

# Human Observation Based Generalized Gamma Functions

Attila Neumann\* Alessandro Artusi\* László Neumann† Georg Zotti\* Werner Purgathofer\*

\* Institute of Computer Graphics and Algorithms  
University of Technology, Vienna, Austria  
email: {aneumann, artusi, gzotti, wp} @cg.tuwien.ac.at

† Grup de Gràfics de Girona, Universitat de Girona, Spain  
Institució Catalana de Recerca i Estudis Avançats, ICREA, Barcelona  
email: lneumann@ima.udg.es

December 30, 2004

## ABSTRACT

This paper describes an accurate method to obtain the Tone Reproduction Curve (TRC) of display devices without using a measurement device. It is an improvement of an existing technique based on human observation, solving its problem of numerical instability and resulting in functions in log-log scale which correspond better to the nature of display devices. We demonstrate the efficiency of our technique on different monitor technologies, comparing it with direct measurements using a spectrophotometer.

**Keywords:** Display Measurement, Human Visual System, Spatial Vision

aspects of human comparison based measurements are detailed with the frame of its application, also with setup of single measurements with the next step and stop criterium, respectively. The chapter about the core mathematical problem discusses 2 different methods: one is defined in a previous paper<sup>1</sup> of the authors, describes a smoothness problem in a linear coordinate system. The other method transforms the problem into a log-log scale. The latter method has to face a numerical problem of the enormously different coefficients of the minimum problem, due to the log-log scaling. We show a two-pass method solving the numerical problem and making the algorithm faster and more reliable at the same time. The paper closes by discussion of the new method and future work.

## 1. INTRODUCTION

The colorimetric characterization of display devices is performed in two steps. The first step is the definition of the *Tone Reproduction Curve* (TRC) which describes the relationship between the input signal of the monitor and the luminance produced on the screen. The second step consists in the definition of a matrix that describes the additive nature of the display device.

Several methods have been proposed to solve both steps in the recent past, and of particular interest is how to define the TRC without the use of a spectrophotometer. These methods use properties of human perception to achieve the final goal. This paper introduces a method to solve the first step without the necessity to use a spectrophotometer. It is based on the work presented in a previous paper,<sup>1</sup> solving its problem of occasional numerical instability and unequal weighting.

The paper is structured as follows. After overviewing related work, we introduce our topic by explaining the principle of a human comparison based measurement system. Following human vision aspects are investigated as spatial vision system and contrast sensitivity. Then mathematical

## 2. RELATED WORK

Several display characterization models with different characteristics have been presented in the past. These models can be classified into two basic categories: *measuring device based* and *human vision based* models.

Many works have been presented for the first category,<sup>2-7</sup> trying to model the real Tone Reproduction Curve (TRC) characteristic of a CRT display device. In many cases these models are still not accurate enough to acquire the real TRC, but just an approximation of it. In addition, even if they reach sufficient accuracy for CRT displays, this accuracy is not achieved for LCD displays. In consequence, the users are not able to gain the high precision required for many applications. On the other hand a model for LCD displays has been proposed<sup>8</sup> which introduces a spline interpolation in order to estimate the TRC.

In any case, these models require a spectrophotometer to get the information necessary to describe the TRC.

The models of the second category are based on interaction with the user and based on human vision and observation. One example is used in commercial software such as Adobe Gamma<sup>TM</sup> which comes with Adobe Photoshop<sup>TM</sup>.

While in the first category acceptable quality can be achieved, in the second one this goal has not been achieved until now for two reasons. First, the current models are not able to estimate or compute the real TRC, but only a simplification of the model used in the first category. Second, the applied mathematical background in these models is typically restricted to describe a kind of simple exponential gamma correction function. In order to solve this problem, a more accurate characterization model of display devices based only on human observation has been presented.<sup>1</sup> In this paper the human vision is used to compare a series of dithered color patches against interactively changeable homogeneously colored display areas, obtaining the TRC of display devices.

### 3. PRINCIPLE

Characterization of a display based only on human perception has advantages as well as disadvantages. The obvious disadvantage is implied by the adaptation mechanism of human perception, making it impossible to define absolute values: the result of human observation always depends on the environmental circumstances. Only a direct comparison of two adjacent regions can give a reliable result, especially the detection of their identical appearance, which is used also by our method.

We perceive an arbitrary pattern as a homogeneous field if the spatial frequency is high enough, i.e., if the view angle of dots in a pattern is below a certain threshold. We can not distinguish luminance contrasts of two homogeneous patches if luminance difference is below a given ratio of absolute luminance, according to the Weber-Fechner law or to the more accurate color appearance models. This ratio depends on absolute luminance level (given in  $[cd/m^2]$ ) and whether photopic or scotopic vision is used. Fortunately, display devices are typically observed in the photopic range. In this case the 1% difference of *Weber-Fechner's law*<sup>9</sup> roughly holds. The human visual system can perceive a 0.5% contrast at 100  $cd/m^2$  average ambient luminance level and at 8 *cycles/degree* spatial frequency. This ratio, and also the error, will increase for the darkest regions.

In the visual comparisons used in our method, limits of spatial vision in the perception of homogeneous regions will be used, but we would like to work with the possible minimal just noticeable luminance difference in order to ensure the highest accuracy. We therefore have to use an optimal balance of conflicting requirements for the sake of spatiochromatic features of observation.

