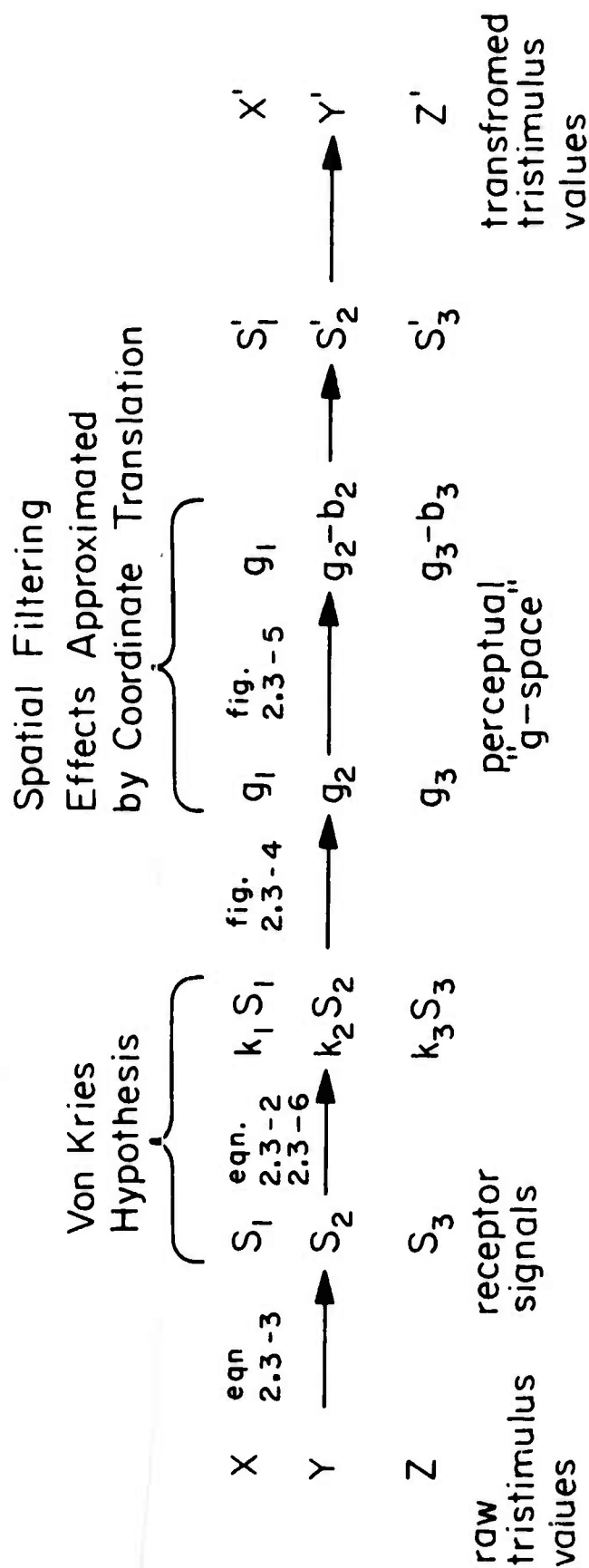


Figure (2.3-11) Block diagram of adaption transform

2.3

subtractive bias mechanism. In this manner, the parameter n can be used to adjust the proportion of the adaption effect contributed by either mechanism.

By using an n -parameter of 2 (square root), a family of color shifts results which corresponds closely with those found experimentally (see fig. 2.3-10). A block diagram of the steps used in simulating adaption-caused color shifts is provided in fig. 2.3-11. This is the chromatic adaption transformation that will be referred to in later chapters.

REFERENCES

1. Wyszecki, G. and Stiles, W.S. Color Science, New York: Wiley and Sons. 1967
2. Sobagaki, H., et. al. Chromatic Adaption Study by Subjective Estimation Method, J. Opt. Soc. Am. Vol.64, p.743, 1974

2.3

3. MacAdam, D.L. A Nonlinear Hypothesis for Chromatic Adaption, Vision Research, Vol. 1, p. 9, 1961
4. Frei, W. Modeling Color Vision for Psychovisual Image Processing, USCEE Rept. 459 p. 112. 1973
5. Frei, W. A New Model of Color Vision and Some Practical Implications, USCIPR Report 530, p. 128, 1974
6. Stockham, T.G. Image Processing in the Context of a Visual Model, Proc. IEEE, Vol. 60, p. 828, 1972
7. Savoie, R.E. Bezold-Brucke Effect and Visual Nonlinearity, J. Opt. Soc. Am. Vol. 63, p. 1253, 1973

ADDITIVE COLOR SYSTEMS

The previous chapter indicated how a wide gamut of colors could be evoked by viewing linear superpositions of three spectral distributions. Color reproduction systems based on the summation of spectra are termed additive, as opposed to systems based on the mixing of dyes, which are known as subtractive. Although subtractive systems are more common, additive systems are substantially easier to analyze, and will be discussed first.

THE ADDITIVE PRINCIPLE AND COLOR TELEVISION MONITORS

The most common example of an additive display system is the color television monitor, with the most popular variety of tube being the three-gun aperture (or shadow) mask type. In this, the screen contains three types of phosphors, (red, green, and blue) in the form of individual "dots" deposited on the inner surface. The dots are arranged in a mosaic of triangular groups or "triads," and obscured from the three electron guns by a

3.1

stencil, or mask [1]. The geometry of the guns, mask, and dots prevents the red gun from exciting any but the red phosphor dots, the green from exciting any but the green dots, etc. Thus, by superposition, any given triad on the screen can be made to elicit the spectral distribution

$$C(\lambda) = \sum_{j=1}^3 P_j p_j(\lambda) \quad (3.1-1)$$

where the $p_j(\lambda)$ are the spectral characteristics of the red, green, and blue phosphors, and the P_j are weights, roughly proportional to the drive signals on the guns. One drawback to this type of tube is the substantial amount of energy lost in the mask, causing inefficient operation. A second and perhaps more serious disadvantage is the inherent difficulty of maintaining alignment between the phosphor dots and the electron beam as it is swept across the screen (line convergence). In quantitative work, a significant problem is the management of the intrinsic nonlinearities of the cathode ray tube. Chief among these is the relationship between the drive signals to the guns and the resulting brightnesses of the stimulated phosphors. A typical set of such relationships is shown in fig. 3.1-1. The curves can be approximated by

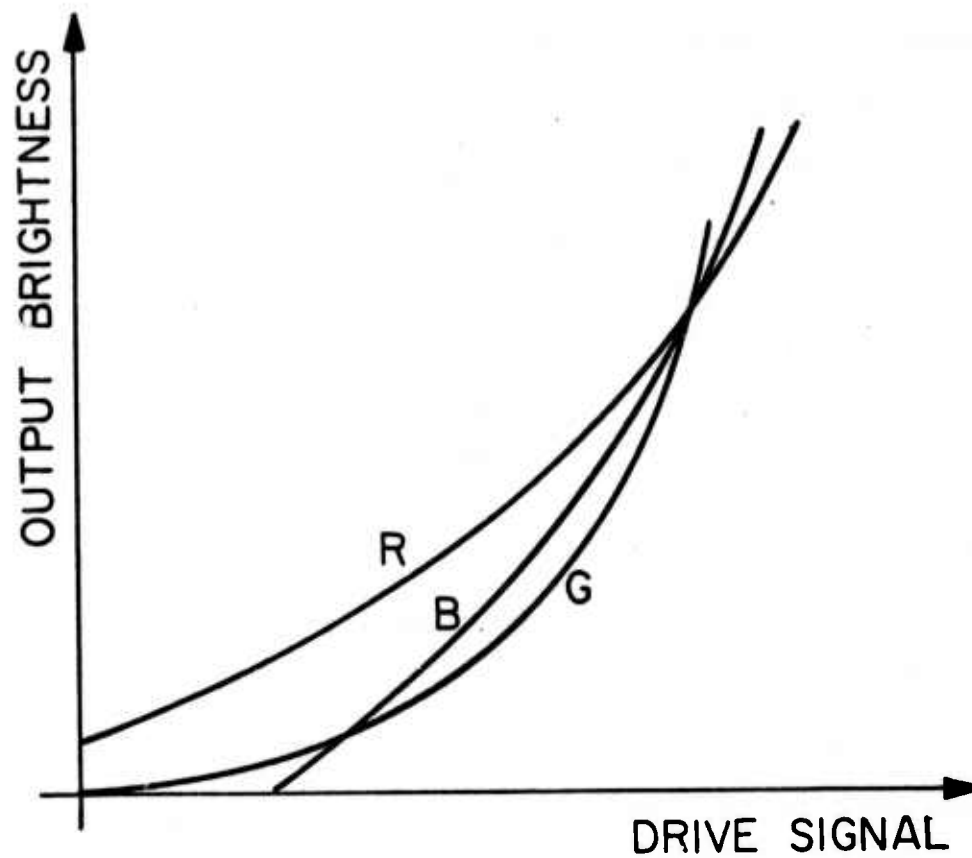


Figure (3.1-1) Nonlinearities of color monitor display

3.1

the power law <*>

$$b = s^{\gamma} \quad (3.1-2)$$

where s is the input drive signal, γ is an exponent and b is the resulting brightness. Such power law relationships are usually called gamma curves. Clearly, if the three curves do not overlap each other exactly, then the grey scale "tracking" will be poor. That is, what may be intended as a uniform grey scale, will not only be non-uniform, but will exhibit color shifts as well. The gamma curves can be modified in three ways. First, the brightness control of the monitor can be used, which has the approximate effect of raising or lowering all three curves on a pedestal. Secondly, the three "cutoff" controls (one for each gun) can be manipulated, causing a similar effect in an individual curve. Thirdly, the red, green, and blue gain controls can be adjusted singly, changing the curve's slope. The effects of the cutoff and

<*> Eqn. 3.1-2 assumes that the brightness control has been adjusted such that zero input gives zero output.

3.1

gain controls on a single curve are shown in fig. 3.1-2. Since there is always a substantial variance between different monitors, the gamma curves for an individual monitor and a particular set of adjustments cannot be specified in general, and must be measured with a photometric device. Unfortunately, it is impossible to linearize the gamma curves by monitor adjustments, and the only recourse is a pre-distortion of the drive signals by the inverse of the gamma curves. In a digital processing environment this can be easily accomplished with a "look up table." Such a procedure is known as gamma correction and is of great importance if the image is to be photographed. This is because the gamma curve has the effect of exaggerating the image contrast, and thus straining the limited dynamic range of color film. This results in a loss of detail in the shadows, and a "washing out" of the highlights (see fig. 3.1-3 in color plate).

A second consideration in quantitative work is the character of the phosphors. If the exact spectral distributions of a particular set of phosphors are needed, they must be measured with a spectro-radiometer. A set of typical curves for modern red, green, and blue phosphors is shown in fig. 3.1-4. It should be noted that the

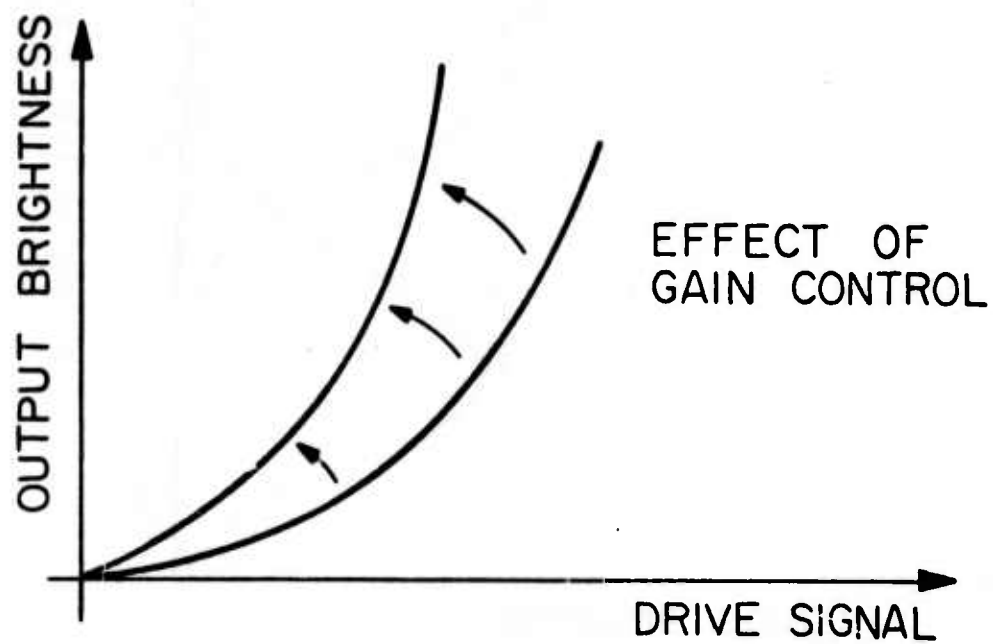
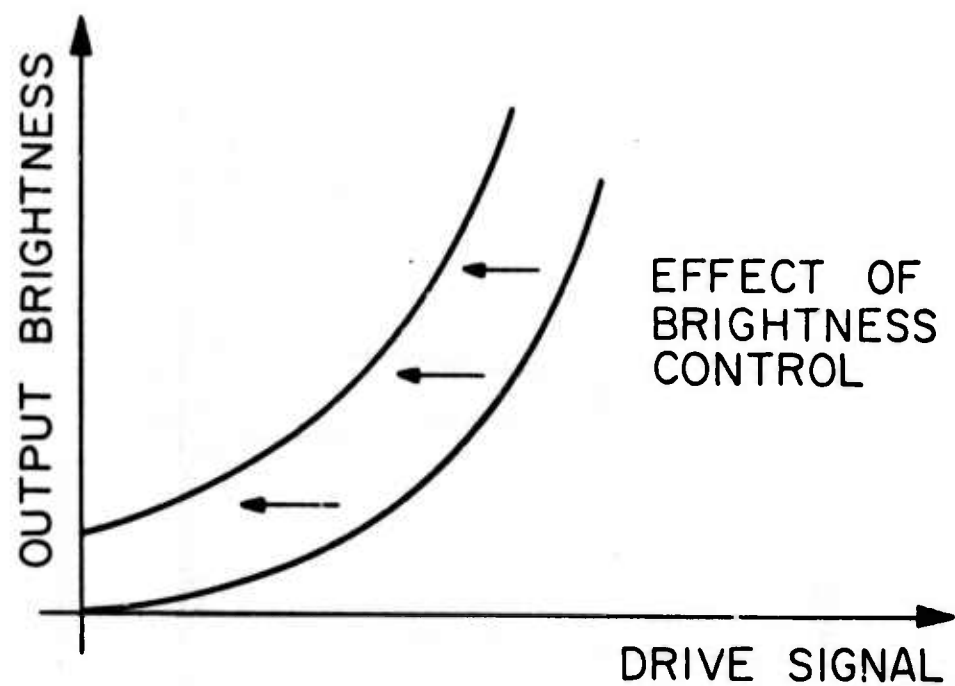


Figure (3.1-2) Effect of brightness and gain controls

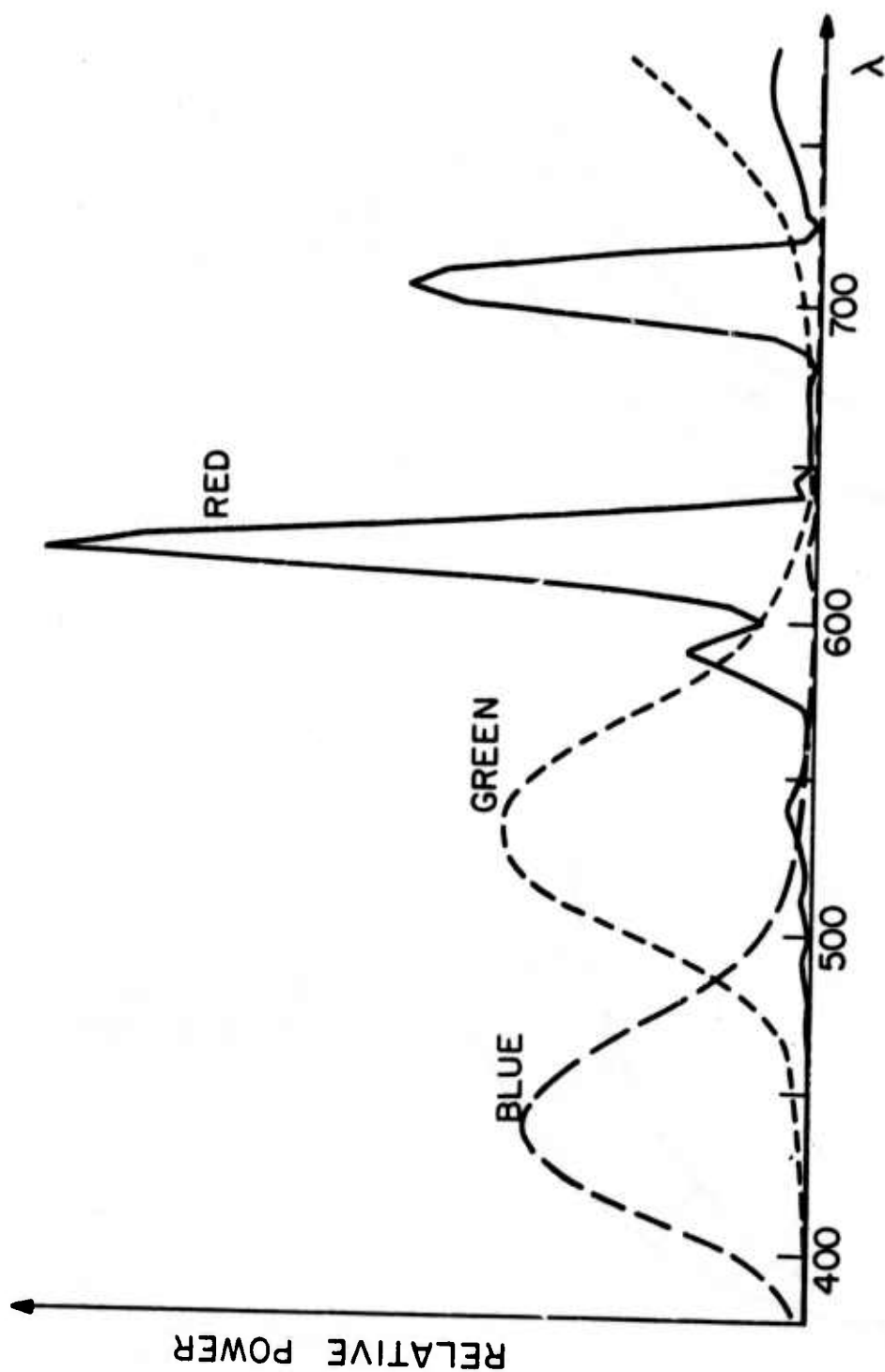


Figure (3.1-4) Typical television phosphor characteristics

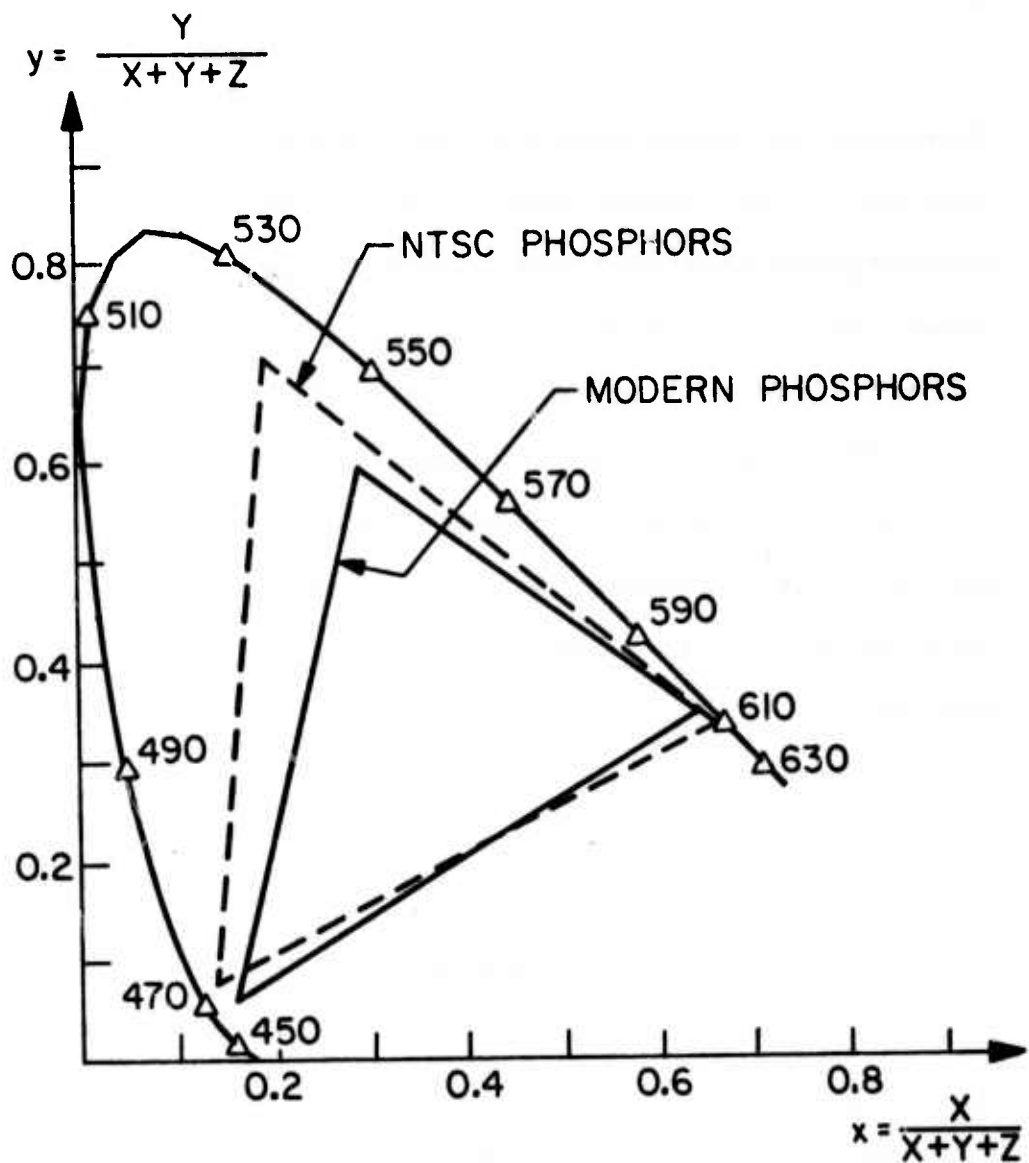


Figure (3.1-5) Comparison of modern phosphors and N.T.S.C. phosphors

3.1

chromaticity coordinates of present day phosphors differ somewhat from those specified by the N.T.S.C. [2]. A chromaticity plot contrasting the two sets of phosphors is shown in fig. 3.1-5. The phosphor distributions are of concern since they are needed to predict the response of eye and film to the various colors generated by the monitor. The range of colors which the phosphors are capable of generating is known as the gamut of reproducible colors, and is the subject of the next section.

REFERENCES

1. Hunt, R.W.G. The Reproduction of Colour, London: Wiley and Sons. 1967
2. DeMarsh, L.E. Color Reproduction in Color Television, in the publication; Color: Theory and Imaging Systems, Society of Photographic Scientists and Engineers, p. 184, 1973

GAMUT AND LUMINANCE RESTRICTIONS

An important attribute of any color reproduction system is the range, or gamut of colors which it can generate. In a qualitative sense, a display limited to only drab, muddy, colors would be considered to have a narrow gamut, while one capable of bright, pure colors would be said to possess a wide gamut. In the case of additive systems, a quantitative description is quite straightforward. Consider, for example, an additive display using red, green, and blue primaries. An arbitrary distribution generated by such a device would be

$$C(\lambda) = R r(\lambda) + G g(\lambda) + B b(\lambda) \quad (3.2-1)$$

where R , G , and B are weights, and $r(\lambda)$, $g(\lambda)$, and $b(\lambda)$ are the spectral distributions of the primaries. Since a negative amount of primary is physically meaningless, and power limitations impose a maximum limit on the weights, (say unity), the gamut can be represented as a unit cube in RGB space (see fig. 3.2-1). That is, only points interior to the cube can be generated by the display. However, such a representation indicates nothing about the

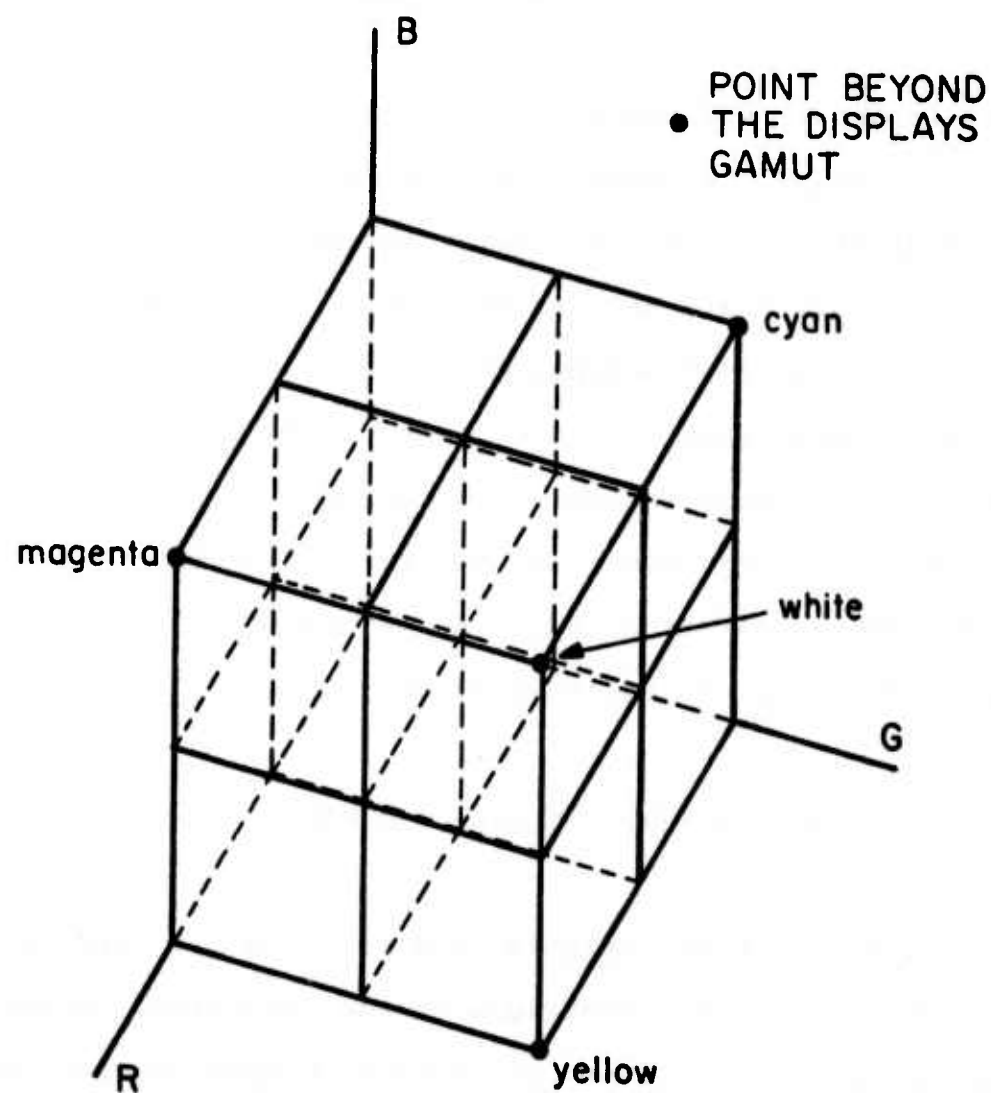


Figure (3.2-1) Gamut of additive system in RGB space

purity of the primaries. Thus, in order for different display systems to be compared meaningfully, their gamuts must be expressed in some common coordinate system, such as the XYZ space. This can be accomplished by the linear transformation

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} & & \\ & a_{ij} & \\ & & \end{bmatrix} \begin{bmatrix} R \\ G \\ B \end{bmatrix} \quad (3.2-2)$$

where the coefficients a_{ij} are determined by the spectral characteristics of the primaries $r(\lambda)$, $g(\lambda)$, and $b(\lambda)$ (see appendix A). The transformation for typical modern TV phosphors balanced to a white point of 6500°K is

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} .489 & .324 & .161 \\ .264 & .672 & .064 \\ .014 & .134 & .806 \end{bmatrix} \begin{bmatrix} R \\ G \\ B \end{bmatrix} \quad (3.2-3)$$

A point-by-point transformation of the unit cube of fig. 3.2-1 using eqn. 3.2-3 yields the distorted solid of fig. 3.2-2. This is the gamut of reproducible colors in XYZ space. A second, and perhaps more meaningful

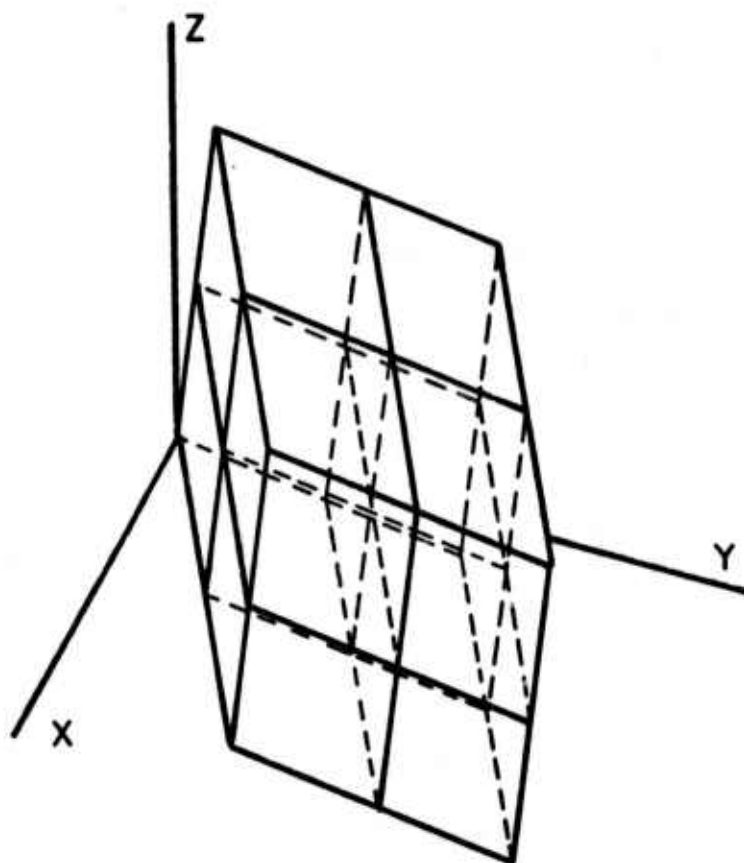


Figure (3.2-2) Gamut of additive system in XYZ space

representation is attained by transforming the RGB cube to a uniform color space, such as the Lab system (see eqn. 2.2-17 and figs 2.2-6, and 2.2-7). The resulting color solid is shown in fig. 3.2-3. The warped and distorted faces of the solid are due to the fact that the RGB to Lab transform is nonlinear. Note that the Lab cube tapers to a point at both low and high lightness regions. That is, highly saturated colors are possible only in the middle lightness ranges. For dim colors, the saturation is limited by the constraint that the weights must be positive. In the case of bright colors, the limits are imposed by the constraint that the weights be less than unity. Four constant lightness planes from the Lab solid are shown in fig. 2.2-5 (color plate). Note that the darker colors are red, green and blue, while the lighter colors are cyan (blue plus green), magenta (red plus blue), and yellow (red plus green). Although in natural scenes the occurrence of very bright, or saturated colors is statistically small (see fig. 3.2-4), the problem of dealing with colors beyond the display's gamut is of some consequence in film recording. This is because a great deal of saturation is lost in the photographic process. As a result, the pre-distorted image must contain colors

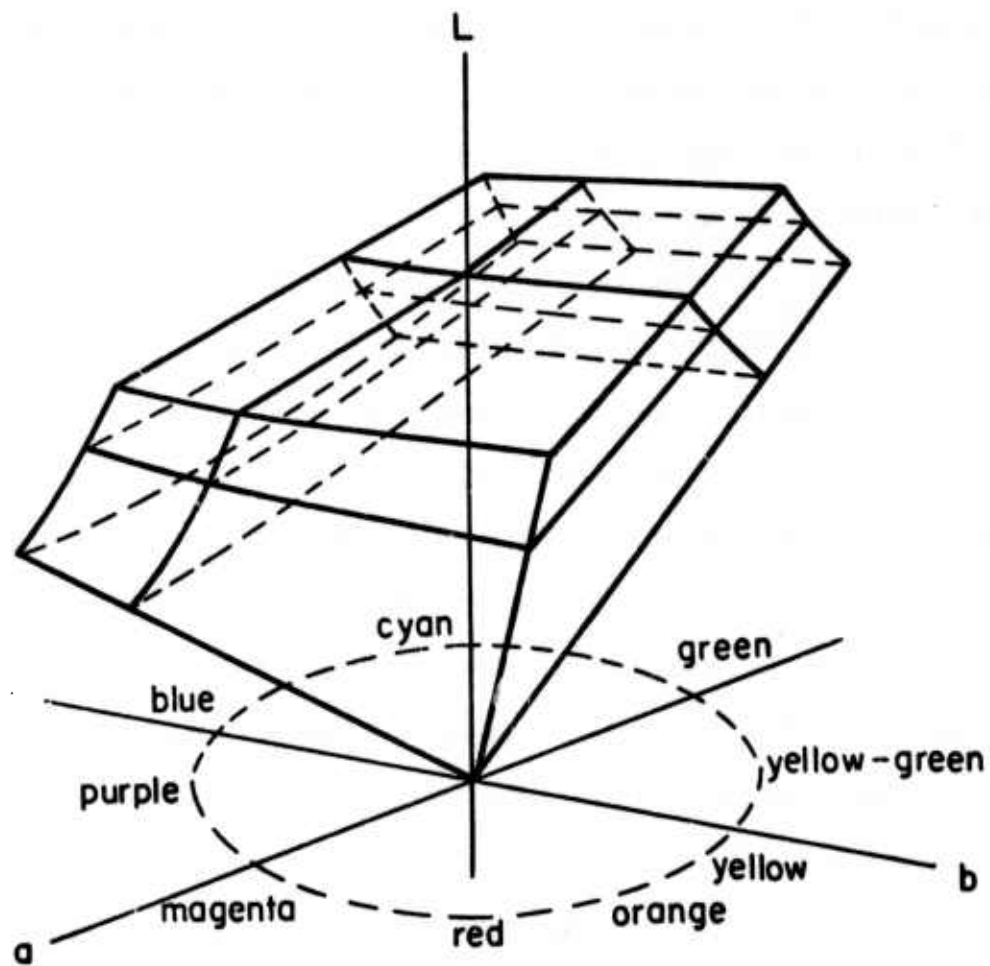


Figure (3.2-3) Gamut of additive system in Lab space

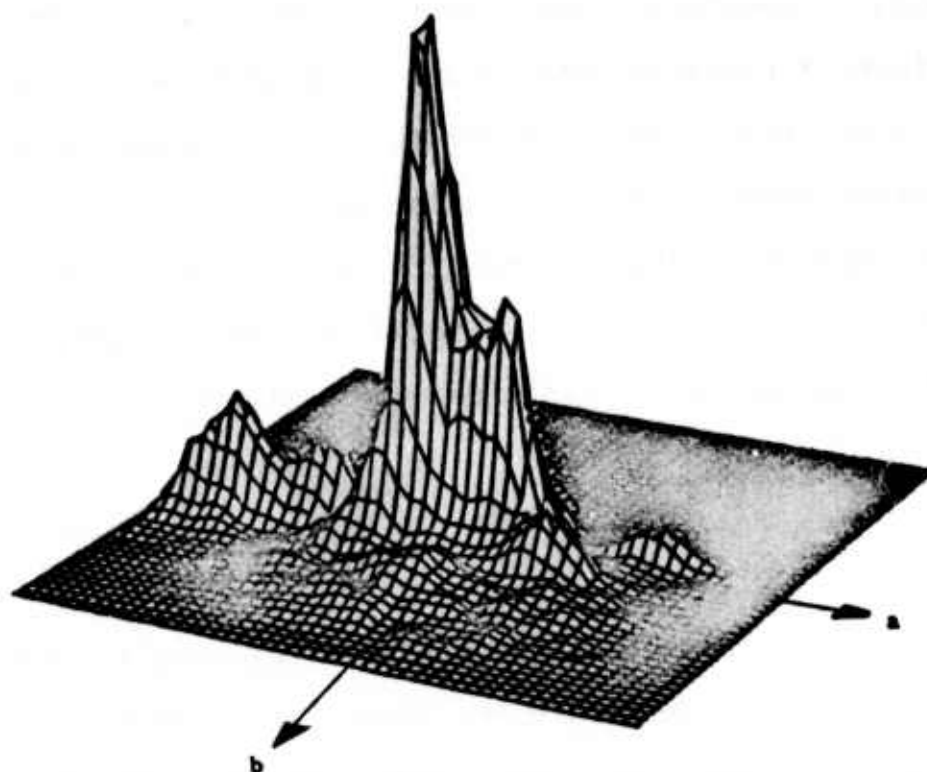


Figure (3.2-4) Histogram of image hues in ab plane
for test image of fig. 6.3-1

3.2

of very high saturation in order to generate the desired results on the film.

Postponing for the moment the question of film recording, consider the question of how to "best" approximate a required color which is beyond the display's gamut (see fig. 3.2-1). Under the naive assumption that the optimal color is located at the point on the surface of the display solid closest to color in question, the problem is easily solved. Stated mathematically, the optimal color vector \underline{T} , minimizes the distance

$$d^2 = \sum_{i=1}^3 (T_i - T'_i)^2 \quad (3.2-4)$$

where \underline{T}' is the color vector beyond the display's gamut, and the \underline{T} vector is within the gamut. For the unit cube color solid of fig. 3.2-1, the constraint on \underline{T} reduces to the stipulation

$$0 < T_i < 1 \quad (3.2-5)$$

for $i=1,2,3$. Because the solid is a cube, the constraints on each element, \underline{T} , are independent of the values taken by the other elements, each term in the summation of 3.2-4

can be minimized without regard to the other terms. Some thought will show that the solution is a zero-one clamp. That is; if T_i' exceeds unity, set $T_i=1$; if T_i' is negative, set $T_i=0$; otherwise, set $T_i=T_i'$. Ironically, this is precisely the fate of a signal "clipped" between a maximum and minimum value (1 and 0 in this case), and is normally what happens to any signal which exceeds the dynamic range of the equipment processing it. If the RGB space, in which the color solid is exactly a cube, were a uniform color space (see Sec. 2.2), the above strategy would indeed be optimal. Unfortunately, this is not the case; the color solid is typically far from rectangular when transformed to a uniform space (see fig. 3.2-3). If the variables in the minimization are not "separable," as is the case with the cubic solid, the problem is substantially more difficult to solve. However, by employing some elementary concepts of Differential Geometry [1], it will be shown that on any continuous surface, a line drawn from the point in question to the closest point on the surface, is always perpendicular to the surface. In order to prove this, it is necessary to introduce the parametric equation of a surface, $\underline{S}(u,v)$. Only two parameters are needed since a surface is

3.2

intrinsically two dimensional. For instance, the parametric representation of a unit sphere could be [1]

$$\underline{S}(u,v) = \underline{e}_1 \sin(u) \cos(v) + \underline{e}_2 \sin(u) \sin(v) + \underline{e}_3 \cos(u) \quad (3.2-6)$$

where the \underline{e}_j are the traditional Cartesian vectors, and u, v are the standard spherical coordinates. The problem is, given a vector \underline{x} , not on the surface, find u, v such that

$$\epsilon = \|\underline{S}(u,v) - \underline{x}\|^2 \quad (3.2-7)$$

is a minimum. Taking partials with respect to u , and v gives

$$(\underline{S}(u,v) - \underline{x}) \cdot \partial \underline{S}(u,v) / \partial u = 0 \quad (3.2-8a)$$

$$(\underline{S}(u,v) - \underline{x}) \cdot \partial \underline{S}(u,v) / \partial v = 0 \quad (3.2-8b)$$

The partials with respect to u and v are vectors tangent to the surface in the directions of increasing u and v (see fig. 3.2-5). Since the line connecting the surface to \underline{x} is perpendicular to both these vectors, it must be perpendicular to the surface. The problem is complicated by the fact that the color solid is usually composed of an

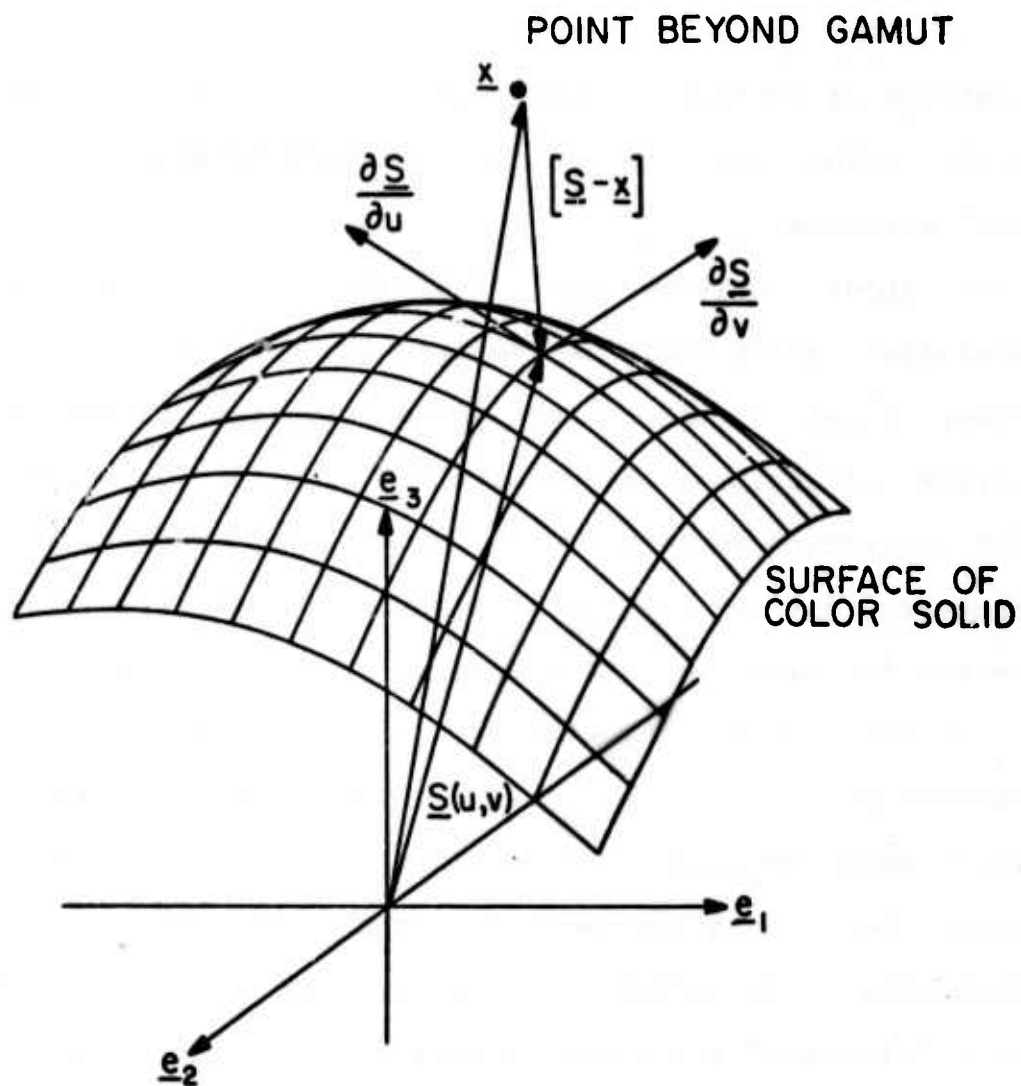


Figure (3.2-5) Geometrical interpretation of the gamut correction problem

3.2

assemblage of surfaces. Hence, it is possible that the closest point may lie on the boundary between two (or three) surfaces.