Fortunately, observations of the apparent identity of two neighbouring color patches will yield constant results for

widely changing environmental lighting conditions. It is arguably the only accurate human observation, while changing circumstances can drastically change any absolute values. For example, a badly lit white table will still appear white, in spite of its absolute color appearing dark gray. The environmental lighting affects only the accuracy of the observation of emissive color patches, i.e., in case of pure emission the observation is more accurate than in case of added ambient lighting comparable to the investigated self emission.

Another point is that using inhomogeneous ambient illumination on the area of the display itself, the increase of illumination on the support area of the comparison reduces the *accuracy* of the comparison, but not its *result*, since this variation modifies the perceived difference between the two small neighbouring areas only slightly. In addition, when the observer recognizes that the lighting is not uniform, the comparison will be subject to a self-correction rather than yield a bad result, and since the comparison does not deal with separate border points, but the whole borderline, equality will be perceived when all "trustworthy" border sections seem balanced. This means, the observer automatically ignores the parts of the border where the non-uniform environmental illumination appears to be changing too rapidly. Any difference in the remaining, "valid" part of the display will be visible and noted by the observer.

Commercial graphics software sometimes includes applets which allow to find the gamma value of the display based on the same principle, so our method can be seen as an extension of them. All TRCs yielded by this method can be seen only as relative, depending on the maximum and minimum screen brightness values, so they need a few additional measurements to complete them by these absolute values. However, for many applications it is sufficient to know the relative curves, so the real question is the accuracy that can be reached with our method, compared to an instrument-based method. Now, the advantage of this method is that its accuracy corresponds to that of human perception, which cannot be asserted for all the measuring devices used for such purposes. In other words, we can achieve just the required accuracy if the same environment is given during the measurement process and during some correction operation depending on its results. In addition, it is easy to choose more appropriate circumstances for the measurement process, e.g., a dark environment and observation and judgement by more persons, which results in more accurate curves.

Contrary to the above mentioned definition of a single value (gamma) as the exponent of a power function, our method uses several elementary measurements with different parameters but with the same basic process for all three

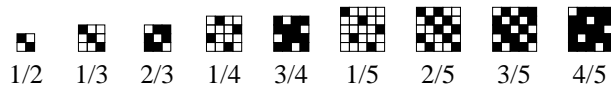


Figure 1: Dither patterns for the “chess board”

color channels.

Measurements are to be performed for the 3 color channels separately, so that series of observations will be performed on dithered patches embedded in homogeneous background color. During a single measurement, the user tunes the background luminance level so that visual difference shall disappear between the central dithered area and the homogeneous background. Such a measurement step is completed when a background luminance level is accepted. Note that human perception is most sensitive when the luminance of the background colour is approximately equal to the average luminance of the dithered patch, which is just our case. It is similar to the crispening effect.<sup>10</sup> A measurement input can be defined as the triplet consisting of the pattern used in the central area (see Fig. 1) and its bright and dark luminance levels, and it results in a background luminance level. Practically, only a ratio between the number of bright and dark pixels is taken into consideration.

The “next triplet”, i.e., the “next measurement”, can be defined depending on, or independent of, the results of the already existing data. To complete the process, a data analysis and optimization step computes the curve in question, which can be specified by combining certain criteria like the minimization of the differences at each point or the overall smoothness of the resulting curve.

The overall evaluation of the measurements as well as the definition of their series can be controlled by different criteria such as minimizing the number of measurements, having a fixed procedure or gaining the most exact characterization of the display possible.

## 4. HUMAN VISION ASPECTS

### 4.1. Spatial Vision

Based only on human observation, our method cannot determine the absolute contrast, i.e., the ratio of max/min output luminance levels,  $L_c(255)/L_c(0)$ . With other “calibration applets” based on just-noticeable contrast using visually equidistant gray series, we can estimate a rough value for this contrast, but this is not the topic of this paper. We recommend a practical approach to calibrate this singular free parameter after applying the mathematical model described below. The practical range of the perceptual contrast of a typical monitor is 20...100; for a high quality display under special conditions, it can be also be more than

1000, while in the worst case of common CRTs it can fall also below 10 in case of a bright environment.

Why do we speak about perceptual difference? According to the color appearance models<sup>10</sup> we have to take at least the ambient luminance level into consideration for a CRT or LCD screen in commonplace environmental settings (office, etc.). Also the rough gamma correction LUT approaches, using a simple gamma value have to be changed for a dark, dim or average adaptation luminance level. Namely, the real gamma value has to be divided approximately by 1.6...1.7, 1.2 and 1, respectively. Thereby we perceive linear contrast e.g. in darkness only as 60 instead of 1000.

But this perceptual aspect of visible contrast influences only the accuracy of measurement, not the original “pre-perceptive” gamma curve itself, which is not depending on ambient luminance level, only on the display. We will compute objective emission values of a display for different input signals and for the 3 basic colors. Of course, if we use an arbitrarily measured or human vision based gamma curve to generate color images, we have to take additionally the influence of ambient light and its effect for the gamma LUT into consideration. But this task would be already a problem of correct color appearance or image reproduction.

The core of the comparison process is a single comparison of an area uniformly filled by a single color against another area, the “chess board”, filled by a repeating pixel pattern. The pattern contains only two definite colors or luminance levels (*low* and *high*), assembled according to a simple, fixed ratio, e.g. 1/2 or some other rational number with small denominator (Figure 1). The homogeneous area encloses the chess-board, acting as background with luminance *back*.