The above analysis has indicated the type of calculation which would be required in order to optimally correct gamut errors. For gamut problems involving hundreds of colors, the computation required for optimal gamut correction would be quite modest for medium-sized computers, making the methods outlined above practical. However, for high resolution images containing thousands of points, the computation time required becomes prohibitive. Consequently, a sub-optimal, but expedient, method must be used. The zero-one clamp discussed above appears to be an attractive (and certainly fast) possibility. In practice however, the results of such signal "clipping" are rather objectionable. For example, a section of an image requiring an amount of red greater than unity would undergo a clipping of the red signal over the entire section. That is, the clamped red signal would become a constant. The visual effect would be a noticeable loss of detail, since all edge and textural information in the red signal would be destroyed. Thus, even though the zero-one clamp might perform well on a

point-by-point basis, when its effects on the textural details are considered, it becomes a less attractive alternative for natural images (as opposed to bar-charts and the like). An alternative sub-optimal approach, is a contraction of the color space required by the image in order insure that it is contained within the allowable gamut. One method to accomplish this is to decrease the contrast of the image (i.e. multiply each tristimulus value by the same constant which is less than unity). From the standpoint of point-by-point mean square error, this approach performs poorly as compared with the zero-one-clamp. Nevertheless, for small losses of contrast, the subjective deterioration of the total image is not discernable. The reason for the visual system's tolerance of contrast errors is perhaps evolutionary. The most common variations in natural scenes are caused by fluctuations in illumination. Since image formation is a multiplicative phenomenon, these fluctuations are essentially contrast variations. Because the contrast reduction technique gave the best performance (considering total image quality), of the fast methods, it was selected for use with large images where computational speed is of paramount importance.

SUBTRACTIVE SYSTEMS

The majority of color reproduction systems are based on the subtraction of spectral energy from white illumination. In order to quantitatively control a subtractive system such as transparency film, the physics and optics of the situation must be understood. The first section of this chapter discusses the mechanisms by which dye mixtures modulate spectral distributions generating colors. The second section analyzes subtractive systems from a rigorous, mathematical standpoint.

SUBTRACTIVE COLOR SYSTEMS

The distinction between additive and subtractive color reproduction is frequently one which causes a great deal of confusion. Although additive systems were the first to be understood in terms of the three color theory of human vision, most practical color reproduction systems today are of the subtractive variety. Historically, it was the misguided persistence of photographic scientists to

REFERENCES

1. Lipschutz, M.M. Differential Geometry, (Schaum's outline), New York: McGraw-Hill. 1969

develop additive color photography which proved to be the greatest impediment to the development of color film. Some forms of additive color photography have enjoyed commercial success, but the complexities of these methods have made them impractical, and they have been relegated to the status of historical curiosities. the best known of these is Maxwell's method, which he demonstrated in his famous lecture of 1861. Maxwell made three separate positive transparencies of the same scene through red, green, and blue filters. Then, using three projectors equipped with red, green and blue filters, the projected images of the three slides were overlayed in exact register. The quality of the resulting reproduction, and the wide variety of colors produced were compelling evidence in favor of the three color theory of color vision.

A second approach known as the mosaic method, is identical in principle, but avoids the three projectors required in the Maxwell technique by use of a microscopic grid, or mosaic of red, green and blue filters, (the same arrangement as found in the shadow mask television tube). The photograph is taken and viewed through the same screen, which is, in effect, a multitude of tiny Maxwell

4.1

projectors. When viewed at a sufficient distance, the eye averages the tiny triads, seeing them as an additive superposition.

Modern photography owes its relative simplicity to the subtractive technique which exploits the absorptive properties of pigments. While an additive system superimposes the spectra of its primaries, a subtractive system generates colors by subtracting or absorbing certain spectral bands from a single white light source. For instance, when white light is transmitted through yellow dye, the blue portion of the white light is blocked, leaving red and green. Thus, the color yellow can be generated by either adding together red and green, or by subtracting blue from white. Consequently, yellow and blue are considered to be complementary colors. (if added together, they give white). Similarly, the absence of green is magenta (red-blue), and a deficiency of red yields cyan (green-blue). The colors cyan, magenta and yellow are often referred to as the subtractive primaries <*> (see fig. 4.1-1). These dyes can be thought of as "notch filters" which block the red, green and blue portions of the spectrum respectively. A sandwich

<*> as opposed to red, green, and blue, which are considered to be the additive primaries.

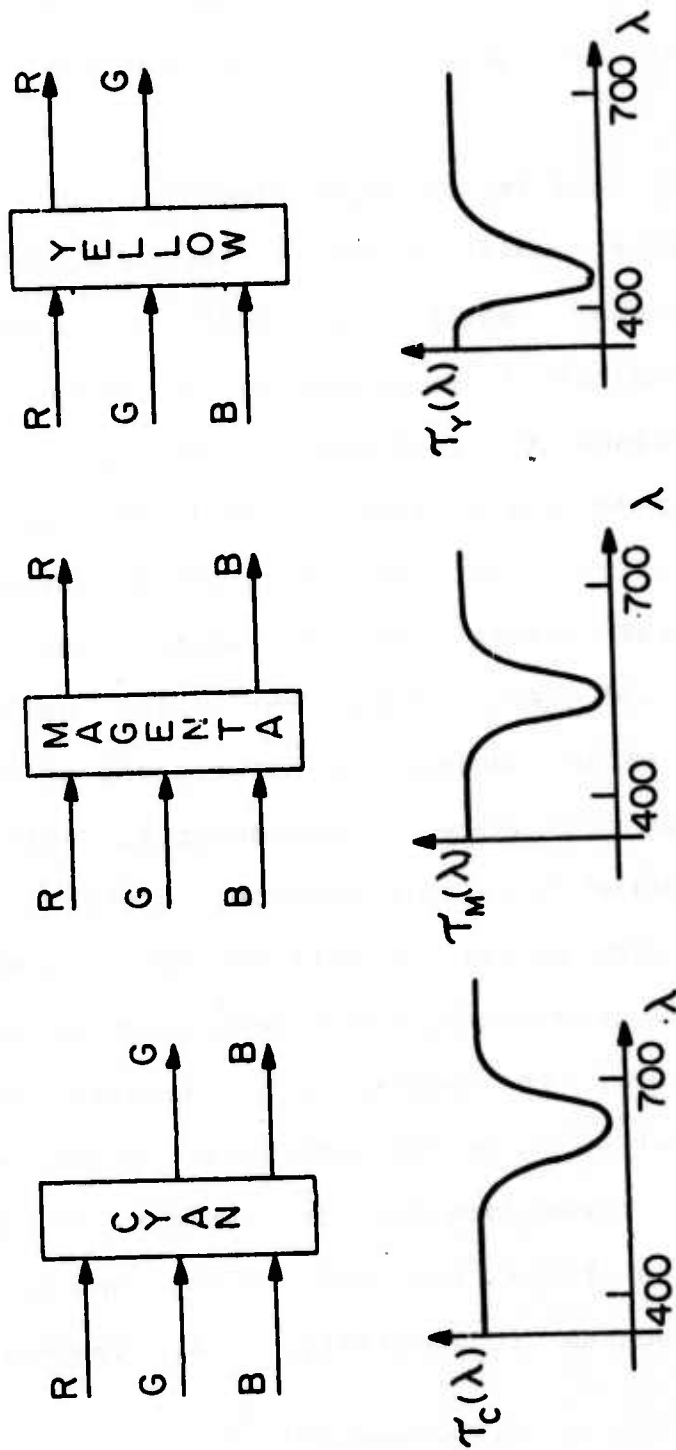


Figure (4.1-1) The subtractive primaries

composed of cyan and magenta filters would thus "notch out" red and green from white, leaving only blue. By mixing various proportions of these dyes, a wide variety of colors can be generated. It should be noted that there is nothing arbitrary about the choice of cyan, magenta and yellow for the subtractive primaries. A subtractive system using red, green and blue primaries would perform quite poorly. This is because red, green and blue dyes are essentially "band pass" filters (as opposed to notch filters) and thus permit only narrow bands of the spectrum to pass through them. Since these bands barely overlap, sandwiching such filters would be equivalent to cascading bandpass filters, and blocking nearly all wavelengths. As a consequence, only very dark colors could be generated by mixing red, green and blue dyes. A quantitative description of subtractive techniques is substantially more involved than that for an additive approach, and will be treated in the next section.

SUBTRACTIVE SYSTEMS; QUANTITATIVE DESCRIPTION

Because of the nonlinearities inherent in subtractive systems, a quantitative description is substantially more involved than is the case for additive systems. As indicated in the last section, the basis of color reproduction with dyes is absorption. As an idealized mathematical model of absorptive material, consider a homogeneous medium, composed of infinitesimal lamina as depicted in fig. 4.2-1. Assume that each of the lamina, of thickness dx , absorbs a fixed proportion of the incident radiant flux, $I(x)$. That is,

$$\frac{dI(x)}{I(x)} = -D dx \quad (4.2-1)$$

where D is a constant determined by the opacity of the material. The solution of the differential equation in 4.2-1 is,

$$I(x) = I_0 \exp(-Dx) \quad (4.2-2)$$

where I_0 is the input radiant flux, and x is the distance penetrated into the medium. This means that the radiant

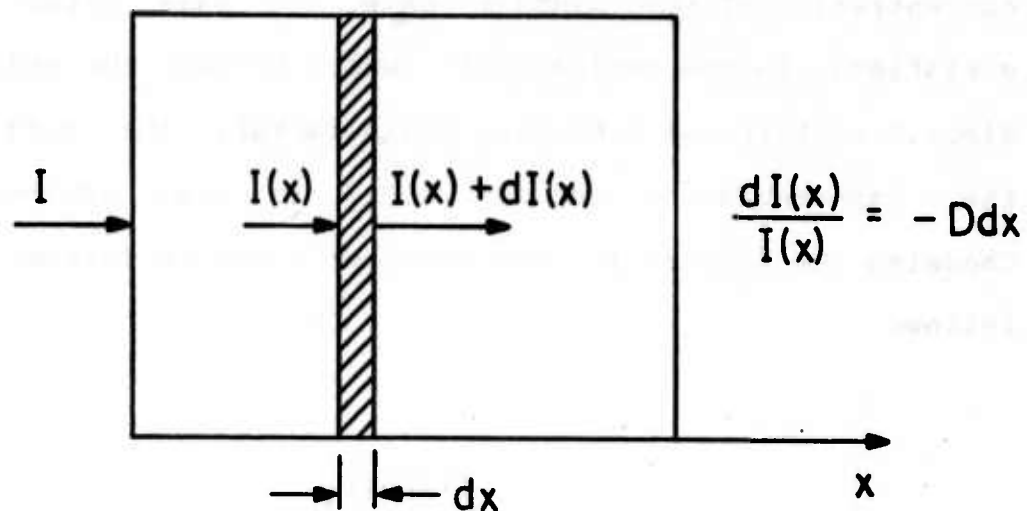


Figure (4.2-1) Mathematical model for absorption

4.2

flux, $I(x)$, decays exponentially as it propagates through the absorptive medium. Furthermore, variations in the concentration of the pigment have the same effect as variations in the optical path length through the medium. Since D contains an arbitrary scale factor, the base of the exponentiation in eqn. 4.2-2 is also arbitrary. Choosing the base as 10, the constant D can be defined as follows

$$D = -\log \left(\frac{I(\text{out})}{I(\text{in})} \right) \quad (4.2-3)$$

where $I(\text{in})$ is the input flux density, and $I(\text{out})$ is the output flux density. This is known as the optical density. The ratio of output to input intensities is defined to be the optical transmissivity, τ [1]. Therefore, the relationships between density and transmissivity are

$$D = -\log(\tau) \quad (4.2-4)$$

and

$$\tau = 10^{-D} \quad (4.2-5)$$

In the actual measurement of density, a degree of variance is present, in that a material's ability to block light is affected by the optical geometry used in the measurement. Reflections, scattering, etc., all affect the apparent transmissivity of a material. Therefore, specifications of density should always include the presumed optical geometry used in the measurement [2, p. 838]. Among the typical geometries are specular, diffuse, and double diffuse.

Up to this point, it has been tacitly assumed that the opaque material was spectrally non-selective, i.e., that its absorptivity was not a function of wavelength. In order to extend the definitions of transmissivity and density to colorants (which are by definition spectrally selective), some further conventions must be established. An obvious extension is to make the density and transmissivity wavelength functions. Thus, let $\tau(\lambda)$ be the attenuation imposed on monochromatic illumination of wavelength λ , and

$$D(\lambda) = -\log[\tau(\lambda)] \quad (4.2-6)$$

where $D(\lambda)$ is defined to be the density as a function of wavelength. In photographic literature, $D(\lambda)$ is termed the spectral density distribution. If more than one colorant is present, the optical density of the mixture becomes

$$D(\lambda) = \sum_{j=1}^3 D_j(\lambda) \quad (4.2-7)$$

where the $D_j(\lambda)$ are the optical densities of the individual colorants. The above relationship is known as the Lambert-Beer law, and is strictly correct only for the transmission of light through solutions. However, since dyes suspended in photographic gelatin behave as if they were dissolved in a solution, the Lambert-Beer law is an excellent approximation for the case of photographic materials viewed by transmitted light [3]. The transmissivity of a dye mixture is obtained by combining eqns. 4.2-5 and 4.2-7,

$$T(\lambda) = 10^{-\sum_j D_j(\lambda)} \quad (4.2-8)$$

The transmissivity of a mixture is thus the product of the

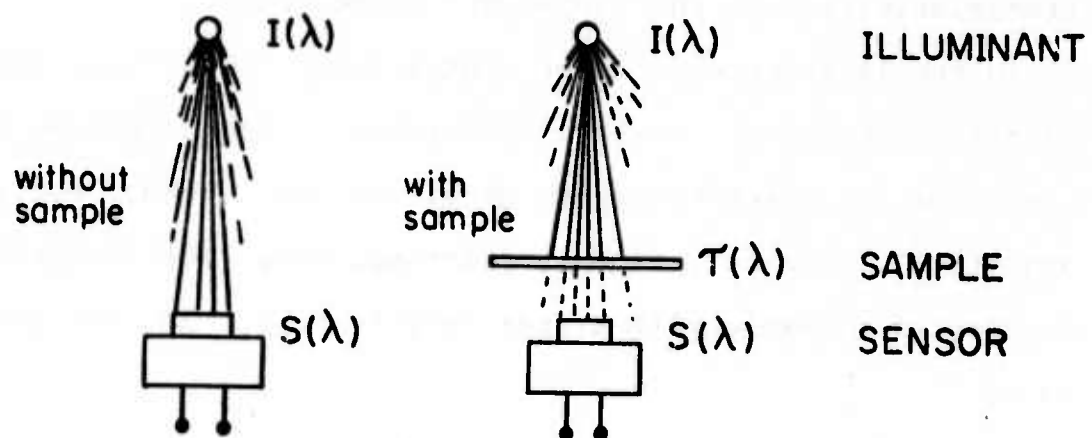
transmissivities of the components taken singly.

In the literature of color photography, densities are usually considered scalar quantities. Consequently, a number of different types of densities are used, e.g., analytical, integral, visual, printing, etc. The integral density of a sample with transmissivity $\tau(\lambda)$, is defined to be

$$D = -\log \left(\frac{\int s(\lambda) I(\lambda) \tau(\lambda) d\lambda}{\int s(\lambda) I(\lambda) d\lambda} \right) \quad (4.2-9)$$

where $I(\lambda)$ is the viewing illuminant, and $s(\lambda)$ is the spectral sensitivity of the sensor used (see fig. 4.2-2). Clearly, this type of density depends on both $I(\lambda)$, and $s(\lambda)$. Consequently, the same sample could be found to have many different densities, depending on the illumination and the sensor used for the measurement.

A more useful (but difficult to determine) density, is the actual concentration of dye (or dyes) present. In order to define such a density, it is first necessary to explain some conventions. As indicated in the derivation of eqn. 4.2-3, the arbitrary scale factor in the



$$\int_{\lambda} S(\lambda) I(\lambda) d\lambda \quad \int_{\lambda} S(\lambda) I(\lambda) T(\lambda) d\lambda \quad \text{OUTPUT}$$

$$\text{TRANSMISSIVITY} \equiv \frac{\int_{\lambda} S(\lambda) I(\lambda) T(\lambda) d\lambda}{\int_{\lambda} S(\lambda) I(\lambda) d\lambda}$$

Figure (4.2-2) Measurement of transmissivity

definition of density was chosen so that density would be equivalent to the negative, base 10, logarithm of transmissivity. For instance, a unit density corresponds to a transmissivity of 0.1. In the case of color dyes, it is not immediately clear how to extend this normalization of monochrome densities. What is typically done in photography, is to define a mixture of cyan, magenta, and yellow dyes which appears colorimetrically neutral, and transmits 10% of the flux incident upon it, as possessing unit dye densities. In other words, the spectral dye densities of the three dyes, $D_j(\lambda)$, are scaled so that

$$\frac{\int 10^{-\sum_j D_j(\lambda)} I(\lambda) t_k(\lambda) d\lambda}{\int I(\lambda) t_2(\lambda) d\lambda} = 0.1 \quad (4.2-10)$$

where $I(\lambda)$ is the viewing illuminant, and the $t_k(\lambda)$ are the XYZ color matching functions, with $t_2(\lambda) = y(\lambda)$, the luminous efficiency function. Note that this convention is consistent with the definition of density used for spectrally non-selective materials. That is, it reduces to eqn. 4.2-5 (for $\tau=0.1$), if $D_j(\lambda)$ and $t_k(\lambda)$ are not functions of wavelength.

In colorimetric literature, tristimulus values of mixtures viewed by transmission are defined as

$$\frac{\int_{10}^{-} \sum_j d_j D_j(\lambda) I(\lambda) t_i(\lambda) d\lambda}{\int I(\lambda) t_2(\lambda) d\lambda} = T_i \quad (4.2-11)$$

where $t_2(\lambda) = y(\lambda)$ is the luminous efficiency function, and d_j are the dye concentrations. However, it is felt by the author that 4.2-11 is unrealistic, in that unity luminance is only attainable with zero density, and hence, will never be achieved. An alternative definition which takes into account the typical viewing conditions for a projected transparency is

$$\frac{\int_{10}^{-} \sum_j d_{min_j} D_j(\lambda) I(\lambda) t_i(\lambda) d\lambda}{\int_{10}^{-} \sum_j d_{min_j} D_j(\lambda) I(\lambda) t_2(\lambda) d\lambda} = T_i \quad (4.2-12)$$

where d_{min_j} are the minimum densities present in the color image being viewed. This is equivalent to defining the lightest areas of an image to be of unity luminance <*>.

The converse problem is less straightforward; i.e. given the desired tristimulus values, what are the necessary dye amounts (or densities), d_j , required to achieve them? The problem is essentially the inversion of eqn. 4.2-12, and will be discussed in Sec. 6.1. A related problem is the determination of the dye concentrations, d_j , from integral densities measured using three different sensors. This problem is described by the following equation

$$\frac{\int_{10}^{-\sum_j d_j D_j(\lambda)} I(\lambda) s_i(\lambda) d\lambda}{\int I(\lambda) s_i(\lambda) d\lambda} = S_i \quad (4.2-13)$$

where the $s_i(\lambda)$ are the spectral sensitivities of the sensors. For the special case in which the sensors are extremely narrow band, and can be approximated by Dirac delta functions, $\delta(\lambda - \lambda_i)$, eqn. 4.2-13 reduces to

$$S_i = 10^{-\left(\sum_j d_j D_j(\lambda_i)\right)} \quad (4.2-14)$$

or

<*> This is felt to be realistic, since projected transparencies are usually viewed in a darkened room, so that the brightest portion of the field of view corresponds to the "thinnest" part of the transparency.

$$\log(S_i) = - \sum_j d_j D_j(\lambda_i) \quad (4.2-15)$$

The above indicates a linear relationship between the log transmissivities (integral densities), and the dye concentrations. It should be emphasized, however, that this relationship is valid only to the extent that the sensor functions, $s_i(\lambda)$, can be approximated by Dirac delta functions. The above linearization is widely used in color densitometry, even though the typical filters used are not sufficiently narrow-band to justify the approximation.

The gamut of colors which such a subtractive system is capable of generating can be calculated by use of eqn. 4.2-12. A set of typical dye densities, $D_j(\lambda)$, used in color photography is shown in fig. 5.2-3. A plot in Lab space of the gamut of colors possible assuming that the dye amounts, d_j , range from 0.5 to 3.0 density units is shown in fig. 4.2-3. The numbers in the plot indicate the lightness (on a scale from 1 to 10), of the third coordinate, L.

The next chapter will discuss the mechanisms which color photography uses to manipulate dye densities.

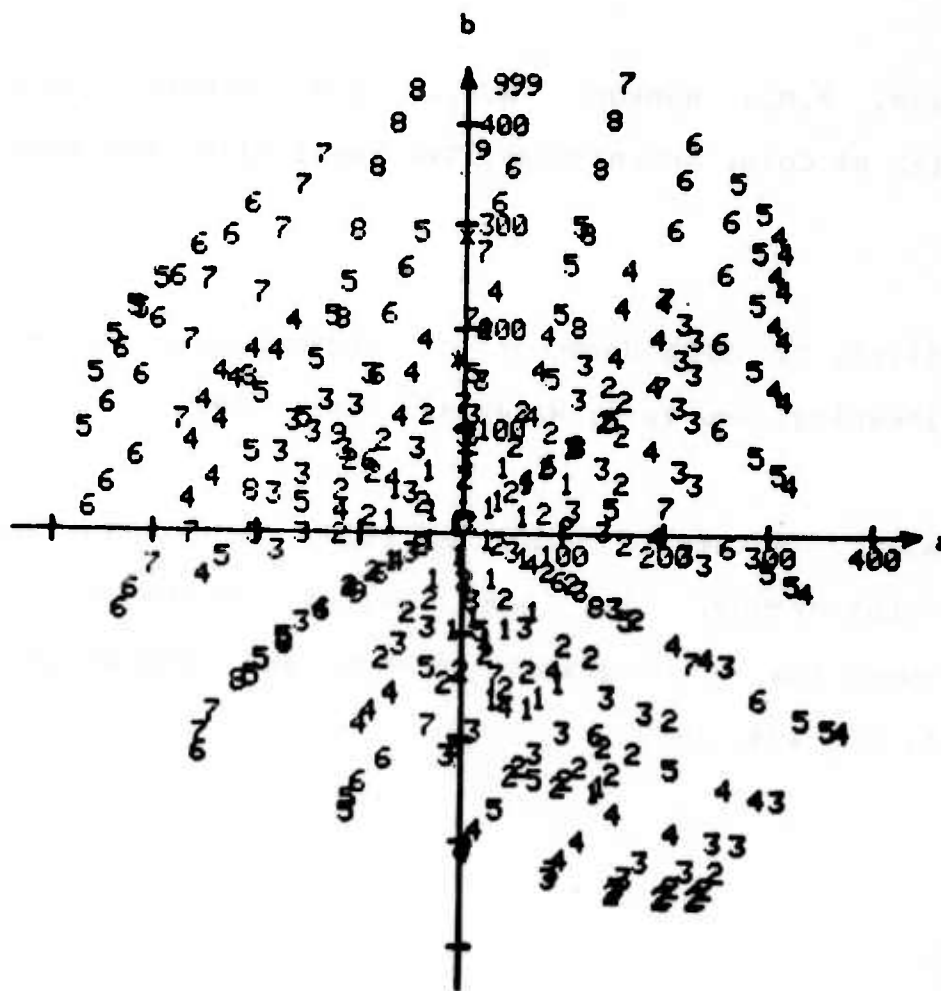


Figure (4.2-3) Gamut of colors in subtractive system

REFERENCES

1. Evans, R.M., Hanson, W.T., and Brewer, W.L. Principles of Color Photography, New York: Wiley and Sons. 1953
2. Woodlief, T. SPSE Handbook of Photographic Science and Engineering, New York: Wiley and Sons. 1973
3. Ohta, N. Metameric Color Matching in Subtractive Color Photography, 1. Dye system obedience to Lambert-Beer law. Photographic Science and Engineering, Vol. 16, p. 136, 1972

COLOR PHOTOGRAPHY

The goal of this chapter is the formulation of a systems model for color photography; i.e. one which can be used to predict the spectral distribution which will result if a given distribution is photographed. The first section describes the mechanics and chemistry involved, the second section, a mathematical model of transparency film, and the third, the experimental verification of the model's predictions.

CHEMISTRY AND PHYSICS OF COLOR PHOTOGRAPHY

The basis of all conventional photography is the light sensitive nature of silver halide compounds, (chlorides, bromides, and iodides). The typical black and white (monochrome) photographic emulsion is prepared by suspending a mixture of silver halides and sensitizing agents in gelatin, coating the resulting suspension on a transparent backing of cellulose acetate, and allowing the material to dry. The properties of the emulsion are

determined by numerous factors, such as the average size of the silver halide grains, and their statistical distribution. When a particular grain is struck by light of sufficient intensity and proper wavelength, changes take place which render the grain more susceptible to chemical reduction by a developing agent. These invisible changes in the emulsion form what is known as the latent image. Upon development, the silver halide grains are transformed to metallic silver in a degree proportional to their exposure to light. Since finely divided silver particles are opaque to light, a negative image results. In order to prevent the unexposed silver halide which remains from eventually darkening the image, the halide grains are removed by treating the emulsion with a fixing agent.

The response of a photographic emulsion to light is usually quantified by using the relationship between the exposure and the resulting density. The exposure can be defined as

$$E = \int_{\lambda} \int_t I(\lambda, t) s(\lambda) dt d\lambda \quad (5.1-1)$$

where $I(\lambda, t)$ is the intensity of the exposing radiation as

5.1

a function of time and wavelength, and $s(\lambda)$ is the spectral sensitivity of the film. Generally, equal exposures result in equal densities, however this rule breaks down for very long exposures, a phenomenon known as reciprocity failure [1]. A typical plot of the relationship between the logarithm (base 10) of the exposure and the resulting optical density (eqn. 4.2-4), is shown in fig. 5.1-1. Such plots are known as Hurter-Driffeld curves, or characteristic curves.

Color transparency film exploits the above principles, using a "sandwich" of three monochrome films, each of which responds to a different region of the visual spectrum (red, green, and blue). The structure used in typical modern color film is shown in fig. 5.1-2. Through a complex series of operations, color film produces colored dyes in response to exposure, in each of its three layers. The red, green, and blue sensitive layers produce cyan, magenta, and yellow dyes in inverse proportion to the degree of exposure received (see Sec. 4.1). This is accomplished as follows:

- (1) Each layer is exposed and developed in the same manner as monochrome film.

- (2) The residual, unexposed silver halides are now

CHARACTERISTIC CURVES

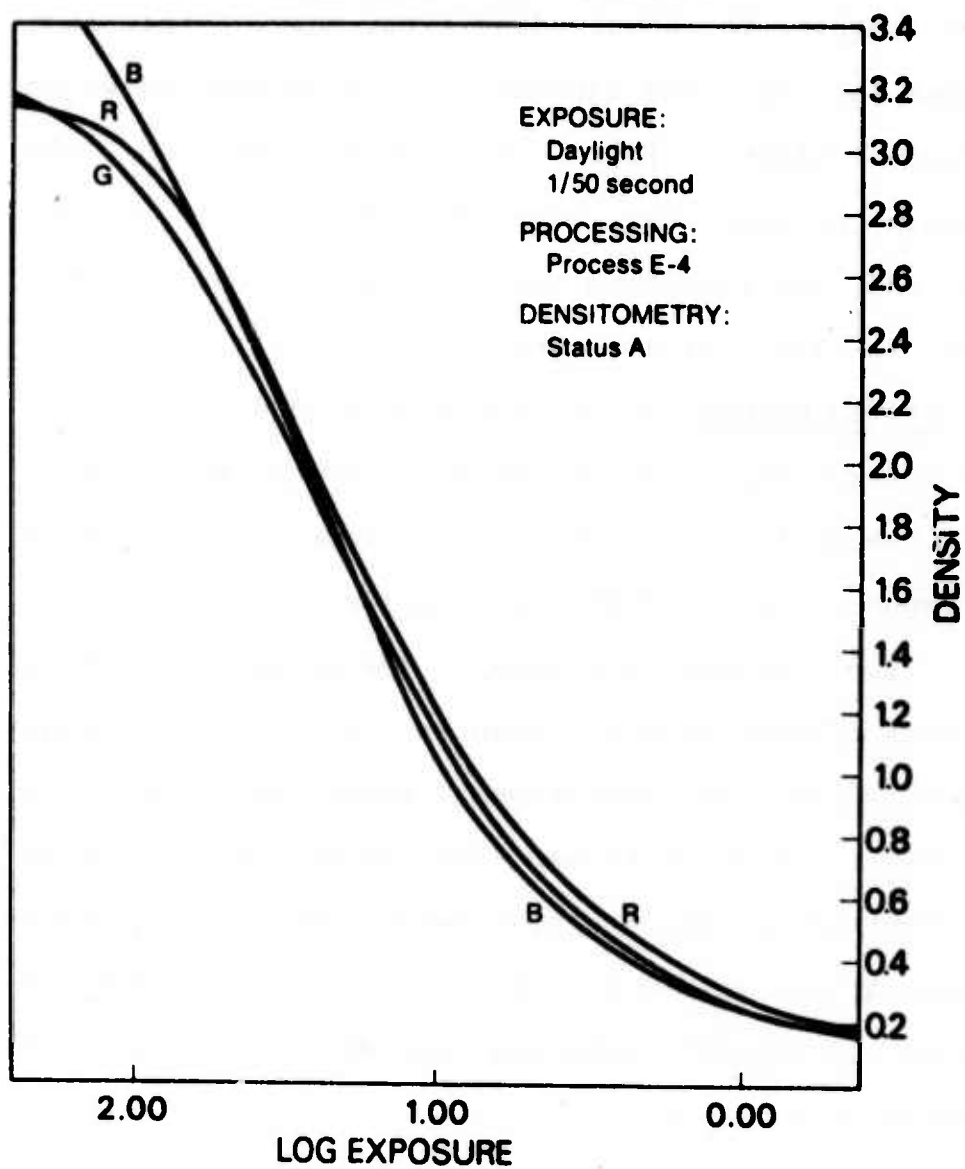


Figure (5.1-1) Hurter-Driffeld curves

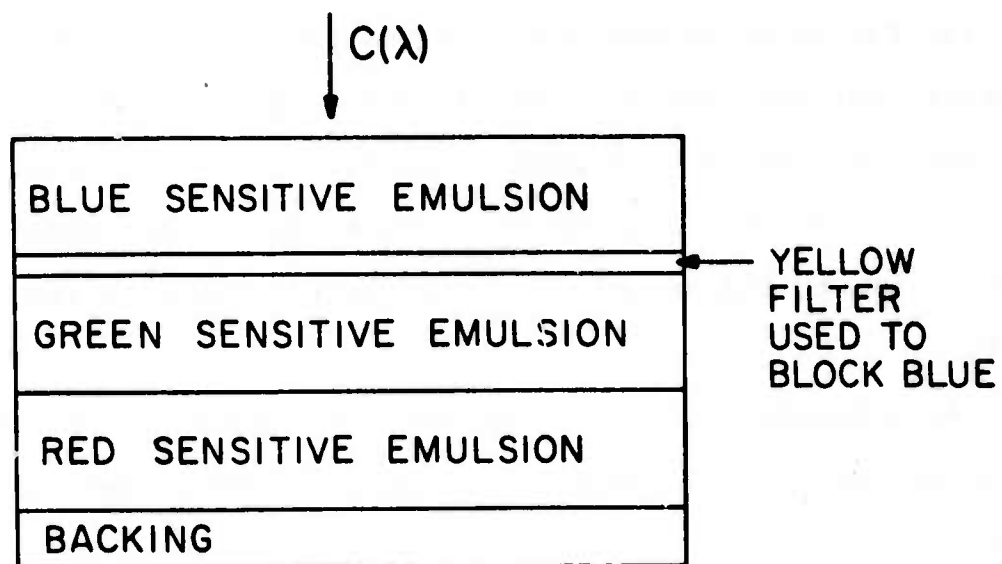


Figure (5.1-2) Typical construction of color film

5.1

deliberately exposed by either a bright light, or chemical means.

(3) The three layers are re-developed, this time the silver halides being reduced are precisely those not affected by the initial exposure. This is done with a special "color developer" which generates dyes in proportion to the amount of silver halide which it reduces <*>.

(4) A bleach is used to extract the metallic silver, leaving an image composed of cyan, magenta, and yellow dyes.

(5) The film is treated with a stabilizer and dried. A schematic representation of the entire process is shown in fig. 5.1-3. A more mathematical treatment of color transparency film is provided in the next section.

<*> The dyes are formed by a chemical reaction between agents in the emulsion known as "couplers," with the oxidation products of the developer.

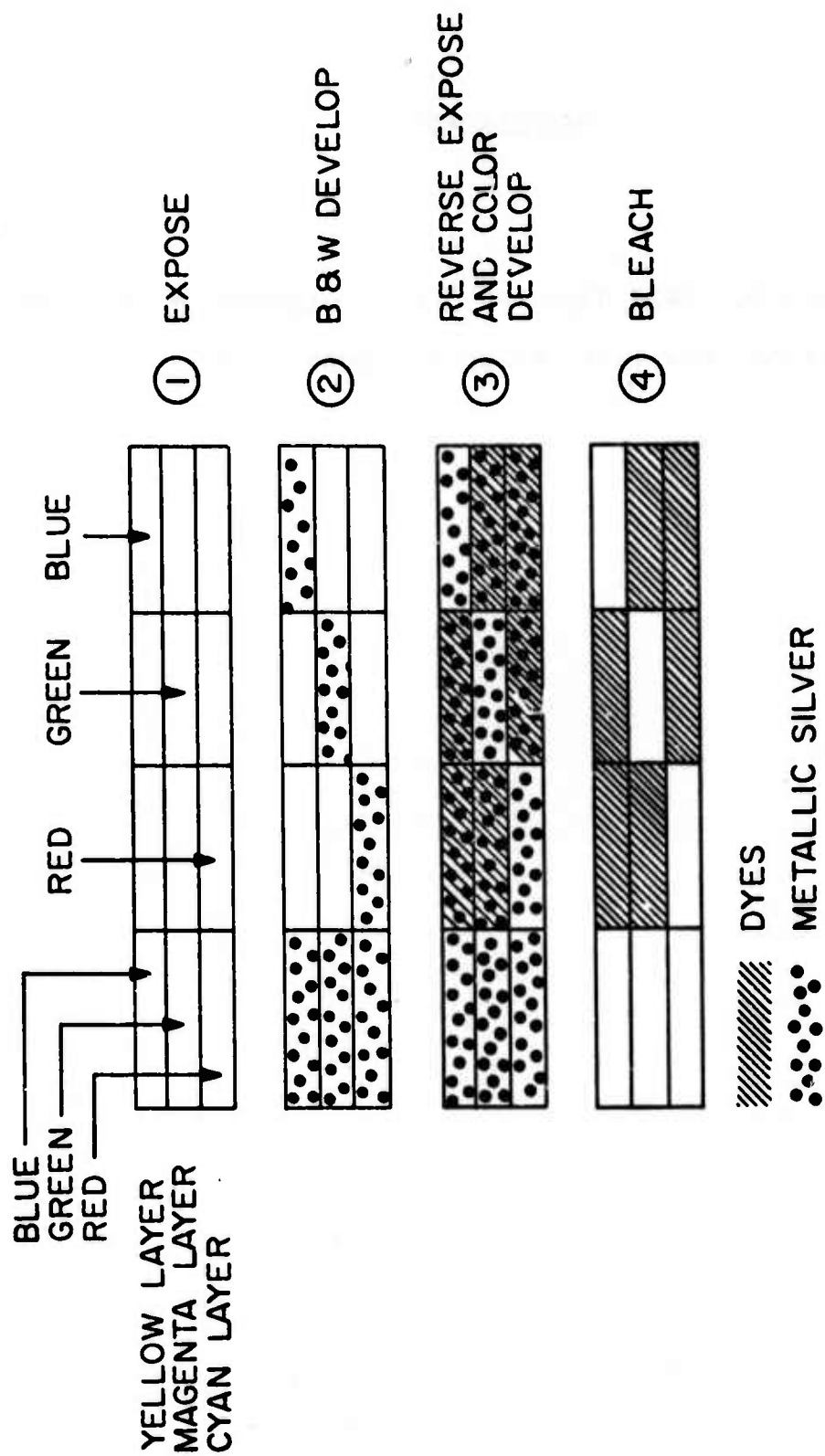


Figure (5.1-3) Schematic representation of modern color film processing

5.1

REFERENCES

1. Woodlief, T. SPSE Handbook of Photographic Science and Engineering, New York: Wiley and Sons. 1973

A MATHEMATICAL MODEL OF TRANSPARENCY REVERSAL FILM

The qualitative description of color film in the previous section can be formulated into a mathematical model, providing predictions in excellent agreement with experiment. The structure of the model shown in fig. 5.2-1 will be discussed component by component. The input spectral distribution, $C(\lambda)$, can be arbitrary, but for the purposes of this study will be assumed to be of the form

$$C(\lambda) = \sum_{i=1}^3 P_i p_i(\lambda) \quad (5.2-1)$$

That is, the input color being photographed is derived from an additive system with primaries $p_i(\lambda)$, and weights P_i (for example a color television monitor). The red, green, and blue sensitive emulsions are presumed to be simultaneously and independently exposed to the input distribution $C(\lambda)$, generating exposures given by <*>

$$x_j = \int L_j(\lambda) C(\lambda) d\lambda \quad (5.2-2)$$

<*> The exposure time is assumed constant, and has been ignored in eqn. 5.2-2.

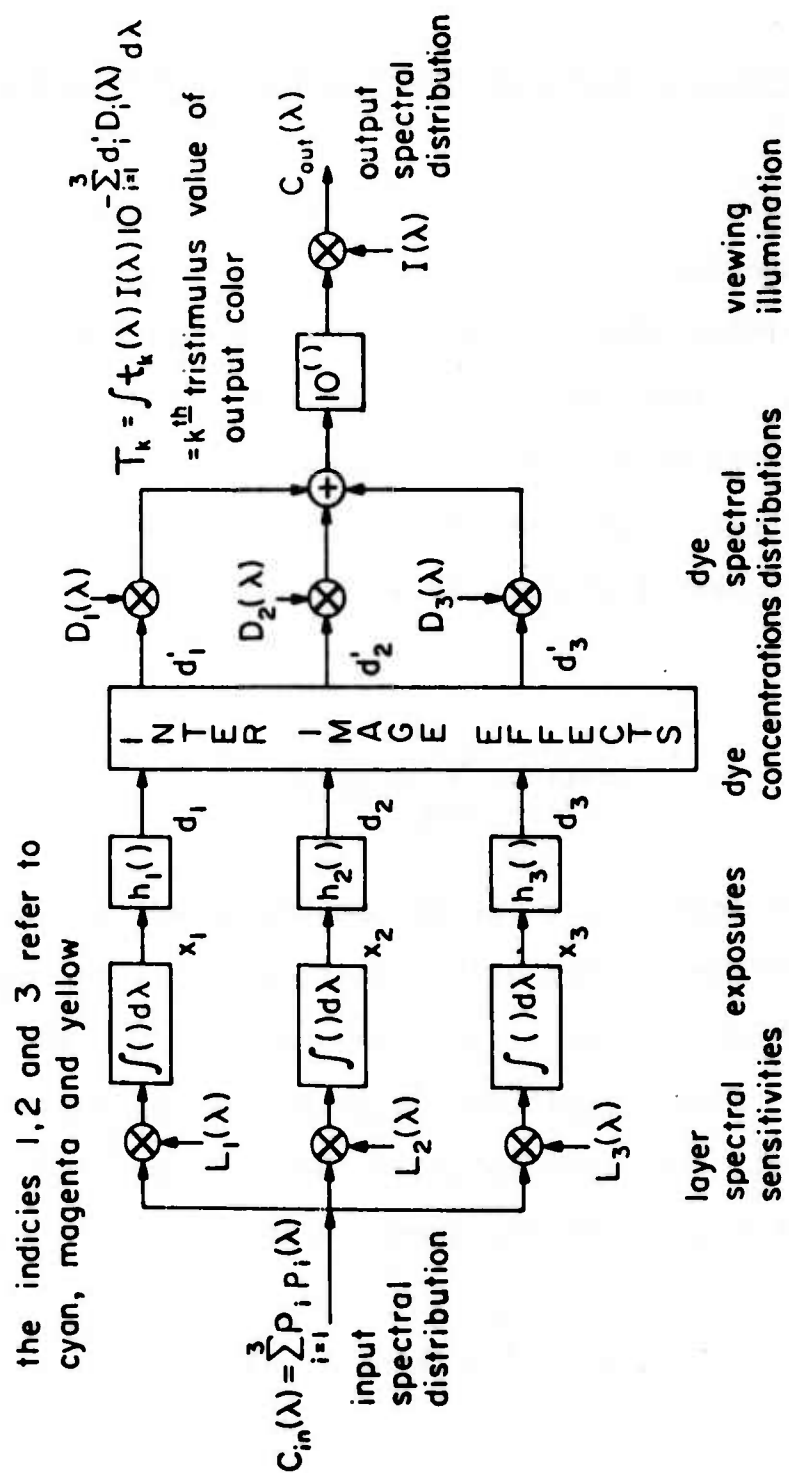


Figure (5.2-1) Mathematical model of color transparency film

where $L_j(\lambda)$ is the spectral sensitivity of the j th layer, and x_j is the resulting exposure. The convention will be to refer to the red sensitive (cyan forming), green sensitive (magenta forming), and blue sensitive (yellow forming) emulsions by the indices 1, 2 and 3. Combining 5.2-1 and 5.2-2 results in a matrix relationship between the vector of exposures, \underline{x} , and the vector of display weights, \underline{p} .

$$\underline{x} = E \underline{p} \quad (5.2-3)$$

where

$$E_{ij} = \int L_i(\lambda) p_j(\lambda) d\lambda \quad (5.2-4)$$

The implicit assumption in eqn. 5.2-2 is that the layer sensitivities, $L_i(\lambda)$, are not functions of the exposing radiation, $C(\lambda)$. This assumption is somewhat simplistic. In reality, the spectral sensitivity of film changes slightly with the intensity of the exposing radiation. In order to make an exact determination of the exposure, the spectral sensitivities, $L_i(\lambda)$, must be known for all

5.2

exposure levels, and the exact solution must be obtained by recursive techniques [1, p. 445]. However, the approximation of eqn. 5.2-2 gives adequate accuracy, and good agreement with experimental results. A plot of the layer sensitivities for EKTACHROME-X film is provided in fig. 5.2-2. The exposures, x_j , give rise to the dye densities, d_j , thru the nonlinearities $h_j(*)$

$$d_j = h_j[\log(x_j)] \quad (5.2-5)$$

The $h_j(*)$ are the Hurter-Driffeld curves for the reversed, or positive process (see fig. 5.1-1). In the case of color film, the densities, d_j , refer not to the density of metallic silver but to the concentration of dye formed in the j th emulsion (Sec. 4.2). The next block of the model in fig. 5.2-1 denoted "interimage effects," represents the interactions between the dye forming layers. Since the three different emulsions are in intimate contact (being in the form of a sandwich), each interferes with its neighbors to some extent. Because all three must chemically compete for the developer, the final dye density in a given layer will be affected by the degree of local chemical activity in the adjacent layers. The usual

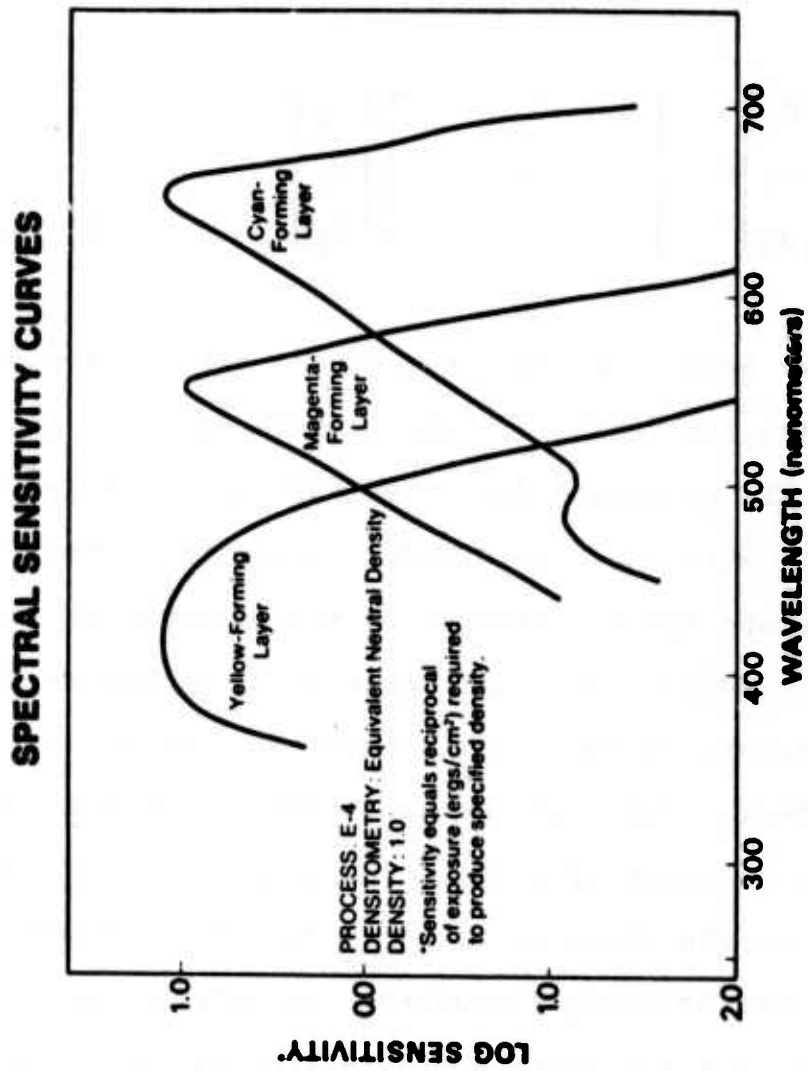


Figure (5.2-2) Layer sensitivities for EKTACHROME-X film

5.2

first order model for interimage effects is a 3x3 matrix multiply [2]

$$\begin{bmatrix} d'_1 \\ d'_2 \\ d'_3 \end{bmatrix} = \begin{bmatrix} & & \\ & G & \\ & & \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ d_3 \end{bmatrix} \quad (5.2-6)$$

Typically, the matrix G has negative off diagonal elements, indicating that if the layers adjacent to a given emulsion are producing dye, then the given layer's dye production will be inhibited. Since the G matrix tends to exaggerate the differences in the elements of the \underline{d} vector, the result of the interimage mechanism is actually an increase in the color saturation. This is a beneficial effect, and can be exploited to enhance the effective gamut of color film [3]. The G -matrix used in the simulation of the model is given in eqn. D-4. As with other formulations assuming linearity, the matrix model of 5.2-6 is only a first order approximation of the actual process.