To be usable, the patterns which are to be observed must appear homogeneous, otherwise the inhomogeneity disturbs our observation of its equivalence with a homogeneous area. The spatial vision models give us quantitative values for the required viewing distance.

The *contrast sensitivity function* (Figure 2) and the appropriate *visible contrast threshold function* (Figure 3) describe the spatial behavior of our vision.<sup>11-14</sup> This figure corresponds to the case of the highest luminance sensitivity. Color vision models<sup>9</sup> are working with the opposite color channel models, using achromatic, red-green and yellow-blue channels following the preprocessing in the retina. The color channels are less sensitive, that means that we do not see really sharp in particular in yellow-blue channels, we have only a feeling of sharpness thanks to the achromatic channel and the huge number of rods in the retina.

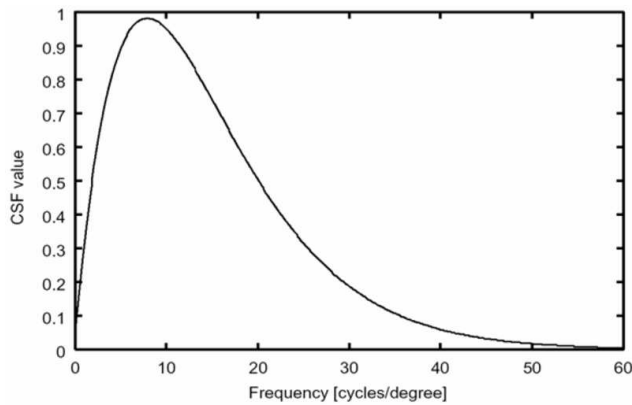


Figure 2: Contrast Sensitivity Function (CSF)

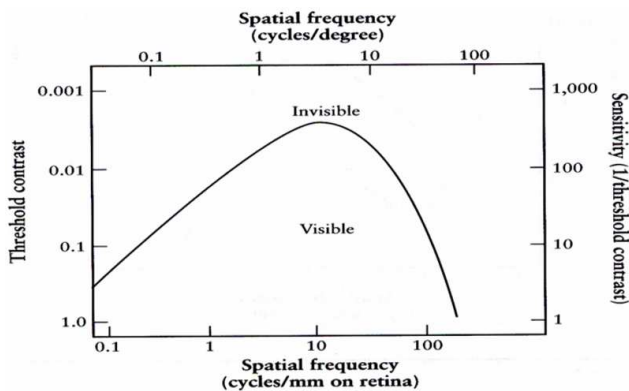


Figure 3: Visible Contrast Threshold

In order to avoid a complicated model and also to ensure the highest accuracy we selected the curve of the highly sensitive achromatic channel for all colors.

Image detail can be measured in *cycles/degree* (*cpd*). In case of a sinusoidal luminance pattern it is one cycle of the sinus function. At a discretized approach it is a pair of a darker and a brighter pixel. We do not see details in images with spatial resolution of more than approximately  $60\text{ cpd}$ , but a homogeneous field for arbitrary contrast. For opponent color channels we reach this visibility at  $30 \dots 40\text{ cpd}$ .

In our case the different “superpixel patters” of Figure 1 represent different *cpd* resolutions. The required view distance can be estimated with  $5 \times 5$  pixels representing the largest patterns of Figure 1.  $60\text{ cpd}$  of Figure 2 and a pattern with  $5 \times 5$  pixels will become invisible for distances farther than about 5m. According to our experience, a distance of 3-4 m is always enough, and for low contrast patterns even 1 m is usually sufficient. Nevertheless, the widely used 40-50 cm distance of gamma calibration tools is not enough for large patterns.

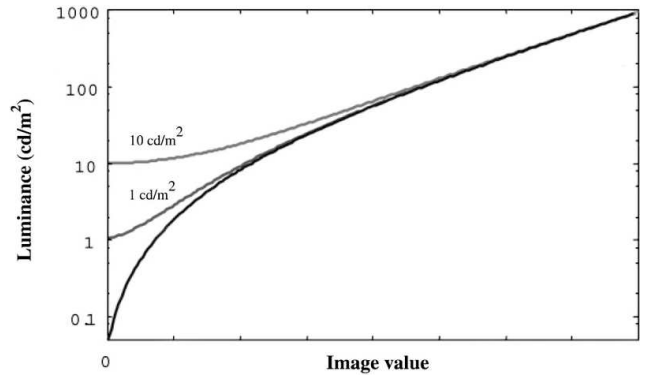


Figure 4. Effect of ambient illumination of  $1\text{cd/m}^2$  and  $10\text{cd/m}^2$  vs. dark ambient on image value.

## 4.2. Displays and Contrast

As mentioned above, CRT, but also LCD have a maximum contrast, which is the ratio of luminance values emitted by the highest and the lowest video signal, e.g., 255 and 0, and is finite instead of infinite, caused by the display technology and also unwanted ambient effects.

A CRT always has unwanted electron backscattering. Reflectivity or BRDF of the surface of a display is always positive. Reflections can be represented, e.g., by addition of a specular and a diffuse component. The specular component is often a visible mirror image of the ambient room.

All of these unwanted effects reduce the contrast of a display to a certain finite value. The answer of a zero input video signal will be a positive value at each color channel. We have not mentioned cross effects of color channels as another unwanted factor, e.g., the excitation of phosphors by the “wrong” cathode ray.

Fig. 4 illustrates how the shape of the display response curve changes depending on ambient luminance level.

Unwanted emission despite zero signals at the different color channels decrease the displayable gamut. Contrasts of color channels influenced by these unwanted emission can be roughly estimated also by visual observation. Special calibration test images can be used for this observation, which contain equidistant gray, red, green and blue scales. Linear contrast values of the color channels can be estimated from the distinguishable dark patches of these scales.

These contrast values can complete results of gamma functions obtained also by human observation. Although contrast estimation is less accurate than the gamma function, fortunately it is sufficient for tone mapping, because the visual quality of displayed images match the accuracy of measurements based on human observation.

## 5. SETUP OF MEASUREMENTS

### 5.1. Relative measurement by human observation

Our method deals with isolated color channels, therefore the computation is applied only on one-dimensional problems, independent from each other. We ignore cross effects between the  $r$ ,  $g$  and  $b$  color channels. In fact, for CRTs the cross effect is 1-3%. In lack of other colorimetric information we assume that the CIE xy chromaticity values are according to the sRGB recommendation, the ITU-R BT.709 standard:  $r = (0.64, 0.33)$ ,  $g = (0.30, 0.60)$  and  $b = (0.15, 0.06)$ , and white point  $D65 = (0.3127, 0.3290)$ . As mentioned in the previous section, only visual comparisons are used as relative inputs, so the result is also a relative function which describes the relative luminance values of the independent channels, from 0 to 255. The relative luminance values  $l_r, l_g, l_b$  can be converted to absolute values  $L_r, L_g, L_b$  by

$$L_c(val) = L_c(0) + \frac{L_c(255) - L_c(0)}{l_c(255) - l_c(0)} \cdot l_c(val) \quad (1)$$

where  $c = r, g, b$  and  $val = 0 \dots 255$ . For the sake of simplicity we will work with  $l_c(0) = 0$  and  $l_c(255) = 1$ , so equation (1) is simplified to

$$L_c(val) = L_c(0) + (L_c(255) - L_c(0)) \cdot l_c(val) \quad (2)$$

The measured and computed values  $l_c(val)$  are independent of the values  $L_c(0)$  and  $L_c(255)$ . However, these absolute values and their ratio defining the absolute contrast range of the display device in question can be interesting for the overall appearance, but finding them is outside the scope of this method.

Showing a dither pattern realized by luminances  $low$  and  $high$ , the observer is requested to tune the third, homogeneous luminance level ( $back$ ) until the luminance difference disappears. Now we have 3 values ( $low$ ,  $high$ ,  $back$ ), and a ratio of the number of the  $low$  luminance level pixels within the pattern,  $ratio = \frac{N_{low}}{N_{low} + N_{high}}$ . The following approximation can be written for the absolute luminance values

$$L_c(back) \approx L_c(low) \cdot ratio + L_c(high) \cdot (1 - ratio) \quad (3)$$

for channels  $c = r, g, b$ . Using equation (1):

$$\begin{aligned} L_c(0) + Q \cdot l_c(back) &\approx \\ &\approx (L_c(0) + Q \cdot l_c(low)) \cdot ratio \\ &+ (L_c(0) + Q \cdot l_c(high)) \cdot (1 - ratio) \end{aligned} \quad (4)$$

with  $Q = \frac{L_c(255) - L_c(0)}{l_c(255) - l_c(0)}$ , and reordering

$$l_c(back) \approx l_c(low) \cdot ratio + l_c(high) \cdot (1 - ratio) \quad (5)$$

shows the independence of the measurements from values of  $L_c(0)$  and  $L_c(255)$  and also that  $l_c(0)$  and  $l_c(255)$  can be predefined arbitrarily:

$$l_c(0) = 0 \quad \text{and} \quad l_c(255) = 1 \quad (6)$$

A single measurement gives a single  $back$  value for the measurement input triplet ( $low, high, ratio$ );  $back$ ,  $low$  and  $high$  are bytes, and  $ratio$  is a simple rational number. The goal is to define the function  $f = l_c$ , that is, 256 separate values for inputs  $0 \dots 255$ .

We now face a practical and two mathematical problems

1. Conditions of the measurements have to be defined
2. Having a list of measurements, a curve is to be defined
3. Having a list of measurements and perhaps a preliminary curve defined by them, either define the next measurement's setup, and/or recommend to stop the process

The first question belongs to the setup problem, while the last is concerning partly also to the setup or preparation and partly to a special control problem. So first these two points will be explained in detail, then the mathematical core of the method will be presented in section 6.

### 5.2. The stop criterium

The measurement shall be accurate and at the same time the number of the measurements should be minimized to achieve good results with low effort for the user. These requirements lead to another optimization problem: defining the next measurement step in an optimal way, and notifying the user when reliability has reached a certain limit.

The main problem for generating the input for the next measurement is that the expected effect of a next measurement on the reliability of the function should be evaluated depending on the behavior of the still unknown function itself. How could this function behave between the points (byte values of colors) where we already have some information, and how can we rely on the previous measurements? These questions are connected to the principles of the optimization method. We use some heuristics relying on experiments, in order to get a compromise between simplicity and accuracy, where the observer can overrule the recommendation to stop.