A second family of phenomena, related to interimage effects, which will be neglected in the film model is known as adjacency effects. This is a term which

describes the interference of photographic image points spatially adjacent to each other (as opposed to interference between layers). Adjacency effects are also considered beneficial, since they result in an apparent sharpening of the image by exaggerating spatial discontinuities in the image.

The spectral density resulting from the dye concentrations, d_j' , is given by the weighted summation

$$D(\lambda) = \sum_{j=1}^3 d_j' D_j(\lambda) \quad (5.2-7)$$

where $D_j(\lambda)$ is the spectral density distribution of the j th dye (see fig. 5.2-3). The transmissivity is therefore

$$\tau(\lambda) = 10^{-\sum_j d_j' D_j(\lambda)} \quad (5.2-8)$$

yielding tristimulus values given by (see eqn. 4.2-12)

$$\frac{\int_{10}^{-\sum_j d_j D_j(\lambda)} I(\lambda) t_1(\lambda) d\lambda}{\int_{10}^{-\sum_j d_{\min_j} D_j(\lambda)} I(\lambda) t_2(\lambda) d\lambda} = T_i \quad (5.2-9)$$

SPECTRAL DYE DENSITY CURVES

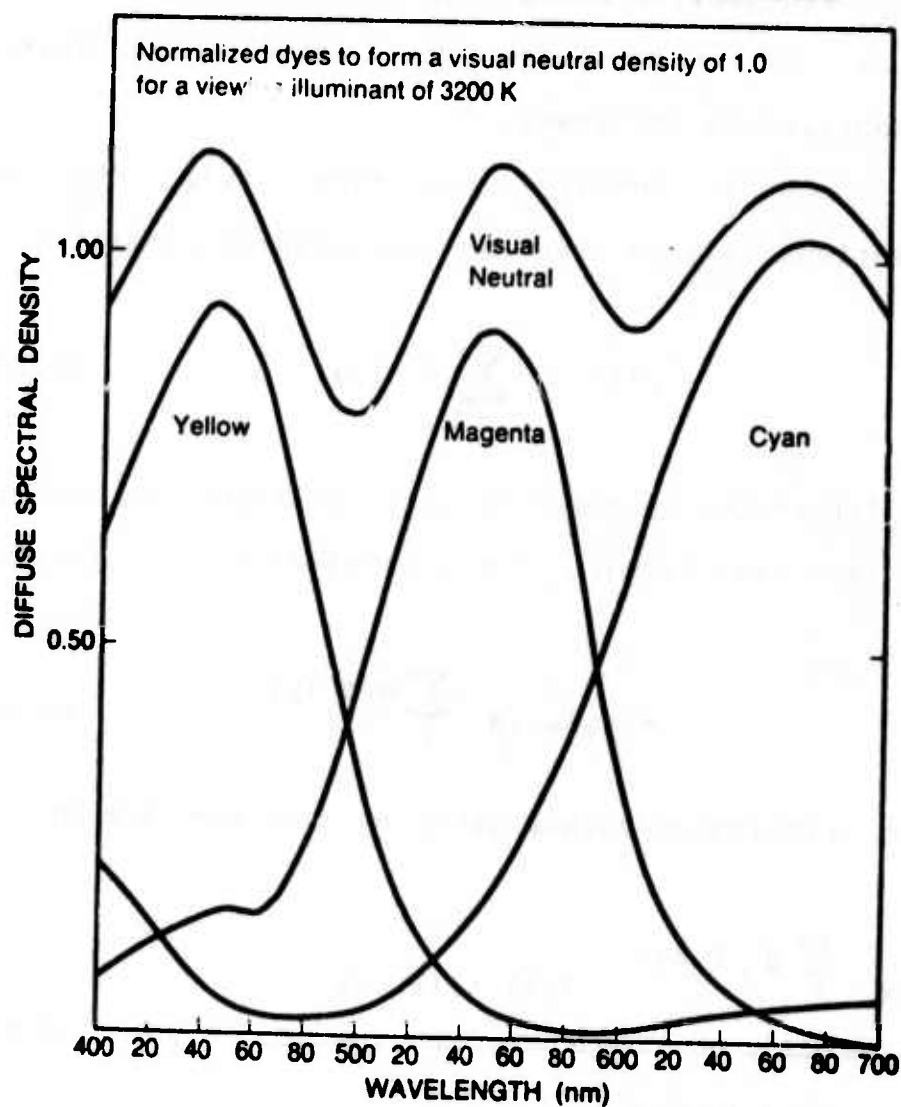


Figure (5.2-3) Dye spectral densities for
EKTACHROME-X film

As before, $I(\lambda)$ denotes the viewing illuminant, and the $t_j(\lambda)$ are the XYZ color matching functions.

In order to assemble a mathematical model for a specific type of film, the necessary spectral sensitivities, $L_j(\lambda)$, characteristic curves, $h_j(*)$, and dye spectral distributions, $D_j(\lambda)$, must be obtained from the manufacturer. A typical set of such specifications for KODAK EKTACHROME-X transparency film is provided in figs. 5.1-1, 5.2-2, and 5.2-3. In the case of interimage effects, however, little information is usually available. Nevertheless, an estimate of the matrix G (eqn. 5.2-6) can be obtained using empirical techniques as follows:

(1) Knowing the photographed input distribution, $C(\lambda)$, and the layer sensitivities, $L_j(\lambda)$, the exposures x_j can be predicted.

(2) With a knowledge of the characteristic curves, $h_j(*)$, the dye concentrations, in the absence of interimage effects can be calculated (the d_j).

(3) Working backwards from the experimentally measured output spectral distribution, $C(\lambda)$, the dye concentrations, d_j' , which gave rise to it can be computed, (see Sec. 6.1).

(4) Using an ensemble of experimentally determined \underline{d}

and $\underline{d'}$ vectors, the best fitting G-matrix can be found using statistical, minimum mean square error techniques (see APPENDIX B).

If the spectral characteristics of the phosphors, and/or the layer sensitivities of the film, are not known, a similar approach can be used to determine the matrix relating the weights on the primaries, P_j , to the resulting film exposures, x_j . In this manner, the necessary parameters for an accurate model can be established. The actual experimental techniques used for the determination of dye densities and tristimulus values will be described in the next section.

REFERENCES

1. Woodlief, T. SPSE Handbook of Photographic Science and Engineering, New York: Wiley and Sons. 1973
2. Hanson, W.T. and Horton, C.A. Subtractive color reproduction: interimage effects, J. Opt. Soc. Am.

5.2

Vol. 42, p.257, 1952

3. Clapper, F.R., Gendron, R.G., and Brownstein, S.A.
Color gamuts of additive and subtractive color
reproduction systems, J. Opt. Soc. Am. Vol. 63, p.625,
1973

A COMPARISON OF EMPIRICAL FINDINGS AND THEORY

The utility of a mathematical model is only as great as its accuracy. Hence, it is of some practical consequence to compare the theoretical predictions of the model described in the previous chapter with the behavior of actual film. In order to do this, an experimental method for the measurement of the dye densities, d_j' , and/or the tristimulus values generated by the film must be available. In this study, most of the the photometric measurements were carried out with the aid of three photoelectric sensors, of composite glass-silicon photocell construction. Their spectral sensitivities were designed to be as close as possible to the XYZ color matching functions $\langle * \rangle$. A plot contrasting the CIE functions $x(\lambda)$, $y(\lambda)$, and $z(\lambda)$, with the sensor characteristics is shown in fig. 5.3-1.

Inasmuch as all sets of color matching functions are linear combinations of one another, the most effective use of the sensors is achieved by estimating the XYZ tristimulus values by linear combinations of the sensor

 $\langle * \rangle$ A set of sensors, or instrument designed to measure tristimulus values is known as a colorimeter.

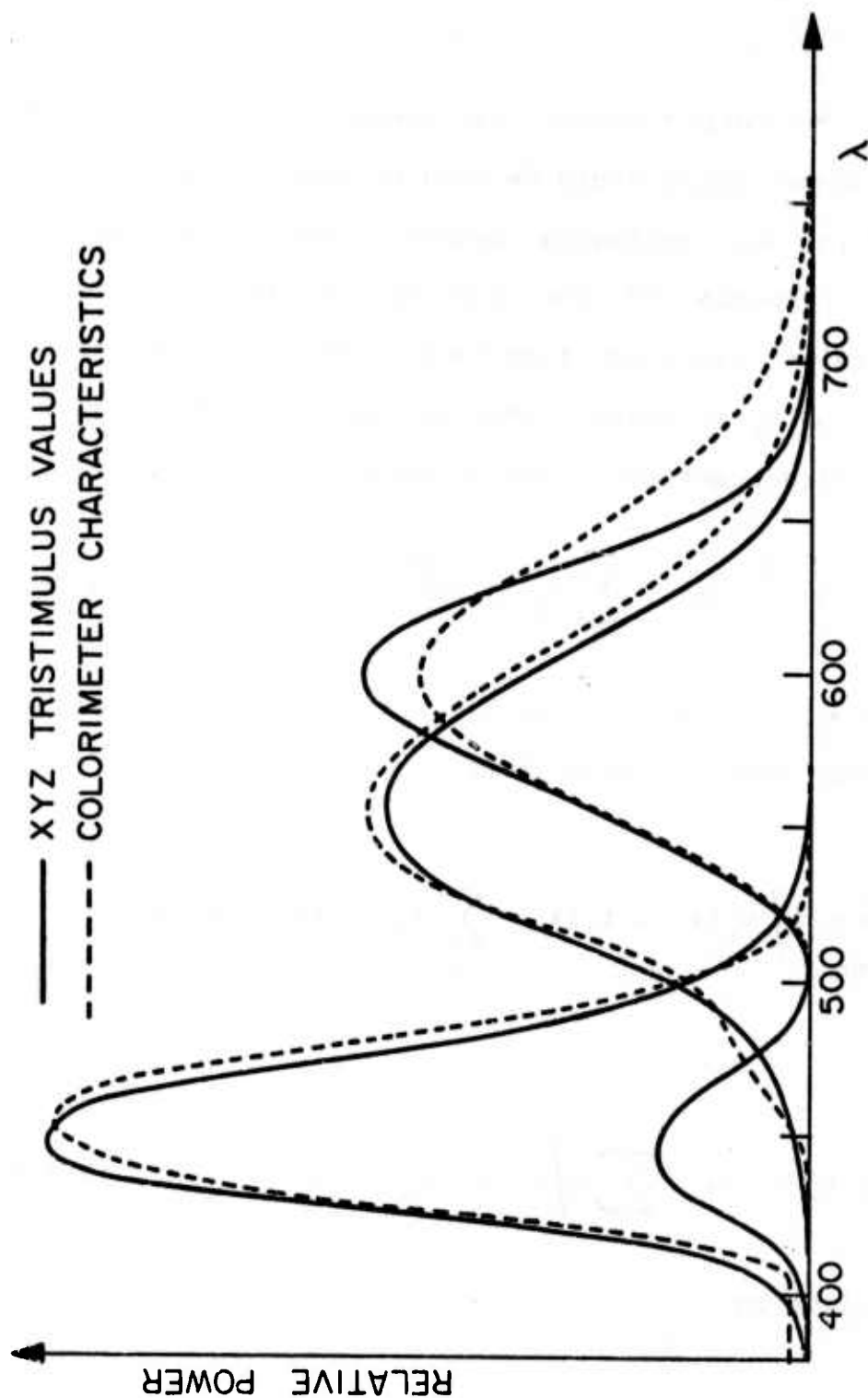


Figure (5.3-1) Comparison of XYZ color matching functions and colorimeter characteristics

5.3

outputs. The "best" linear transformation (in the mean square sense) which could be used to achieve this, can be obtained in the following manner. Let $t'_j(\lambda)$ be the spectral response of the j th sensor, and $t_j(\lambda)$ be the desired color matching functions. The problem is to estimate the $t_j(\lambda)$ using linear combinations of the $t'_j(\lambda)$. The most direct method is the minimization of the integral

$$e_j = \int \left[t_j(\lambda) - \sum_k a_{jk} t'_k(\lambda) \right]^2 d\lambda \quad (5.3-1)$$

over all a_{jk} for $j=1,2,3$ and $k=1,2,3$. Setting the partial derivatives equal to zero gives

$$\frac{\partial e_j}{\partial a_{nm}} = 2 \int t'_n(\lambda) \left[t_j(\lambda) - \sum_k a_{jk} t'_k(\lambda) \right] d\lambda = 0 \quad (5.3-2)$$

or

$$\int t_j(\lambda) t'_n(\lambda) d\lambda = \sum_k a_{jk} \int t'_k(\lambda) t'_n(\lambda) d\lambda \quad (5.3-3)$$

The solution is

$$A = [B'] [B''] \quad (5.3-4)$$

where

$$B_{ij}'' = \int t_i'(\lambda) t_j'(\lambda) d\lambda \quad (5.3-5)$$

and

$$B_{ij}' = \int t_i(\lambda) t_j'(\lambda) d\lambda \quad (5.3-6)$$

A plot contrasting the $\underline{t}(\lambda)$ vectors with their estimates, $A\underline{t}'(\lambda)$, is shown in fig. 5.3-2. In other words, given only the outputs of the three sensors with spectral sensitivities $\underline{t}'(\lambda)$, the best estimate of the actual tristimulus values is given by

$$T_i = \sum_{j=1}^3 a_{ij} \int t_j'(\lambda) C(\lambda) d\lambda \quad (5.3-7)$$

with the a_{ij} defined above. This is a useful technique for significantly improving the precision of an imperfect colorimeter.

A second measurement of great value is the determination of the dye densities in the film sample, d_j' . If the spectral density distributions, $D_j(\lambda)$, are

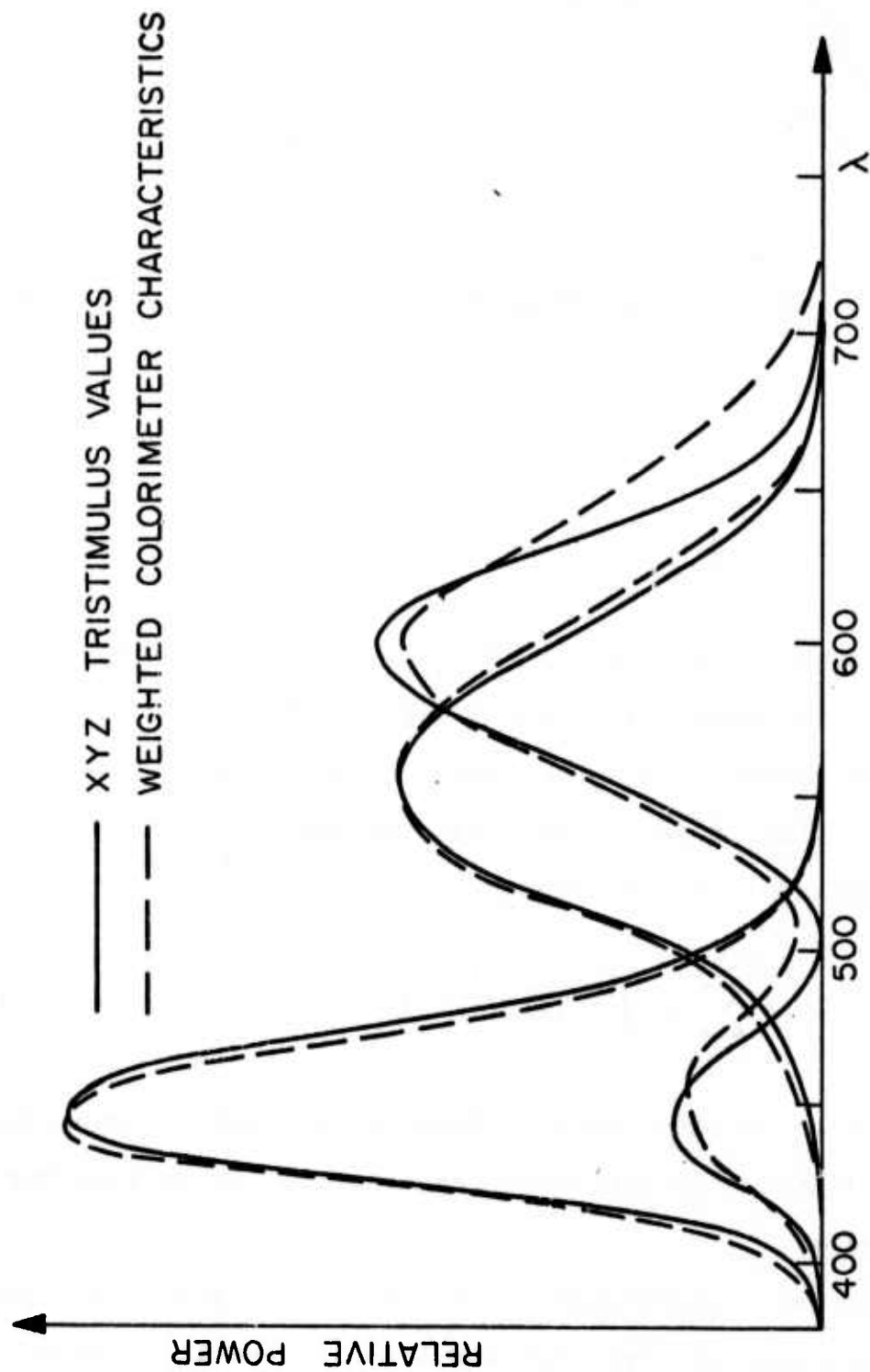


Figure (5.3-2) Actual vs. estimated XYZ color matching curves, using colorimeter

known, both the tristimulus values and the dye densities can be computed in one measurement. From eqn. 4.2-13, the transmittances measured using the colorimeter are

$$\frac{\int_{10}^{-\sum_j d_j D_j(\lambda)} I(\lambda) t_i'(\lambda) d\lambda}{\int I(\lambda) t_i'(\lambda) d\lambda} = \tau_i' \quad (5.3-8)$$

The problem is, given the τ_i' , find the d_1, d_2 , and d_3 which gave rise to them. This problem can be solved by iterative techniques, and will be discussed fully in Sec. 6.1. Assuming for the moment that such techniques are available, it is clearly possible to empirically determine the actual dye densities in a sample of color film. Therefore, using a color monitor with known phosphor characteristics, and a colorimeter as described earlier, the model of Sec. 5.2 can be tested against experiment. This was performed using a CONRAC color monitor, and KODAK EKTACHROME-X film. Pure red, green, and blue fields were generated on the monitor, and photographed over a wide range of exposures. The exposures were varied in increments of one f-stop $\langle * \rangle$, from extreme underexposure to extreme overexposure. The

model was then used in a computer simulation of the same experiment. A comparison of the actual tristimulus values with those predicted by the model is shown in fig. 5.3-3.