A measurement means a definition of a value by the user, where the  $low$  and  $high$  values used in the pattern are given as well as their mixing  $ratio$ , represented by the pattern. Then the  $low$ ,  $high$  and also the resulting colour value would be taken as support points. The method restricts the possible

low and high points for the new measurement only on the set of already used support points. The central players of the process are the *Reliability* and its opposite and reciprocal, the *Uncertainty*. The simple rule is that the *Reliabilities* can be added, its meaning and usability is of course depending on the definition of the *Reliability*, which is following here on.

The process defining the next measurement consists of three successive steps, using functions  $f_1$ ,  $f_2$  and  $f_3$  balancing the process properly:

1. Two kinds or layers of *Reliabilities* are defined. First the reliabilities of existing support points are defined as the sum of all the *Reliabilities* coming from the measurements in which these points appear (as *low*, *high*, or result point). Obviously, the points 0 and 255 are assigned absolute *Reliabilities* (i.e., zero *Uncertainty*). An individual measurement's *Uncertainty* is equal to the error of its condition (10) to be minimized, multiplied by the effect of the pattern's non-homogeneity and the contrast sensitivity belonging to the luminance to be set.

$$\begin{aligned} Unc_{supp}(j) &= |M(j)| \\ &\times f_1(\text{nonhomogeneity}(\text{pattern}(j))) \\ &\times f_2(\text{contrast}(\text{sensitivity}(\text{luminance}(j)))) \end{aligned}$$

2. Then, each point is assigned the other kind of *Reliability* as a sum of the 'effects' of the support points on the regular point in question. These 'effects' are also some sort of *Reliabilities*, but *Uncertainties* are directly computed by adding the *Uncertainty* of the support point in question and another *Uncertainty* which is a square function of the distance between the two points. This value characterizes an existing set of measurements, so when their compound value, actually their maximum, reaches a certain threshold, the process is suggested to be terminated.

$$Unc_{reg}(i) = \frac{1}{\sum_{j \in \mathcal{M}, i @_{Mj}} \frac{1}{Unc_{rel}(i,j)}} \quad (7)$$

where  $Unc_{rel}(i, j) = Unc_{supp}(j) + f_3((i - j)^2)$ .

In Eq. (7),  $i @_{Mj}$  means no  $j' \in \mathcal{M}$  lies between  $i$  and  $j$ . Of course there cannot be more than 2 such  $j$  indices for any given  $i$ .

3. Finally all the possible triplets (*low*, *high*, *ratio*) are evaluated so that the estimated effect of the triplet is added to the *Reliabilities* of the ordinary points and then their maximum values are compared in order to select the best triplet. In order to estimate the effect of

an incomplete measurement series on a regular point, which is embodied in an additional *Reliability*, its correspondent *Uncertainty* is computed as a product of the potential  $q$  quality of the measurement, and a sum of the *Uncertainties* of the given low and high points and a distance dependent component.  $q$  quality factor is estimated from the pattern and the luminance levels playing taking part in it.

$$\begin{aligned} Effect_{estim}(i, j_1, j_2, ratio) &= \\ q(j_1, j_2, ratio) &\times \left( f_3((i - j_3)^2) + f_3((i - j_x)^2) \right. \\ &\left. + \frac{1}{Unc_{supp}(j_1)} + \frac{1}{Unc_{supp}(j_2)} \right) \quad (8) \end{aligned}$$

In Eq. (8)  $j_3$  is the estimated result of the new measurement following the current solution, neighbours of  $i$  in  $\mathcal{M}$  are  $j_1$  and  $j_2$ , that is  $i @_{Mj_1}$  and  $i @_{Mj_2}$ , so  $i @_{M \cup \{j_3\}} j_3$  shall be fulfilled anyway.  $x$  is identical to 1 or 2 depending on the position of  $i$ , which lies between  $j_3$  and  $j_x$ , or by other words  $i @_{M \cup \{j_3\}} j_x$  is satisfied.

The main idea of the process is that the estimated effect of a tentative measurement is computed as sum of changes of *Reliabilities* of regular points as in Eq. (7), and their estimated *Uncertainties* are computed also as in Eq. (7). Their expected uncertainties are preestimated in such way, and then the overall change as their sum, which drives selection of the new triplet needed to minimize the predicted overall uncertainties.

## 6. THE MATHEMATICAL PROBLEM

Given a list of measurements with input controlled by a predefined or automatically derived series, a function is to be defined corresponding to them.

First of all properties have to be defined which the function has to fulfill. A set of approximations as in eq. (5) is obviously not enough to give a definition of the function, simply because of the degree of its freedom. In order to reduce it, the resulting function should fulfill certain additional conditions, defining its general behaviour. These conditions can restrict the set of the possible functions to a certain class and/or give additional criteria corresponding to the likely nature of the result.

We worked out two different approaches in this paper, both of them assume the function being smooth or, more exactly, having a small second derivative. The first approach

works in the original coordinate system while the second one translates its axes, so the problem is defined in another coordinate system.

### 6.1. Problem in linear coordinate system

First, we restrict the domain of the function only to the possible inputs, i.e., integers from 0 to 255. The second derivative on a domain consisting of a discrete set of points can be written in the form of finite differences:

$$S(i) = f(i+1) + f(i-1) - 2f(i) \quad (i = 1 \dots 254) \quad (9)$$

We transform the approximations (5) for the  $N$  measurements ( $j = 1 \dots N$ ):

$$\begin{aligned} M(j) = & f(\text{low}_j) \cdot \text{ratio}_j \\ & + f(\text{high}_j) \cdot (1 - \text{ratio}_j) \\ & - f(\text{back}_j) \end{aligned} \quad (10)$$

Now we have two separate sets of conditions which impact the shape of the desired function: the *smoothness conditions*  $S(i)$  and the *measurement conditions*  $M(j)$ . The latter can be used also by two different modes. One is to take them as hard constraints, i.e.,  $M(j) = 0$ , the other is to minimize them together with the other conditions.

It can be argued that there is no exact measurement, at least because setting their values should give an exact real number, but the measurements can be chosen only from a discrete set of numbers actually. On the other hand, the user can introduce errors by his/her estimation as well, so in addition there can even be more or less contradictory conditions. The problem is solved as compound minimum problem of the smoothness and measurement conditions, and their importances are accounted for by weight factors  $s_i$  and  $m_j$ . The optimal result would have all of the expressions  $S(i)$  and  $M(j)$  equal to zero ( $i = 1 \dots 254, j = 1 \dots N$ ), so we have to minimize the expression

$$F = \sum_{i=1}^{254} s_i \cdot S(i)^2 + \sum_{j=1}^N m_j \cdot M(j)^2, \quad (11)$$

where by (6)  $f(0)$  and  $f(1)$  are constant, and  $F$  is a 254-variable function. As a result we get a smooth function conforming well to the measurements, as expected. All in all there are  $256 - 2 + N$  minimum criteria and 2 constraints despite the original 256 variables. These obviously cannot be made zero at the same time, so the solution will be a compromise depending on the weights  $s_i, m_j$  and the content of  $M(j)$ .

There are several efficient methods to solve the quadratic minimum problem, two of which are mentioned here. One

is solving  $F$  using a system of 254 linear equations with a sparse matrix. The other is directly solving the minimum problem by an appropriate descent method.<sup>15</sup>

We have chosen a conjugate gradient method, which is a modification of the steepest descent method, and is in our case faster by one magnitude than the ordinary steepest descent (gradient) method. An optimization needs 20–50 elementary steps, each of them consisting of an evaluation of  $F$  and its derivative.

A problem not yet mentioned is the definition of the weights  $s_i$  and  $m_j$ . Considering the equal importance of the different smoothness conditions and the different measurement conditions respectively, we assume all  $s_i = s$  and  $m_j = m$ . Multiplying the whole problem by a constant we can set  $s = 1$ , so it is enough to define  $m$ .

To define this value, let us consider the overall behaviour of the solution. Optimizing the total smoothness leads to spline-like curves, where the magnitude of the average value of the second derivative is  $1/255^2$ , and its overall sum is  $1/255$ . Minimizing its distribution, we get a slowly changing second derivative, i.e., a curve behaving locally similar to a polynomial of about 3rd degree. A sudden jump causes an  $O(1)$  constant anomaly, so if the magnitude of  $m$  is between 1 and  $1/255$ , we have a locally and also globally well behaving function. Of course this value can be modified or tuned further by demand.<sup>1</sup> describes this method in more details.

### 6.2. Solution in log-log coordinate system

Considering that the optimization introduced above tries to reach maximal smoothness which leads to functions behaving locally like 3rd order polynomials, and also considering that the display characteristics used to be approximately a power function which differs from this one, another approach was investigated.

We transform the coordinate system, the domain as well as the range, to a log-log scale so that

$$\log y = g(\log(x/255)) \quad \text{where} \quad y = f(x) \quad (12)$$

Linear functions of this coordinate system are corresponding to  $c \cdot x^p$  power functions in the original system, and the minimum problem results in the possibly smoothest functions, that is, the functions most similar to linear functions, so this transformation looks like the appropriate way to get power-like functions in the original coordinate system.

#### 6.2.1. Problem statement

Taking the new variables, the coefficients  $s_i$  of the smoothness conditions  $S(i)$  will change, and also the measurement conditions  $M(j)$  shall be rewritten with the exp functions

of the new variables, since the formula is applicable to the original values. All in all another minimum problem is to be solved by the conditions:

$$\begin{aligned} S(i) &= s_{i,i+1} \cdot f(\log(i+1)) \\ &+ s_{i,i-1} \cdot f(\log(i-1)) \\ &+ s_{i,i} \cdot f(\log i) \end{aligned} \quad (13)$$

with  $(i = 2 \dots 254)$ , where, from the finite difference formulas,

$$\begin{aligned} s_{i,i+1} &= \frac{2}{(\log(i+1) - \log i) \cdot (\log(i+1) - \log(i-1))} \\ s_{i,i-1} &= \frac{2}{(\log i - \log(i-1)) \cdot (\log(i+1) - \log(i-1))} \\ s_{i,i} &= \frac{2}{(\log(i+1) - \log i) \cdot (\log i - \log(i-1))} \end{aligned}$$

$S(1)$  has been omitted, since the variable  $y_0 = -\infty$  is omitted anyway, which means that no smoothness condition could be defined using this variable. It is neither a real problem in this coordinate system, nor in its retransformed linear version. Behind this little but strange point can be found the characteristics of the power-like functions. All in all the problem is simpler by one variable, which gives one more degree of freedom. This additional freedom lies behind the phenomenon that an arbitrary power function can be defined by one measurement, as it can be done by e.g. the simple gamma applets.