The tristimulus values are plotted in the *ab* plane of the Lab coordinate system. The *L* coordinate is represented by the integer values 1 thru 9 (with 0 being black and 10 being white). The nonlinear nature of color film is responsible for the "loops" in the diagram. Note that they emerge and return to the center of the *ab* plane as the exposure is increased. Furthermore, note that the resulting color is a very critical function of the exposure; differences of one-half f-stop cause significant changes in the resultant color. For instance, the red loop proceeds clockwise as the lightness (or equivalently exposure) is increased. This corresponds to a color shift from dark red to light orange, even though the color being photographed remains fixed. The agreement between the model's theoretical predictions and the experimental findings is quite good. Similar comparisons using less saturated colors were made, with equally good results.

In photographic nomenclature, the dynamic range of

 <*> Each f-stop corresponds to a doubling of the exposure. Therefore the exposures were incremented as an exponential ramp in intensity, with each step twice the previous value.

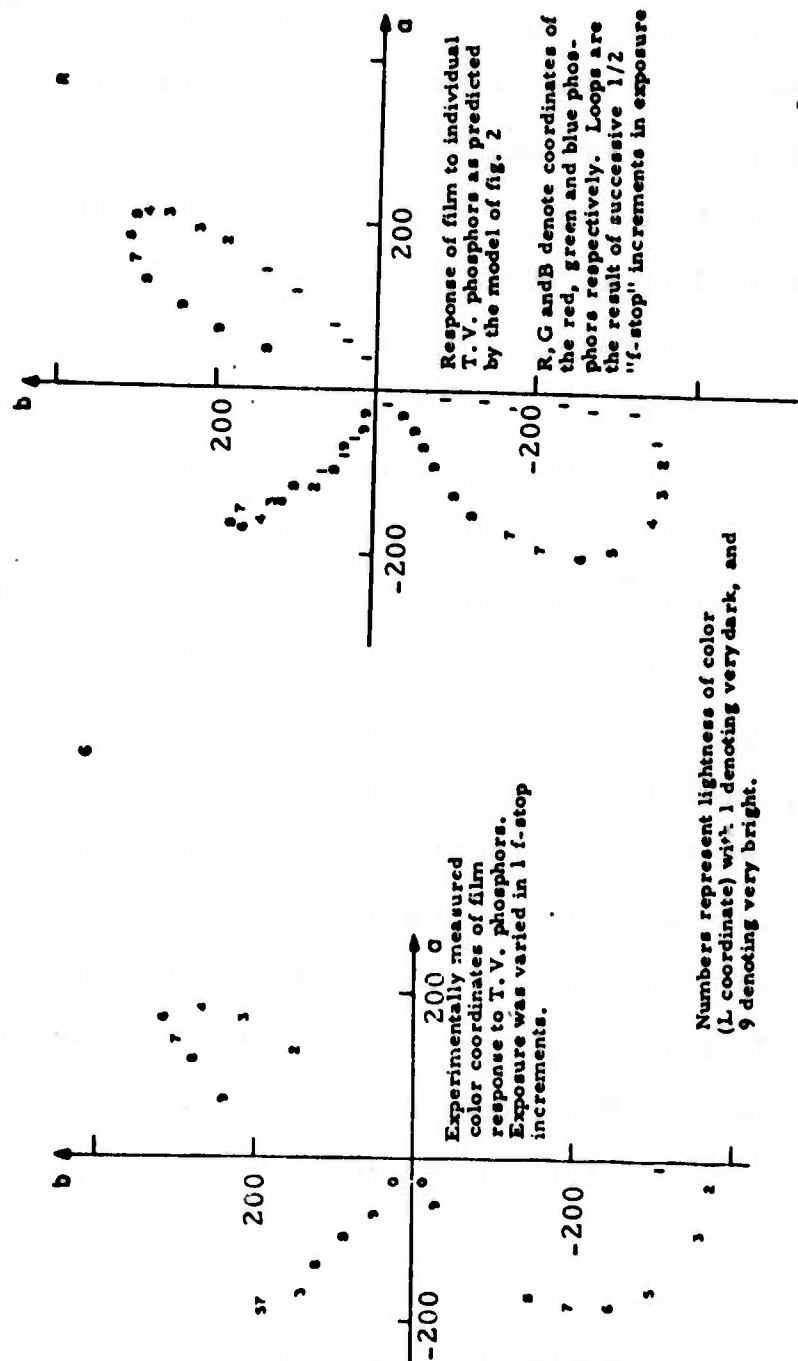


Figure (5.3-3) Comparison of actual vs. predicted tristimulus values, using the film model

acceptable exposure levels for a particular film is known as the film's exposure latitude. The very narrow latitude of most color films is a consequence of a compromise in their design. Specifically, the slope of the Hurter-Driffeld curves (known as the gamma of the film) is deliberately made quite high, typically between 1.5 and 2.0. This is done in order to partially regain some of the color saturation which is inevitably lost in the recording process <*>. The penalty for the added saturation is film's inability to operate over a wide range of exposures (as evidenced in fig. 5.3-3).

Having verified the accuracy of the model, the next step is the inversion of the model. That is, finding the necessary input color (on the monitor), required to produce a desired output color (on the film). This is clearly possible, provided that the requisite output color is one which the monitor is capable of causing the film to produce; i.e. the output color is within the film's gamut, and the necessary input color is within the monitor's gamut. These considerations will be more fully explored in the succeeding sections.

<*> The loss of saturation is predominantly due to the broad, overlapping, absorption bands of the dyes used.

REFERENCES

1. Ralston, A. A First Course in Numerical Analysis, New York: McGraw- Hill. 1965

THE DISPLAY-FILM SYSTEM AS A MAPPING OF TRISTIMULUS VALUES

For the case in which a digitally controlled image is displayed on a monitor and photographed, the recording process can be considered to be a mapping from an input set of tristimulus values to an output set. Therefore, the colorimetric distortions inherent in the recording process can be neutralized by a pre-distortion the image to be photographed, using the inverse of the mapping. Although this presents no conceptual difficulties, the computational difficulties involved in the inversion are substantial. The first section of this chapter describes the individual components of the display-film model, then discusses the mathematical procedures required to invert the model, working directly from the desired output tristimulus vector. This approach will be termed the direct method.

The subsequent section will present two methods which can be used to implement the pre-distortion which do not require the formal inversion of the model. These will be termed indirect methods.

The final section discusses the experimental results, and some practical considerations.

DIRECT INVERSION OF THE FILM MODEL

The motive for the development of an accurate film model was not merely the ability to simulate film, but rather a means towards the control of film, exploiting its predictability. The problem of colorimetrically pre-distorting the image to be photographed, in order to compel the film to yield the desired tristimulus values, is mathematically equivalent to the inversion of the film model. This section will discuss the manner in which this can be accomplished.

Before discussing the inversion process, the "forward model" of fig. 6.1-1 will be described in detail. The actual transformations used for the specific film and display in this study are provided in APPENDIX D.

First, the spectral sensitivities of the display phosphors, and the layer sensitivities of the film must be available. Using eqn. 5.2-4, the 3x3 matrix, E , in the following relationship can be calculated

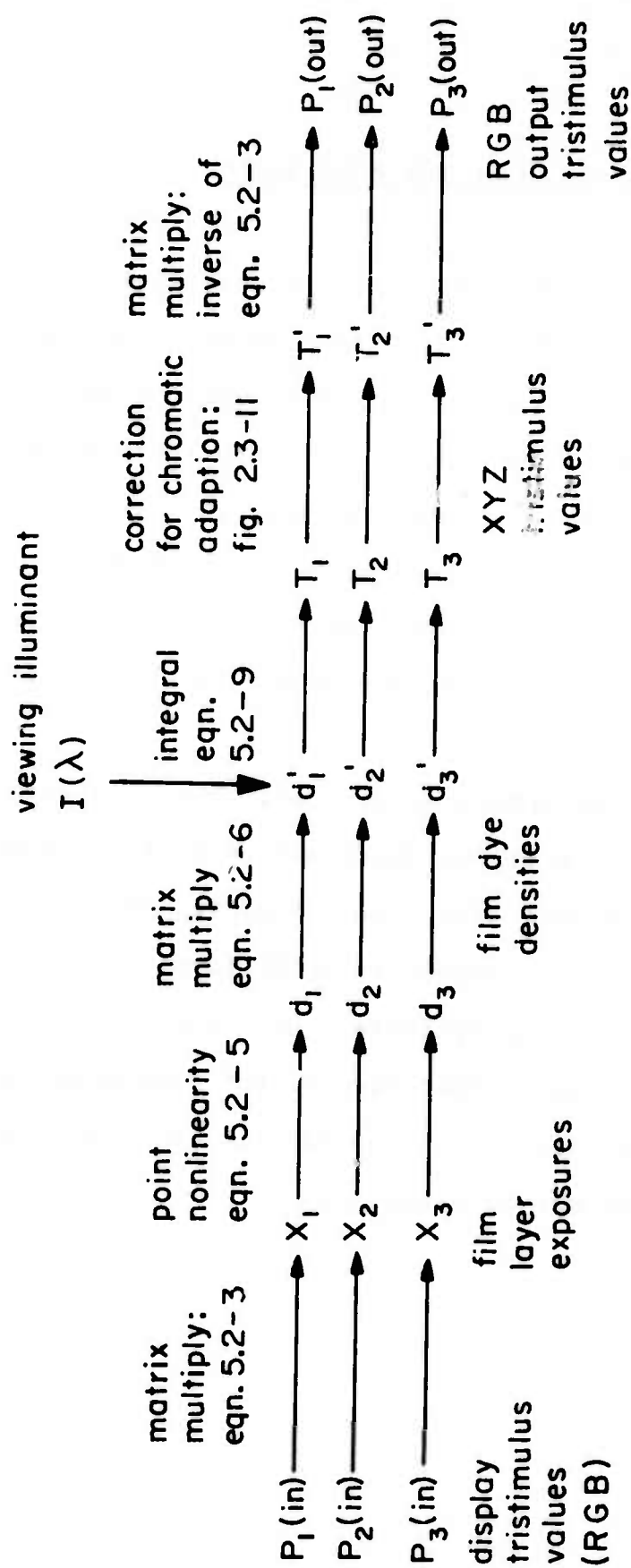


Figure (6.1-1) The display-film cascade as a mapping

6.1

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} & & \\ & E & \\ & & \end{bmatrix} \begin{bmatrix} P_1(\text{in}) \\ P_2(\text{in}) \\ P_3(\text{in}) \end{bmatrix} \quad (6.1-1)$$

where the x_j denote the film layer exposures, and the $P_j(\text{in})$ denote the display phosphor weights. With the aid of the Hurter-Driffeld curves supplied by the film manufacturer, the dye densities (neglecting interimage effects) can be computed

$$d_j = h_j [\log(x_j)] \quad (6.1-2)$$

where the $h_j[\cdot]$ represent the H&D curve nonlinearities (see eqn. 5.2-5 and fig. 5.1-1). The resulting dye densities are transformed using the interimage effect matrix G according to the relation (see eqn. 5.2-6)

$$\begin{bmatrix} d'_1 \\ d'_2 \\ d'_3 \end{bmatrix} = \begin{bmatrix} & & \\ & G & \\ & & \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ d_3 \end{bmatrix} \quad (6.1-3)$$

With the aid of the Lambert-Beer law and eqn. 4.2-12, the tristimulus values (in the absence of adaption

effects) can be computed as

$$T_i = \frac{\int_{10}^{-\sum_j d_j' D_j(\lambda)} I(\lambda) t_i(\lambda) d\lambda}{\int_{10}^{-\sum_j d_{min_j} D_j(\lambda)} I(\lambda) t_2(\lambda) d\lambda} \quad (6.1-4)$$

At this point, if it is assumed that the photograph is viewed under the same state of adaption as the monitor display, the tristimulus values calculated in 6.1-4 can be transformed to RGB space by use of the 3x3 matrix multiply

$$\begin{bmatrix} P_1(\text{out}) \\ P_2(\text{out}) \\ P_3(\text{out}) \end{bmatrix} = \begin{bmatrix} & & \\ & A & \\ & & \end{bmatrix} \begin{bmatrix} T_1 \\ T_2 \\ T_3 \end{bmatrix} \quad (6.1-5)$$

Transformations between sets of tristimulus values, such as the one above, are discussed in APPENDIX A.

If an adaption transform is required, the procedures described in Sec. 2.3 can be used. First, the XYZ tristimulus values resulting from eqn. 6.1-4 must be transformed to a coordinate system determined by the actual receptor sensitivities of the eye.

$$\begin{bmatrix} S_1 \\ S_2 \\ S_3 \end{bmatrix} = \begin{bmatrix} & & \\ & A & \\ & & \end{bmatrix} \begin{bmatrix} T_1 \\ T_2 \\ T_3 \end{bmatrix} \quad (6.1-6)$$

where the matrix A is determined by the assumed fundamental sensitivities (for instance see eqn. 2.3-11). If the simple Von Kries type of scaling is considered adequate, the receptor signals, S_j , defined by eqn. 6.1-6 are scaled so that the reference white yields $S_j = 1$. This is expressed in eqn. 2.3-2. If a more accurate adaption transform is desired, the Von Kries scaling can be supplemented by the subtractive bias procedure illustrated in fig. 2.3-4. Finally, the adaption-corrected XYZ tristimulus values can be obtained from the perceptual g -variables of fig. 2.3-4 by reversing the transformation. A block diagram of the entire adaption transform is shown in fig. 2.3-11. With the output color expressed in the RGB coordinate system (using 6.1-5), the entire process can be viewed as a channel associating an output RGB vector with a given input RGB vector.

The transformations discussed above provide a method for the calculation of an output tristimulus RGB vector, denoted \underline{P} (out), given the input tristimulus vector,

denoted $\underline{p}(\text{in})$. The converse problem (the inversion), is equivalent to the inversion of each of the individual transforms in the cascade. This is a straightforward procedure for all of the above relations, with the exception of eqn. 6.1-4. In order to invert this integral equation, the dye densities, d_j' , necessary to produce a given set of tristimulus values, T_j , must be calculated. The most practical method for the inversion of nonlinear, vector functions such as the above is the Newton-Raphson iteration technique [1]. This method consists of progressively refining an initial guess for the vector \underline{d} $\langle*$ by the following recursion

$$\underline{d}^{k+1} = \underline{d}^k - [\underline{J}]^{-1} \underline{e}^k \quad (6.1-7)$$

Each of the variables is "updated" after each recursion, with the k th estimate denoted by the superscript (k) . The matrix \underline{J} (known as the Jacobian matrix), is computed from the partials

$$J_{ij} = \frac{\partial T_i}{\partial d_j} \quad (6.1-8)$$

 $\langle*$ The prime on the \underline{d}' vector has been omitted for the sake of notational clarity.

or

$$J_{ij} = \frac{-\ln(10) \int_{10}^{-\sum_k d_k D_k(\lambda)} I(\lambda) t_i(\lambda) D_j(\lambda) d\lambda}{\int_{10}^{-\sum_k d_{\min_k} D_k(\lambda)} I(\lambda) t_2(\lambda) d\lambda} \quad (6.1-9)$$

The error vector, \underline{e}^k , is the difference between the current tristimulus vector, \underline{T}^k , (computed from the current \underline{d}^k vector, as in 6.1-1)) and the desired tristimulus values, $\hat{\underline{T}}$.

$$\underline{e}^k = \underline{T}^k - \hat{\underline{T}} \quad (6.1-10)$$

Convergence of the iteration hinges on the accuracy of the initial guess for the \underline{d} vector. That is, if the iteration is started with an estimate of the \underline{d} vector which differs greatly from the "true" value, the procedure will diverge. A simple technique for overcoming this difficulty is as follows:

(1) Treating the relationship between the \underline{d} and \underline{T} vectors as a vector mapping from the three dimensional

6.1

d-space to the three dimensional T-space, compute a collection of T vectors corresponding to a large number of different d vectors.

(2) Considering the resulting ensemble of d and T vectors as a statistical mapping, find the 3x3 matrix, A, which best fits the data (see APPENDIX B).

(3) Having determined the matrix, operate on a desired T vector with the inverse of A to obtain a good estimate of the associated d-vector.

(4) Use this d vector as the initial guess in the recursion of 6.1-2, and recursively obtain the exact T-vector.

The remaining blocks of the model invert directly:

(1) The dye densities before the interimage effects are obtained by the inverse of the matrix G (see eqn. 5.2-6).

(2) The exposures follow from the dye densities via the inverses of the Hurter-Driffeld curves.

(3) The display weights can be obtained from the exposures using the inverse of the E matrix (see eqn. 5.2-3).

The method described above gives accurate results, but necessitates a great deal of computation per image point.

6.1

Each of the recursions demands 12 tristimulus integrations, (9 for the J matrix, and 3 for the e-vector). Typically, at least 5 recursions are needed for convergence. Therefore a minimum of 60 numerical integrations per image point are required. If the number of image points involved is on the order of hundreds, such as would be the case for very low resolution images, the exact inversion is well within the practical limitations of standard computers. If on the other hand, the pre-distortion of medium-to-high resolution images is attempted, the procedure could easily consume thirty minutes of computer time on all but the most powerful machines.

However, there is little to be gained by inverting the model with greater precision than the model itself represents reality. Since the overall speed of the inversion procedure is basically determined by the number of numerical tristimulus integrations performed, one obvious route to greater efficiency is the use of an approximate, but fast, numerical integration procedure. The method to be described below is based on this philosophy.

The conventional techniques used for numerical

integration are Simpson's rule, and the trapezoid rule [1]. These schemes are easy to implement, but lack efficiency. The customary method for fast tristimulus integration, is the selected ordinate method, which was developed before the advent of modern computers [2, p. 317]. The basic idea in this approach is the selection of unequally spaced abscissas. First, the i th tristimulus integral is broken up into $N+1$ contiguous intervals, separated by the points λ_j^i for $j=1,2,3,\dots,N$. Thus

$$T_i = \sum_j \int_{(\lambda_j^i)}^{(\lambda_{j+1}^i)} t_i(\lambda) C(\lambda) d\lambda \quad (6.1-11)$$

For closely spaced λ_j^i 's, the integral can be approximated as

$$T_i \approx \sum_j t_i(\lambda_j^i) C(\lambda_j^i) \Delta\lambda_j^i \quad (6.1-12)$$

If the λ_j^i 's are chosen such that

$$\Delta\lambda_j^i t_i(\lambda_j^i) = k_i \quad (6.1-13)$$

where the k_i 's are constants, then (6.1-13) reduces to

$$T_i \approx k_i \sum_j C(\lambda_j^i) \quad (6.1-14)$$

The obvious advantage of this scheme is the elimination of the multiplications, only one multiplication and n additions are needed. Tables of the λ_j^i 's can be found in various references [2, p. 324]. The penalty for the increase in speed is of course a decrease in accuracy, since the approximation of (6.1-12) is rather crude. However, more sophisticated numerical integration methods are available which can provide great speed with minimal loss of precision. Among the best of these is Gaussian quadrature integration, which will now be described.

As with the selected ordinate technique, the Gaussian method uses unequally spaced abscissas. Consider the numerical integration of the arbitrary function $w(x)f(x)$

$$\int_a^b w(x) f(x) dx \approx \sum_{i=1}^n H_i f(x_i) \quad (6.1-15)$$

where $w(x)$ is some nonnegative weighting function, the x_i are the points at which $f(x)$ is evaluated, and the H_i are weights. In the Gaussian method, both the weights, H_i , and the abscissas, x_i , are chosen in an optimal fashion, to be described shortly. The advantage of this technique

is that it is at least "twice as good" as conventional numerical integration procedures. Specifically, an n -point Gaussian formula gives results identical to that of the best approximations utilizing $2n$ points selected arbitrarily. That is, if $f(x)$ in (6.1-15) can be fitted closely with a polynomial of order $2n-1$, then an n -point Gaussian formula is sufficient to integrate $w(x)f(x)$ accurately.

The weights, H_i , and abscissas, x_i , of (6.1-15) are found by use of the family of polynomials orthogonal to the weighting function, $w(x)$, over the interval (a,b) .

$$\int_a^b w(x) Q_n(x) Q_m(x) dx = \begin{cases} h_n^2 & \text{if } m=n \\ 0 & \text{if } m \neq n \end{cases} \quad (6.1-16)$$

where $Q_n(x)$ is an n th order polynomial. In the n -point formula, the abscissas, x_i , are given by the zero crossings (roots) of $Q_n(x)$, and the weights are given by

$$H_i = \frac{1}{Q'_n(x_i)} \int \left(\frac{w(x) Q_n(x)}{x - x_i} \right) dx \quad (6.1-17)$$

where $Q'_n(x_i)$ is the derivative of $Q_n(x)$, evaluated at

$x=x_i$. Unfortunately, there are only a few weighting functions, $w(x)$, for which the corresponding orthogonal polynomials, $Q_n(x)$, are known analytically [3]. Consequently, the practical usage of Gaussian quadrature integration has generally been quite limited. Nevertheless, the orthogonal polynomials associated with any nonnegative weighting function can be generated numerically, permitting the implementation of the technique for an arbitrary $w(x)$, (for instance a color matching function). Although the orthogonalization can be accomplished by the traditional Gram-Schmidt process, a more suitable method for computer implementation is the recurrence relation [3]

$$Q_{n+1}(x) = (x - B_n)Q_n(x) - \left[\frac{h_n}{h_{n-1}} \right]^2 Q_{n-1}(x) \quad (6.1-18)$$

where

$$h_n^2 = \int_a^b w(x) Q_n^2(x) dx \quad (6.1-19)$$

$$B_n = \left(\frac{1}{h_n} \right)^2 \int_a^b x w(x) Q_n^2(x) dx \quad (6.1-20)$$

$$Q_0(x) = 1 \quad (6.1-21)$$

and

$$Q_1(x) = x - \frac{\int_a^b x w(x) dx}{\int_a^b w(x) dx} \quad (6.1-22)$$

The resulting polynomials will be normalized so that their leading coefficients are equal to unity, i.e.

$$\begin{aligned} Q_n(x) &= x^n + b_n x^{n-1} + c_n x^{n-2} + \dots + r_n \\ &= \prod_{i=1}^n (x - x_i) \end{aligned} \quad (6.1-23)$$

The above equations can be easily programmed on a digital computer, using any weighting function and interval. For instance, with $w(x) = 1$, and interval $(-1,1)$, the algorithm generates the classical Legendre polynomials numerically.

The method can be applied to the color matching problem by letting $f(x)$ be the spectral distribution in question, and $w(x)$ one of the color matching functions [4]. For instance the n -point formula for the second XYZ function,

or luminance, would be

$$T_i = \int t_i(\lambda) C(\lambda) d\lambda \approx \sum_{i=1}^n H_i C(\lambda_i) \quad (6.1-24)$$

The weights H_i , and abscissas λ_i , have been calculated using the 1931 CIE color matching curves as the weighting functions. This has also been done with the CIE illuminants A (2856° K), and C (overcast skylight), incorporated into the weighting functions. The results are given in Table I. For instance, the three-point approximation for the luminance is

$$T_2 = \int t_2(\lambda) C(\lambda) d\lambda \quad (6.1-25a)$$

or

$$T_2 \approx 0.1582 C(487.3) + 0.6616 C(558.4) + 0.1802 C(630.6) \quad (6.1-25b)$$

In the case of a sample with reflectance $R(\lambda)$, viewed with illuminant C, denoted $I(\lambda)$

$$T_2 = \int t_2(\lambda) I(\lambda) R(\lambda) d\lambda \quad (6.1-26a)$$

TABLE I

Weights (H_i), and abscissas, (λ), for Gaussian-quadrature method. Abscissas expressed in nanometers.

Unbiased

| order | T_1 | | T_2 | | T_3 | |
|-------|-------------|--------|-------------|--------|-------------|--------|
| | λ_i | H_i | λ_i | H_i | λ_i | H_i |
| 3 | 441.0 | 0.1573 | 487.3 | 0.1582 | 422.2 | 0.2701 |
| | 571.9 | 0.5345 | 558.4 | 0.6616 | 459.6 | 0.6414 |
| | 638.1 | 0.3081 | 630.6 | 0.1802 | 509.6 | 0.0887 |
| 4 | 433.0 | 0.1173 | 463.6 | 0.0462 | 409.6 | 0.0803 |
| | 512.0 | 0.1506 | 529.4 | 0.4537 | 442.6 | 0.5689 |
| | 595.6 | 0.6029 | 589.6 | 0.4410 | 479.3 | 0.3268 |
| | 655.4 | 0.1293 | 651.6 | 0.0591 | 526.8 | 0.0241 |
| 5 | 422.1 | 0.0553 | 449.3 | 0.0178 | 398.1 | 0.0198 |
| | 461.0 | 0.1190 | 507.5 | 0.2280 | 430.4 | 0.3556 |
| | 566.0 | 0.3743 | 559.0 | 0.4997 | 461.0 | 0.4930 |
| | 618.9 | 0.4093 | 612.9 | 0.2338 | 496.9 | 0.1235 |
| | 672.0 | 0.0421 | 666.5 | 0.0207 | 537.6 | 0.0083 |
| 6 | 413.5 | 0.0244 | 440.6 | 0.0084 | 389.8 | 0.0061 |
| | 448.1 | 0.1278 | 488.3 | 0.0966 | 421.0 | 0.1762 |
| | 530.9 | 0.1185 | 535.6 | 0.3955 | 447.0 | 0.4794 |
| | 587.2 | 0.4661 | 583.1 | 0.3771 | 476.4 | 0.2872 |
| | 634.0 | 0.2454 | 630.6 | 0.1140 | 511.2 | 0.0478 |
| | 682.2 | 0.0178 | 676.5 | 0.0084 | 544.5 | 0.0035 |

TABLE I (CONTINUED)

Biased with illuminant A

| order | T_1 | | T_2 | | T_3 | |
|-------|-------------|--------|-------------|--------|-------------|--------|
| | λ_i | H_i | λ_i | H_i | λ_i | H_i |
| 3 | 452.2 | 0.0624 | 502.2 | 0.1616 | 428.6 | 0.0972 |
| | 582.9 | 0.6907 | 570.0 | 0.6498 | 467.2 | 0.2222 |
| | 645.9 | 0.3454 | 638.9 | 0.1886 | 517.8 | 0.0353 |
| 4 | 439.9 | 0.0415 | 476.2 | 0.0390 | 417.1 | 0.0323 |
| | 545.7 | 0.2224 | 539.5 | 0.4350 | 449.6 | 0.1941 |
| | 607.3 | 0.7021 | 598.5 | 0.4587 | 486.9 | 0.1168 |
| | 663.6 | 0.1325 | 657.6 | 0.0673 | 532.3 | 0.0116 |
| 5 | 431.4 | 0.0257 | 457.8 | 0.0111 | 405.3 | 0.0074 |
| | 484.1 | 0.0465 | 517.1 | 0.2116 | 436.1 | 0.1188 |
| | 574.7 | 0.4778 | 568.0 | 0.4933 | 467.3 | 0.1736 |
| | 625.5 | 0.4959 | 619.9 | 0.2581 | 503.6 | 0.0502 |
| | 676.7 | 0.0525 | 670.7 | 0.0258 | 541.1 | 0.0047 |
| 6 | 421.4 | 0.0105 | 446.4 | 0.0042 | 395.4 | 0.0017 |
| | 456.8 | 0.0414 | 498.0 | 0.0833 | 426.2 | 0.0595 |
| | 550.2 | 0.1900 | 543.9 | 0.3699 | 453.0 | 0.1628 |
| | 598.1 | 0.5642 | 590.5 | 0.3972 | 482.5 | 0.1059 |
| | 641.7 | 0.2707 | 636.1 | 0.1338 | 516.5 | 0.0226 |
| | 686.6 | 0.0217 | 679.3 | 0.0115 | 546.9 | 0.0022 |

TABLE I (CONTINUED)
Biased with illuminant C

| order | T_1 | | T_2 | | T_3 | |
|-------|-------------|--------|-------------|--------|-------------|--------|
| | λ_i | H_i | λ_i | H_i | λ_i | H_i |
| 3 | 442.5 | 0.1849 | 484.3 | 0.1601 | 425.9 | 0.3480 |
| | 569.6 | 0.5093 | 556.4 | 0.6679 | 461.5 | 0.7403 |
| | 636.7 | 0.2867 | 628.9 | 0.1719 | 509.7 | 0.0936 |
| 4 | 435.3 | 0.1407 | 462.6 | 0.0523 | 415.3 | 0.1161 |
| | 510.0 | 0.1545 | 527.8 | 0.4521 | 445.4 | 0.6799 |
| | 593.8 | 0.5650 | 587.0 | 0.4393 | 480.6 | 0.3630 |
| | 654.0 | 0.1207 | 650.1 | 0.0563 | 527.5 | 0.0229 |
| 5 | 425.7 | 0.0706 | 449.2 | 0.0210 | 404.4 | 0.0276 |
| | 462.6 | 0.1343 | 505.1 | 0.2269 | 433.5 | 0.4436 |
| | 565.3 | 0.3677 | 557.4 | 0.5107 | 462.8 | 0.5695 |
| | 618.2 | 0.3698 | 611.2 | 0.2219 | 497.4 | 0.1334 |
| | 670.8 | 0.0385 | 665.2 | 0.0195 | 538.4 | 0.0078 |
| 6 | 419.7 | 0.0352 | 440.9 | 0.0100 | 395.2 | 0.0067 |
| | 450.8 | 0.1475 | 486.6 | 0.1043 | 424.4 | 0.2275 |
| | 532.2 | 0.1241 | 534.6 | 0.3961 | 449.3 | 0.5708 |
| | 586.5 | 0.4395 | 581.1 | 0.3756 | 477.7 | 0.3251 |
| | 633.7 | 0.2188 | 629.5 | 0.1062 | 511.5 | 0.0484 |
| | 681.3 | 0.0158 | 675.5 | 0.0078 | 545.0 | 0.0034 |

or

$$T \approx 0.1601 R(484.3) + 0.6679 R(556.4) + 0.1719 R(628.9) \quad (6.1-26b)$$

For the tristimulus integrations considered in this dissertation, the spectral distribution $C(\lambda)$, is the product of the illuminant used, $I(\lambda)$, and the transmissivity of the film sample, i.e.

$$C(\lambda) = I(\lambda) 10^{-\left(\sum_j a_j D_j(\lambda)\right)} \quad (6.1-27)$$

Since the computation in (6.1-27) is the most time consuming step in the calculation of the weighted summation of (6.1-24), the speed of this type of tristimulus integration is basically governed by the number of points used. Inasmuch as an n -point Gaussian formula is equivalent to a $2n$ -point conventional method, (and usually much better), the technique provides a speed advantage of at least 2:1 over the non-Gaussian schemes.

As an illustration of the relative accuracy afforded by the procedure, tristimulus values for the three spectral

distributions of fig. 6.1-2 have been computed using (1) an 80-point Simpson's rule method; (2) the Gaussian scheme for orders 3,4,5 and 6; and (3) the selected ordinate technique for orders 3,4,5 and 6. The results are given in table II. The tabulated errors, E1, and E2, are with respect to the results of the 80-point Simpson's rule, which were assumed to be the exact tristimulus values. The errors were computed by two different methods. First, the standard Euclidean separation was calculated in XYZ space

$$E1 = \sqrt{\Delta(T_1)^2 + \Delta(T_2)^2 + \Delta(T_3)^2} \quad (6.1-28)$$

Secondly, the T vectors were transformed to Lab space, and the Euclidean metric computed

$$E2 = \sqrt{\Delta L^2 + \Delta a^2 + \Delta b^2} \quad (6.1-29)$$

By either criterion, the Gaussian method is clearly superior. It has been found that the fifth and sixth order quadrature formulas give quite acceptable accuracy when used for tristimulus calculations involving

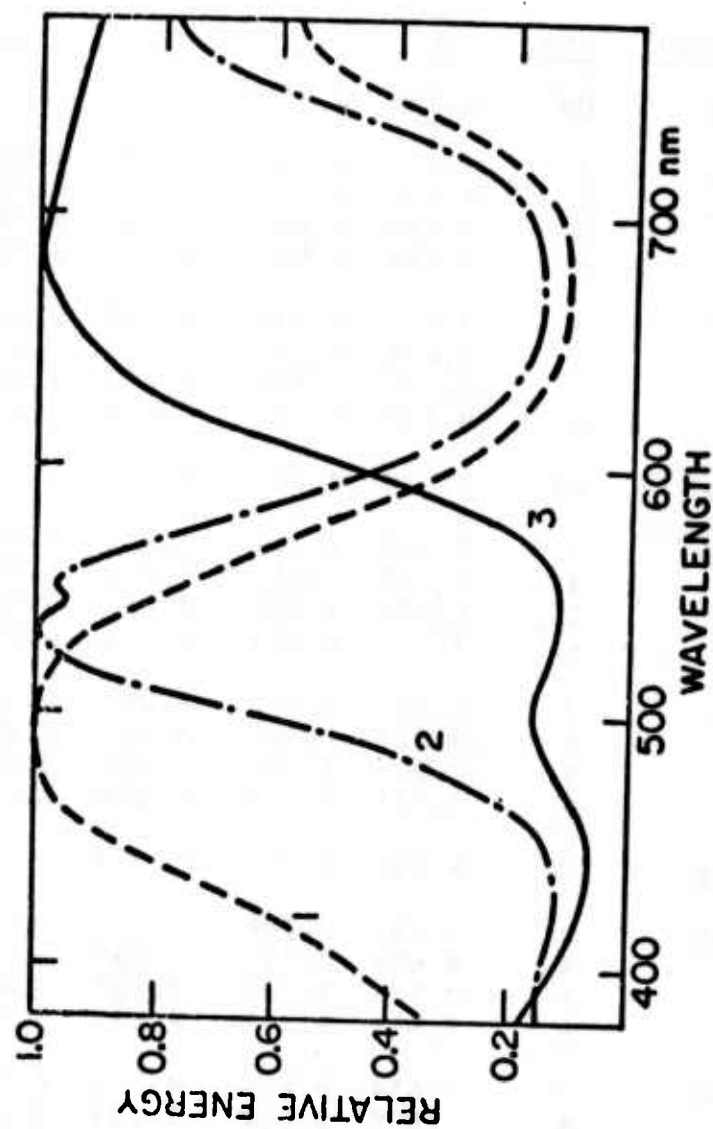


Figure (6.1-2) Spectral distributions of colors used in comparison of tristimulus integration methods

TABLE II

Estimates of tristimulus values for the spectral distributions of fig. 6.2-1, computed by three methods; (1) Simpson's Rule (S-R), (2) Gaussian Quadrature (G-Q), and (3) the Selected Ordinate method (S-O).

| <u>dist.</u> | <u>method</u> | <u>pts.</u> | <u>X</u> | <u>Y</u> | <u>Z</u> | <u>E1</u> | <u>E2</u> |
|--------------|---------------|-------------|----------|----------|----------|-----------|-----------|
| 1 | S-R | 80 | 0.399 | 0.551 | 0.733 | | |
| | | | | | | | |
| | G-Q | 3 | 0.412 | 0.568 | 0.739 | 0.023 | 2.095 |
| | | 4 | 0.386 | 0.555 | 0.732 | 0.014 | 5.962 |
| | | 5 | 0.398 | 0.559 | 0.736 | 0.008 | 2.723 |
| | | 6 | 0.398 | 0.549 | 0.732 | 0.003 | 0.417 |
| | S-O | 3 | 0.442 | 0.540 | 0.720 | 0.047 | 19.660 |
| | | 4 | 0.431 | 0.542 | 0.728 | 0.034 | 14.705 |
| | | 5 | 0.426 | 0.548 | 0.730 | 0.028 | 11.208 |
| | | 6 | 0.430 | 0.551 | 0.730 | 0.031 | 11.499 |
| | S-R | 80 | 0.526 | 0.813 | 0.251 | | |
| | | | | | | | |
| | G-Q | 3 | 0.570 | 0.834 | 0.251 | 0.049 | 9.919 |
| | | 4 | 0.510 | 0.825 | 0.252 | 0.020 | 8.593 |
| | | 5 | 0.520 | 0.826 | 0.252 | 0.015 | 5.554 |
| | | 6 | 0.527 | 0.806 | 0.252 | 0.007 | 2.216 |
| | S-O | 3 | 0.493 | 0.711 | 0.283 | 0.111 | 20.701 |
| | | 4 | 0.518 | 0.726 | 0.279 | 0.091 | 23.223 |
| | | 5 | 0.539 | 0.758 | 0.275 | 0.062 | 20.065 |
| | | 6 | 0.533 | 0.773 | 0.268 | 0.044 | 13.826 |
| 3 | S-R | 80 | 0.696 | 0.707 | 0.124 | | |
| | | | | | | | |
| | G-Q | 3 | 0.704 | 0.704 | 0.124 | 0.008 | 3.980 |
| | | 4 | 0.699 | 0.711 | 0.125 | 0.005 | 0.554 |
| | | 5 | 0.696 | 0.707 | 0.125 | 0.002 | 0.645 |
| | | 6 | 0.698 | 0.706 | 0.125 | 0.002 | 0.911 |
| | S-O | 3 | 0.654 | 0.667 | 0.141 | 0.061 | 6.847 |
| | | 4 | 0.680 | 0.672 | 0.138 | 0.041 | 7.597 |
| | | 5 | 0.695 | 0.683 | 0.135 | 0.027 | 8.256 |
| | | 6 | 0.682 | 0.688 | 0.132 | 0.025 | 3.232 |

photographic dyes. The use of these techniques makes possible colorimetric calculations which would otherwise consume excessive amounts of computer time.

REFERENCES

1. Ralston, A. A First Course in Numerical Analysis, New York: McGraw- Hill. 1965
2. Wyszecki, G. and Stiles, W.S. Color Science, New York: Wiley and Sons. 1967
3. Beckmann, P. Orthogonal Polynomials for Engineers and Physicists, Boulder:Golem Press. 1973

INDIRECT INVERSION AND COMPUTATIONAL CONSIDERATIONS

As indicated in the previous section, the film model can be inverted accurately by the use of recursive techniques. This section presents two entirely different approaches, which avoid the formal inversion completely. As usual, the increase in speed is at the expense of accuracy. However, it is felt by the author that even with these approximate numerical procedures, the dominant sources of subjective errors in the colorimetric corrections are due to inadequacies in the visual model upon which the calculations are based, and not on the numerical errors themselves.

The first method to be described involves the fitting of an easily computable mathematical vector function to the mapping between the input and output tristimulus values.

$$\underline{P}(\text{in}) = \underline{f}[\underline{P}(\text{out})] \quad (6.2-1)$$

where $\underline{P}(\text{out})$ is the desired output tristimulus vector, and $\underline{P}(\text{in})$ is the corresponding input vector. By simulating the model for a large number of input tristimulus vectors,

and associating with these their corresponding output tristimulus vectors, statistical techniques can be invoked to fit a 3×3 matrix to the generated mapping. However, since the film mapping is inherently nonlinear, this approach is quite inaccurate, as might be expected. Nevertheless, by "extending" the output vector to include nonlinear components, a nonlinear vector function can be fitted. This technique is discussed in APPENDIX B. The approach offers a substantial savings in computation time, since all the functions within the model, the Hurter-Driffeld curves, the dye system nonlinearities, etc. are all coalesced into one large matrix multiply. This method is capable of producing an almost arbitrarily high degree of accuracy, as the size of the extended predictor vector is raised. Of course, the point of diminishing returns is soon reached. The transform which was found to give the best compromise between precision and complexity <*> is

 <*> The transform was evaluated in the RGB coordinate system of the television phosphors, denoted by the vector \underline{P} , (as opposed to a system such as XYZ), in order to eliminate coordinate conversions, and hence, maximize speed.

6.2

$$\begin{bmatrix} P_1(\text{in}) \\ P_2(\text{in}) \\ P_3(\text{in}) \end{bmatrix} = \begin{bmatrix} & & & & & & & & \\ & & & & & & & & \\ & & A & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \end{bmatrix} \begin{bmatrix} P_1(\text{out}) \\ P_2(\text{out}) \\ P_3(\text{out}) \\ [P_1(\text{out})]^2 \\ [P_2(\text{out})]^2 \\ [P_3(\text{out})]^2 \\ P_2(\text{out}) P_3(\text{out}) \\ P_1(\text{out}) P_3(\text{out}) \\ P_1(\text{out}) P_2(\text{out}) \\ [P_1(\text{out})]^3 \\ [P_2(\text{out})]^3 \\ [P_3(\text{out})]^3 \end{bmatrix} \quad (6.2-2)$$

where $P_1(\text{out})$, $P_2(\text{out})$, $P_3(\text{out})$ are the weights on the color monitor phosphors corresponding to the desired image, while $P_1(\text{in})$, $P_2(\text{in})$, and $P_3(\text{in})$ are the tristimulus values required for the pre-distortion.

In calculating the matrix of eqn. (6.2-2), 125 colors, distributed throughout $P(\text{in})$ -space were generated, and used as inputs to a simulation of the model. These were then paired with their corresponding $\underline{P}(\text{out})$ -vectors, and the matrix relationship computed by the methods of

APPENDIX B. The average error of the fit, over all 125 colors was 0.02, with the error, E, defined as

$$E = \sqrt{\frac{\Delta P_1^2(\text{in}) + \Delta P_2^2(\text{in}) + \Delta P_3^2(\text{in})}{P_1^2(\text{in}) + P_2^2(\text{in}) + P_3^2(\text{in})}} \quad (6.2-3)$$

where $\Delta P_j(\text{in})$ is the difference between the actual value of $P_j(\text{in})$, and its estimate, computed using eqn. 6.2-2.

The second and fastest alternative investigated was a large look-up table. Since the model is merely a device to associate input vectors with output vectors, it could theoretically be replaced by a large look-up table with "sufficiently close" entries. Such a table could be used to assign to every desired output vector, a required input vector. Unfortunately, since each entry in the table would be three dimensional, $3N$ numbers would be required for N -level quantization in red, green, and blue. Even for 32 level quantization, the table would need to store about 96,000 numbers. The advantage of this approach is clearly its speed, which would be limited only by the input-output capability of the computer used. Although, the generation of the table might expend a great deal of computer time, the cost would be more than compensated for if numerous pictures were processed.

A variation of this scheme might be the combination of a three dimensional interpolation technique with the look-up table. This modification would involve a trade-off between computational demands and storage requirements. A second possible variation might be the use of uniform color scales for the quantization. It has been demonstrated [1] that such a procedure reduces (or at least renders more uniform) the perception of colorimetric errors when tristimulus values are quantized. The drawback here would be the additional amount of computation required, partially defeating the purpose of the look-up table.

Of the methods described in this section, the approach represented by eqn. 6.2-2 is believed to offer the best compromise between accuracy and computational speed for high resolution images. A great advantage of the method is that the channel mapping need never be inverted directly. All that is required, is a statistically significant collection of input tristimulus vectors paired with their corresponding output vectors. This ensemble of vectors can be generated by simulating the process in question with a large number of inputs, or by strictly empirical techniques. That is, if no accurate model is

available, the ensemble of input and output vectors can be generated experimentally by using colorimeters, spectro-radiometers, etc.

The next section describes some experimental results, and gives practical suggestions for the actual implementation of the pre-distortion procedure. The method chosen was that described by eqn. 6.2-2.

REFERENCES

1. Frei, W. A New Model of Color Vision and Some Practical Implications, USCIP Report 530, p. 128, 1974

EXPERIMENTAL RESULTS

A good experiment will often discredit the theory which it attempts to verify; or at least betray its weakest aspects. Hence, the formulation of a theoretical model is typically a recursive process, driven by experimental feedback.

In the case of the film recording problem, the experimental procedure used was as follows. A 256x256 digital color image was displayed on a Conrac color monitor, and photographed on EKTACHROME-X film. A reproduction of the resulting transparency is shown in fig. 6.3-1 (color plate). The original digital image was then processed by the inverse of the display-film mapping, and the resulting pre-distorted image also photographed. The method of eqn. 6.2-2 was used to effect the pre-distortion. The actual numerical transformations used are provided in APPENDIX-D. Both the original slide and the one taken of the pre-distorted image were placed together in a slide mount, projected, and compared with the actual original image on the television monitor, in a darkened room. This process was repeated dozens of times. For each recursion, a modification of some part of the

model was implemented, in an attempt to "converge in" on the subjectively best result. What follows is a summary of what was learned.