The measurement conditions change into

$$\begin{aligned} M(j) &= \exp(f(\log(\text{low}_j))) \cdot \text{ratio}_j \\ &+ \exp(f(\log(\text{high}_j))) \cdot (1 - \text{ratio}_j) \\ &- \exp(f(\log(\text{back}_j))) \end{aligned} \quad (14)$$

where  $(j = 1 \dots N)$ .

The minimum problem (11) can be written with the expressions above in order to obtain the form of the directly transformed measurement conditions, but in this case we face two new problems. One is the convexity of the measurement conditions, corresponding to (14). The square of these expressions will not be convex, which leads to some algorithmical difficulties, especially by considering the original dimensionality of the problem, which has 256 variables.

The other problem is much more crucial. Conditions in (13) can be weighted arbitrarily, their weights expressing their individual importance. If different weights are used, they could distort the overall smoothness depending how unbalanced they are, that is, the function would be smoother in one part and less smooth in another. In a drastical case

it destroys the original expectation, that is, the function will not behave like a power function, which was the argument to apply the axis transformation in the first place. Unfortunately, if the weights are equal, the magnitudes of the coefficients will be enormously different, their maximum value will be

$$\begin{aligned} s_{2,2}/s_{254,254} &= \frac{(\log 255 - \log 254) \cdot (\log 254 - \log 253)}{(\log 3 - \log 2) \cdot (\log 2 - \log 1)} \\ &\approx \frac{1/254.5 \cdot 1/253.5}{1/2.5 \cdot 1/1.5} \\ &\approx 2^{-14} \end{aligned}$$

which leads to unsolvable numerical instabilities when minimizing them, especially by taking their squares the magnitude will be squared as well ( $2^{-28}$ ). So it is obvious that the log-log transformation with the form of the minimum problem cannot work because of numerical reasons.

### 6.2.2. A two pass solution for the log-log problem

There are two sets of points on the transformed independent axis, which play a special role in our problems.  $\mathcal{S}$  consists of all points playing any role in the smoothness conditions (13). The other, set  $\mathcal{M}$ , is a subset of the 256-element set  $\mathcal{S}$ .  $\mathcal{M}$  consists of all points playing any role in the measurement conditions (14). Let us consider the next smoothness conditions written for the triplets of the sparser set  $\mathcal{M}$ .

$$\begin{aligned} S_{\mathcal{M}}(k) &= s_{\mathcal{M}k,j} \cdot f(\log j) \\ &+ s_{\mathcal{M}k,k} \cdot f(\log k) \\ &+ s_{\mathcal{M}k,l} \cdot f(\log l) \end{aligned} \quad (15)$$

where  $(\log j, \log k, \log l)$  is a successive triplet of  $\mathcal{M}$ , and from the finite difference formulas

$$\begin{aligned} s_{\mathcal{M}k,j} &= \frac{2}{(\log k - \log j) \cdot (\log l - \log j)} \\ s_{\mathcal{M}k,k} &= \frac{2}{(\log l - \log k) \cdot (\log k - \log j)} \\ s_{\mathcal{M}k,l} &= \frac{2}{(\log l - \log k) \cdot (\log l - \log j)} \end{aligned}$$

It can be seen that solving the minimum problem of the smoothness criteria (13) over  $\mathcal{S}$ , and taking its values only over its subset  $\mathcal{M}$ , this restricted solution can be taken as a good quasi minimum for the problem of the smoothness criteria (15) defined directly over  $\mathcal{M}$ . And vice versa, solving the problem just over  $\mathcal{M}$  and fixing its results, then solving the previous problem by these constraints, the obtained overall solution will also be a good quasi minimum for the non-constrained original problem.

It can be seen as well that the situation is similar if other, new conditions are added to the original problem. As its



main example, the measurement conditions (14) can be added to both versions of the smoothness problem above, that is, to the original condition system (13), and also to its extension by (15) over  $\mathcal{M}$ . their solutions will approximate each other well.

Let us suppose that the minimum problem of this extended condition system has already been solved. Now, let  $f(\log j)$  values be fixed for each  $\log j$  from  $\mathcal{M}$ , and the minimum problem be solved again with this additional constraints. Let us recognise that the solution of this second problem cannot differ from the non-constrained one, since this solution could be a better solution also for the first problem, contradicting its minimality! Therefore the same solution must be obtained for both of them.

After this recognition let us consider the series of problems:

- A Smoothness condition (13) over  $\mathcal{S}$   
Measurement condition (14) over  $\mathcal{M}$   
no constraints
- B Smoothness condition (13) over  $\mathcal{S}$   
Smoothness condition (15) over  $\mathcal{M}$   
Measurement condition (14) over  $\mathcal{M}$   
no constraints
- C Smoothness condition (15) over  $\mathcal{M}$   
Measurement condition (14) over  $\mathcal{M}$   
no constraints
- D Smoothness condition (13) over  $\mathcal{S}$   
Smoothness condition (15) over  $\mathcal{M}$   
Measurement condition (14) over  $\mathcal{M}$   
constraints of solution C over  $\mathcal{M}$
- E Smoothness condition (13) over  $\mathcal{S}$   
constraints of solution C over  $\mathcal{M}$

A is the original problem, it gives an equivalent solution with B, and B gives also an equivalent solution over  $\mathcal{M}$  with C. D gives also an equivalent solution with C, so that they are equivalent over  $\mathcal{M}$ . Finally, as proven above, the solution of E is identical to the solution of D.