(1) Correction for chromatic adaption is of primary importance. In spite of this, adaption is completely ignored in most of the colorimetric analyses of which the author is aware. Simply normalizing the XYZ tristimulus values (as is done in the Lab system), is totally inadequate for tungsten illumination.

(2) It was originally assumed that a transparency exposed to produce a neutral density $\langle * \rangle$, would appear neutral when projected. Surprisingly, this is not so. The reason is that the eye never completely adapts to the "white" of the slide projector at the typical brightness levels encountered. Tungsten illumination will usually retain a slightly yellow-red character, even when viewed in a darkened room. Consequently, the quality of the white balance is usually improved by a slight cyan-blue bias in the transparency being viewed. One of the most significant, yet difficult problems encountered, was a definition of the sensation of "white" in mathematical

 $\langle * \rangle$ A neutral density is one which merely attenuates the illuminant, and does not affect its chromaticity or color balance.

terms. Most analytical descriptions of projected transparencies define white to be the projector's illuminant. This was found to be unrealistic, in that a transparency would need to produce zero density in order to generate such a white. A more suitable definition, takes this into account, using the minimum density of the slide in question as the reference white (this is expressed quantitatively in eqn. 4.2-12).

(3) The nonlinear nature of film cannot be ignored. Any inversion of the display-film mapping must take into account the "operating point" of the film. In other words, a particular color balance on the monitor will result in a particular set of dye densities. In the mathematical simulation of this process, precisely the same points on the Hurter-Driffeld curves must be used if there is to be a good equivalence with reality. For quantitative work, merely assuming that the film is operating in the linear portion of the H&D curve is inadequate. If the film is exposed to a certain white point, the simulation must reflect this exactly. Furthermore, the exposure is critical. Variations of less than one-half f-stop <*> result in quite noticeable color shifts.

6.3

(4) At the outset of this research it was thought that the limiting factor in reproduction fidelity would be the variance in the photographic process, due to either inconsistent film and/or processing. However, this notion was unfounded. With reasonable care (i.e. keeping film refrigerated and using fresh processing chemistry) excellent consistency resulted.

(5) Gamut restrictions also proved to be less of a problem than was anticipated. Since the eye is insensitive to slight changes in image contrast, scaling the red, green, and blue pixels by a constant less than unity reduced the required dynamic range so that neither the film or display gamuts were exceeded. This procedure is discussed in Sec. 3.2.

A reproduction of the transparency obtained from photographing the pre-distorted image is shown in fig. 6.3-2 (color plate). The most noticeable improvements are in the facial tones and the more saturated colors. Furthermore, a comparison of the "before and after" transparencies with the original image displayed on the color monitor, showed that the slide

<*> An f-stop is a photographic term, corresponding to a doubling of the exposure.

taken of the pre-distorted image gave a very good match, while the colorimetric errors in the slide taken of the original image were quite considerable.

Since the errors in film recording are typically losses in saturation, the photograph of the pre-distorted image tends to have a higher average saturation. Thus the procedure gives the added bonus of a brighter, subjectively more pleasing image, than would be the case if the distortions went uncorrected.

The last section will describe some additional projects being contemplated which would utilize the techniques which have been developed.

GENERALIZATIONS AND TOPICS FOR FUTURE RESEARCH

In a more abstract sense, this dissertation has described the modeling and inversion of a three dimensional nonlinear channel, with the intent of removing the degradations induced upon the "signals" passing through the channel. Although the specific "channel" discussed was color television-to transparency film, the conceptual framework could easily encompass other applications. For instance, in the case of reflection type color prints [1], or reproductions from an ink-jet plotter, the only modification required would be a generalization of the the Lambert-Beer law. Recall that the relationship between the incident illumination, $I(\lambda)$, and the transmitted spectral distribution, $C(\lambda)$, for a transparency with dye densities d_j is given by

$$C(\lambda) = I(\lambda) 10^{-\sum_j d_j D_j(\lambda)} \quad (7-1)$$

With dye mixtures viewed by reflected light, the equivalent relation is much more complicated. Viewing

angle, illumination geometry, and multiple internal reflections must be considered [1]. However, the spectral distribution will remain a function of the dye densities.

$$C(\lambda) = I(\lambda) f(d_1, d_2, d_3) \quad (7-2)$$

Once this relationship has been determined, it can be substituted for eqn. 7-1, enabling the procedure to be used for reflection-type materials.

Heretofore, it has been assumed that the image in question was available in digital form. Clearly, the original photograph must have been scanned, or digitized in some fashion initially. As with the process of "writing out" an image from computer to photograph, the converse process of "reading in" an image from photograph to computer can be treated as a three dimensional channel. Recall that the relationship between the "scanner signals," S_i , and the dye densities, d_j , is given by

$$\frac{\int 10^{-\sum_j d_j D_j(\lambda)} I(\lambda) s_i(\lambda) d\lambda}{\int I(\lambda) s_i(\lambda) d\lambda} = S_i \quad (7-3)$$

where the $s_i(\lambda)$ are the spectral sensitivities of the scanner's color separation channels. Given the scanner signals, the dye densities can be computed from the inversion of 7-3. Working backwards through the film model, the layer exposures which gave rise to the dye densities can also be calculated. Hence, correction for moderate overexposure or underexposure can be achieved <*>. For instance, in the case of underexposure, the image can be scanned, the layer exposures deduced and scaled up to the correct level, and a new photograph written out in such a way as to force the correct layer exposures.

A non-deterministic variation of this procedure can be used to correct for the fact that the layer sensitivities of film are not color matching functions. Using the concepts of estimation theory, and the statistics available on the spectral distributions of natural scenes,

<*> If the horizontal portions (shoulder or toe) of the HD curve are involved, some of the signal has been "clipped," and therefore irretrievably lost. Hence, the limitation to moderate exposure errors.

estimates of the actual tristimulus values which gave rise to the layer exposures in question can be calculated <*>.

Having gained some degree of control over the film recording process, there is no reason to limit colorimetric manipulations to the duplication of an image. Various types of enhancement are also easy to achieve. A simple modification of the film recording procedure which has been described could be used to increase the saturation of the image colors beyond that required for mere fidelity of reproduction. This could be advantageous in aerial imagery, or in applications where the visibility of subtle transitions is important.

An interesting prospect in the category of enhancement is the generalization of homomorphic filtering to color images. A possible approach is the use of the visual model (see APPENDIX C), in an effort to exaggerate the type of edge enhancement performed by the eye. This would result in the added benefit of reducing the dynamic range of a color image, without sacrificing the saturation of its hues. Also, the problems associated with clipping due to the limited dynamic range of high gamma films would be

<*> This is an underdetermined problem, discussed in [2].

minimized.

In conclusion, it must be admitted that the type of systems approach advocated in this dissertation does not always give results superior to those produced by "ad hoc" techniques. It is the opinion of the author that this is due to the inadequacy of the mathematical models on which some of the more scientific approaches are based. For instance, a colorimetric analysis based solely on raw tristimulus values (ignoring adaption effects) will usually give absurd results, distictly inferior to those attainable by use of common sense and qualitative techniques. In the parlance of the computer community, "garbage in, garbage out." If a model is only a first approximation to reality, then exacting calculations carried out to ten decimal places of accuracy are often no more than an academic exercise.

However, there is little doubt that our universe is an orderly one, even though its laws may elude us. In image processing, the laws of optics and photometry are well understood, but those governing the human visual system's interpretation of imagery retain a degree of mystery. It is in this area of a visual fidelity criterion that great strides can be made. The work described in this

dissertation has reflected this philosophy, and it is hoped, has made some progress towards the above goal.

REFERENCES

1. Ohta, N. Color Reproduction Quality with Color Films and Color Prints, The Journal of Photographic Science, Vol 20, p. 239, 1972
2. Mancill, C. Color Image Restoration by Linear Estimation Methods, USCIPi rept. 530, p. 80, 1974

APPENDIX A

TRANSFORMATIONS BETWEEN SETS OF TRISTIMULUS VALUES

As shown in Sec. 2.2, the tristimulus values associated with a particular color can be expressed in many different coordinate systems, all of which are linearly interrelated. For example, the television phosphor weights, \underline{P} , and the XYZ tristimulus values, \underline{T} , are related by a 3x3 matrix multiply

$$\underline{T} = \underline{A} \underline{P} \quad (\text{A-1})$$

In order to determine the coefficients of the matrix \underline{A} , the spectral characteristics of the three phosphors, $p_j(\lambda)$, must be known. The relationship then follows directly

$$T_i = \int t_i(\lambda) C(\lambda) d\lambda \quad (\text{A-2a})$$

$$T_i = \int t_i(\lambda) \sum_j P_j p_j(\lambda) d\lambda \quad (\text{A-2b})$$

$$T_i = \sum_j a_{ij} P_j \quad (\text{A-2c})$$

APPENDIX A

$$a_{ij} = \int t_i(\lambda) p_j(\lambda) d\lambda \quad (A-3)$$

Typically however, the $p_j(\lambda)$ distributions are known only in relative units, and arbitrary but unknown scale factors on each of the $p_j(\lambda)$ prevent the above determination. Alternatively, the chromaticities of the individual phosphors may be known, and the A matrix may be desired for the case in which the white point of the monitor is balanced to a particular color temperature. For instance, in Sec. 2.3 the matrix required was one such that unity values for the P_j resulted in a white point of 6500° K. Also, the transformation might be known for a particular white point, but the transformation for a second white point may be desired. The problem then, is to scale the columns of A in such a way as force the desired white point. Assuming that the unscaled matrix A has been computed using eqn. (A-3), the scale factors sought are the k_1 , k_2 , and k_3 which satisfy

$$\begin{bmatrix} a \end{bmatrix} \begin{bmatrix} k_1 \\ k_2 \\ k_3 \end{bmatrix} = \begin{bmatrix} W_1 \\ W_2 \\ W_3 \end{bmatrix} \quad (A-4)$$

APPENDIX A

where the W are the XYZ coordinates of the desired white point. If the matrix were already perfectly scaled, the required values of k would all be unity. Since both W and A are known, the \underline{k} vector can be solved for by inverting the A matrix

$$\underline{k} = [A]^{-1} \underline{W} \quad (A-5)$$

The correctly scaled A matrix is then obtained by absorbing the scale factors, k_m into the matrix

$$a'_{mn} = a_{mn} k_m \quad (A-6)$$

That is, the m th column is scaled by the m th scale factor. This insures that unity phosphor weights yield the desired white point. Using this technique, the transform from the phosphor tristimulus values to the XYZ tristimulus values was calculated as

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} .489 & .324 & .161 \\ .264 & .672 & .064 \\ .014 & .134 & .806 \end{bmatrix} \begin{bmatrix} R \\ G \\ B \end{bmatrix}$$

APPENDIX B

MULTIDIMENSIONAL MAPPINGS, AND LEAST SQUARES CURVE FITTING OF SURFACES

A common problem in applied science is the fitting of a mathematical curve to empirical data. In colorimetry, since data is usually three dimensional, curve fits must be generalized to higher dimensionalities. The simplest problem is, given an ensemble of vectors, \underline{x} , and \underline{y} , which are related by some (possibly stochastic) vector function \underline{f}

$$\underline{y} = \underline{f}(\underline{x}) \quad (\text{B-1})$$

find the matrix, A , which minimizes the error

$$e = E \|\underline{Ax} - \underline{y}\|^2 \quad (\text{B-2})$$

where E denotes the expected value, or average over all samples. This is a classical problem in Estimation Theory, and is usually treated in a statistical context. The optimal A matrix (in the mean square sense) is the solution to

APPENDIX B

$$E\{[A\underline{x} - \underline{y}] [\underline{x}']\} = 0 \quad (B-3)$$

where \underline{x}' is the transpose of the column vector \underline{x} . This is the celebrated orthogonality principle, which is easily proved. The error in (B-2) can be expressed as

$$e_i = \sum_{n=1}^N \left[\sum_{j=1}^3 a_{ij} x_j(n) - y_i(n) \right]^2 \quad (B-4)$$

where the indexes (n) refer to the nth vector sample, and the subscripts denote the element within the vector. Taking partial derivatives of the error with respect to the a_{ij} coefficients and setting the result equal to zero gives

$$\frac{\partial e_i}{\partial a_{km}} = \sum_j a_{ij} \sum_n [x_j(n) - y_i(n)] x_m(n) \quad (B-5a)$$

$$0 = \sum_{j=1}^3 a_{ij} \left(\sum_n x_j(n) x_m(n) - \sum_n y_i(n) x_m(n) \right) \quad (B-5b)$$

This is equivalent to the orthogonality condition stated above in matrix notation. The required matrix, A, is thus given by

APPENDIX B

$$A = [E\{\underline{y} \underline{x}'\}] [E\{\underline{x} \underline{x}'\}]^{-1} \quad (B-6)$$

One of the advantages of this technique, is that it is not limited to square matrices, i.e. the lengths of the \underline{x} and \underline{y} vectors need not be equal. This property can be used to great advantage in the fitting of nonlinear functions to the data. Consider the estimation of the vector \underline{y} using only the vector \underline{x} . In analogy with the one dimensional case, more flexibility in the curve fit can be achieved by allowing higher powers of the predictor variables, x_j , in the estimation of the predicted variables, y_j . An approach which has been found fruitful [*] is the postulation of the following structure between the \underline{x} and \underline{y} domains

[*] A similar approach, known as a phi-machine, is used to generate nonlinear separating surfaces in pattern recognition algorithms [1].

APPENDIX B

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} & & & & & & & \\ & & & & & & & \\ & & & & & & & \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ (x_1) \\ (x_2) \\ (x_3) \\ x_2 x_3 \\ x_1 x_3 \\ x_1 x_2 \\ \cdot \\ \cdot \end{bmatrix} \quad (B-7)$$

The extended \underline{x} vector can include any nonlinear functions of the x_j , such as reciprocals, logarithms, etc. Note however that the relationship in (B-7) is still linear in the coefficients a_{ij} . Therefore, the A matrix can be found using (B-6). The accuracy of the fit depends of course on the number and type of nonlinearities used. By making the extended \underline{x} vector sufficiently long, the desired degree of accuracy can usually be achieved.

This method can be used to advantage whenever multidimensional nonlinearities are involved. If great accuracy is needed, the technique can be used to calculate

APPENDIX B

an estimate of the desired vector, to use as the initial guess in a recursive procedure (Sec. 6-1).

Since the matrix multiply of (B-7) is a completely deterministic (as opposed to recursive) operation, it provides a substantial savings in computation time. Therefore, if a degree of accuracy can be sacrificed, a significant gain in computational efficiency can be attained for nonlinear colorimetric problems.

REFERENCES

1. Nilsson, N.J. Learning Machines, New York:McGraw-Hill, 1965

APPENDIX CSPATIAL FILTERING OF COLOR IMAGES

Although the subject of this thesis is the recording of color images, with emphasis on colorimetric fidelity, certain enhancement techniques can be used to achieve a subjective improvement in the appearance of an image (as opposed to mere duplication). One such technique is spatial filtering. It is clear that since the human visual system attenuates the lower spatial frequencies, the higher frequency components must be of greater utility. Also, from the standpoint of information theory, the less correlated signal components possess greater entropy, and hence more information. High emphasis spatial filtering has long been used to enhance and/or restore monochrome imagery. In the case of color images, very little work has been done; perhaps because it is unclear how the extra dimension(s) of color should be handled. The most obvious (and naive) technique is to filter the red, green, and blue sub-images independently, as if each were a separate monochrome image. This approach leads to very disappointing results, because the color balance of the original image is invariable

APPENDIX C

destroyed. Clearly, the processing must assure that the hues of the individual color patches in the image do not suffer large shifts. A very satisfactory procedure for accomplishing this is based on the visual model of fig. 2.1-3. Note that the filtering in the model takes place subsequent to the cross-coupling of the channels. As discussed in the chapter on chromatic adaption (Sec. 2.3), the effect of this on the chromatic channels is basically an origin shift, compensating for changes in illuminant. Concerning the actual shapes of the filter functions in each of the channels, there is good evidence that the luminance channel provides most of the high emphasis. The plot of fig. C-1 indicates the relative frequency responses that are thought to exist in the chromatic and luminance channels <*>. The hypothesis that the chromatic response drops off before that of the luminance frequency response is supported by the fact that Mach-band phenomenon are generally stronger for luminance edges than for chrominance edges. This phenomenon is exploited in the N.T.S.C. color television system, by devoting most of the signal bandwidth to the luminance

<*> The exact shape of the chromatic response function(s), especially at low spatial frequencies, is still open to question at the time of this writing.

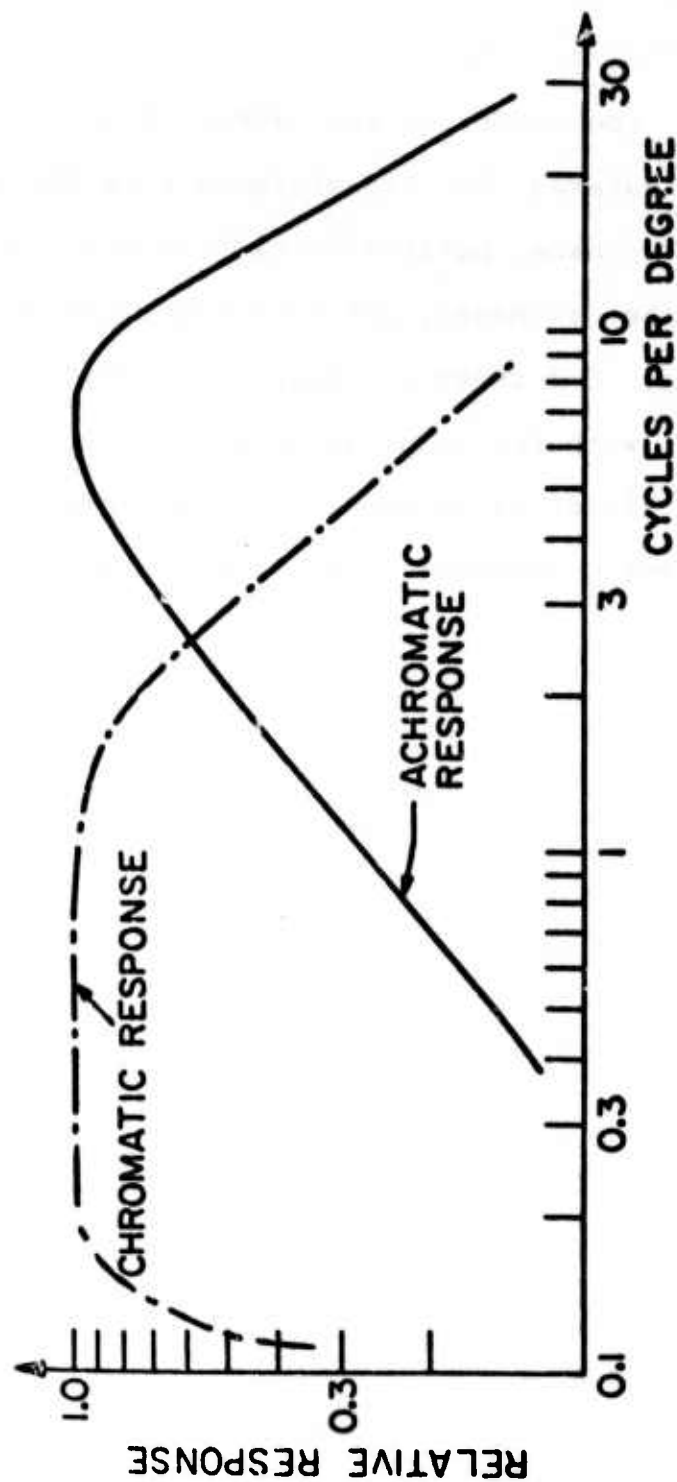


Figure (C-1) Frequency response of chromatic channels in visual system model

APPENDIX C

component. Consequently, the visual system's processing can be simulated by transforming from RGB space to the perceptual q-space, performing high emphasis filtering on each of the q-images, and then transforming back to RGB space. Fig. C-2 shows an image processed in this manner, contrasted with the original (fig. 6.3-1). Note that the subjective effect is an increase in sharpness coupled with an apparent increase in the uniformity of the illumination.

APPENDIX D

TRANSFORMATIONS USED IN THE FILM MODEL

RGB phosphor weights to XYZ tristimulus values (6500° K white point)

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} .489 & .324 & .161 \\ .264 & .672 & .064 \\ .014 & .134 & .806 \end{bmatrix} \begin{bmatrix} R \\ G \\ B \end{bmatrix} \quad (D-1)$$

XYZ tristimulus values to RGB phosphor weights

$$\begin{bmatrix} R \\ G \\ B \end{bmatrix} = \begin{bmatrix} 2.99 & -1.34 & -.490 \\ -1.06 & 1.97 & .056 \\ .111 & -.260 & 1.03 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} \quad (D-2)$$

Phosphor weights to layer exposures

$$\begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} = \begin{bmatrix} 1.79 & .668 & 1.46 \\ .770 & 2.00 & .112 \\ .325 & 1.09 & 1.74 \end{bmatrix} \begin{bmatrix} R \\ G \\ B \end{bmatrix} \quad (D-3)$$

APPENDIX D

Interimage effect matrix

$$\begin{bmatrix} d_1' \\ d_2' \\ d_3' \end{bmatrix} = \begin{bmatrix} 1.18 & -.17 & -.01 \\ -.16 & 1.35 & -.19 \\ -.12 & -.19 & 1.31 \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ d_3 \end{bmatrix} \quad (D-4)$$

Matrix of eqn. 6.2-2 (in transposed form). Used to perform pre-distortion.

$$\begin{bmatrix} 2.982 & -0.136 & 0.886 \\ -1.369 & 1.744 & -1.480 \\ 0.113 & 0.178 & 2.453 \\ -1.654 & -0.302 & -1.259 \\ 0.771 & -1.276 & 0.710 \\ -0.580 & -0.261 & -1.903 \\ 0.211 & -0.004 & -0.008 \\ 0.082 & -0.038 & 0.015 \\ -0.514 & 0.040 & 0.453 \\ 0.897 & 0.122 & 0.454 \\ -0.291 & 0.616 & -0.406 \\ 0.171 & 0.140 & 0.894 \end{bmatrix}$$

(D-5)