It implies a two pass method resulting in a well-approximating solution of the original minimum problem. First, problem C shall be solved which is a minimum problem in few variables, containing also slightly non-convex components, but its technical problems can be overcome by a moderate effort.

Second, problem E shall be solved, which includes all the numerical instabilities, but in this form it is a pure smoothness problem at the same time. This problem is identical with defining a spline across a finite set of given points. Independently from the number of given points, it leads to a 2-dimensional optimum problem defining interval-wise 3rd order polynomials. These polynomials are defined by their

coefficients instead of a set of points of the domain that would promise numerical difficulties.

For the last step a reconversion shall be performed, transforming the  $\log y_i = f(\log i)$  that is  $\log y_i = f_i$  values back to  $y_i = \exp f_i$ . The obtained curve is a piecewise transition of power-like base functions, so it coalesces the behaviour of the power function and changes one power function to another in the smoothest way in sense of the log-log coordinate system.

### 6.2.3. Discussion of the methods

The original linear problem is a 256 variable,  $256 + N$  condition minimum problem, with  $N \ll 256$ . It is solved practically by a steepest descent method which is a stable and fast algorithm.

The behaviour of the approximation highly depends on the balance between the smoothness conditions and measurement conditions. The solution of a compound quadratic optimum problem has been selected in order to compensate for the possible errors of the measurements, which are actually human perception based, but errors would come into account in any sort of measurements. In case of absolute weights of measurements, that is, assuming their absolute accuracy, a spline problem is implied, similar to the second pass of the two pass method. This case shows the nature of the obtained curve better, which is an interval-wise 3rd-order polynomial in a general case too. Interval-wise 3-order polynomials of a spline realize the smoothest transition between their supporting points or, in other words, the smoothest transformation of linear segments in order to link them together smoothly. This character is watermarked also in a softer weighting similar to the weights given in the linear coordinate system.

The second method using the log-log coordinate transformation also keeps the smoothness, since the transformation of the coordinate system has been introduced just in order to conform the preferred shapes to the displays. It is only a practical question, but just because it is a practical question, it can result in a better practical solution.

After the coordinate transformation, a two pass separation of the algorithm has been applied in order to overcome the numerical instabilities. In fact this method could have been applied also for the original lin-lin coordinate system, and this could reveal that the solution is also a spline in any case because of the equality of the solutions of the constrained and non-constrained problems as it has been proven in the previous section. However, it is clear that the quasi-equivalences of the previous section are not strict equivalences, so this solution is not the theoretical optimum.

Nevertheless, applying this two-step process to our case we have decreased the problem of the weights, which is addressed in both steps to the same  $\mathcal{M}$  domain where they are naturally commensurable. On the other hand the second pass is also obvious, and both are fast. In addition one can introduce additional freedom into the spline computation, allowing the values over  $\mathcal{M}$  not to be exact, but the function being smoother. This leads to a similar question as it has appeared with the lin-lin method's weighting, but in a much simpler form and in a faster algorithm. All in all, it means that the theoretical exactness has been sacrificed for getting a faster, simpler result which is still extendable to a similarly flexible solution.

In addition it can be seen that the applied methods and forms are secondary compared to the original problem, which is not exactly formalized, but a matter of decision.

## 7. RESULTS

To evaluate the ability of our method to capture the real TRC curve of a monitor, we compare it with results from direct measurements of the TRC curve using a GretagMacbeth Spectrolino spectrophotometer. In order to demonstrate the general usability of the method, the experiments are performed on two different monitor technologies: Cathode Ray Tube (CRT, Nokia 446Pro) and Liquid Crystal Display (LCD, EIZO FlexScan L985EX). The Spectrolino measured 52 color values per channel, with color values 0 to 255 increasing in equal steps of 5.

Figure 5 shows the results for the CRT Monitor in the left half, for the LCD monitor in the right half, for all three channels. In this figure the TRCs obtained as result of the spectrophotometer measurements are taken as reference (solid line). The TRCs obtained as result of our method (dashed line) practically are able to reproduce the reference curve, or a close approximation. This is valid for both monitors used in the experiments. Note the jaggy in the Spectrolino CRT curves, caused by the different colour depth setting of 16-bit.

## 8. CONCLUSION AND FUTURE WORK

A significant improvement for our model for gamma display characterization<sup>1</sup> has been presented in this paper which is able to overcome the limitations of the previous version, in particular, solving the problem of instability and unequal weighting. This method is simplified and made more reliable by changing the simple smoothness principle to a log-log scale based one. A trade-off between the smoothness and observation conditions is analysed, improving the reliability of the results of the method.

The flexibility of the model allows it to be used in many applications without the necessity to use a spectrophotometer. Also a fast and simple (re)characterization is possible, just using an interactive process with the end user. The main benefit of this approach is that it is a cheap and simple solution, either to define a relative TRC of a display or to verify a given one.

There are a couple of open possibilities in the method. On one hand its usability could be improved by taking information on the absolute values, that is, about the contrast value and the smallest (or the largest) value of the absolute luminance, either staying at the human perception based input or using other methods.

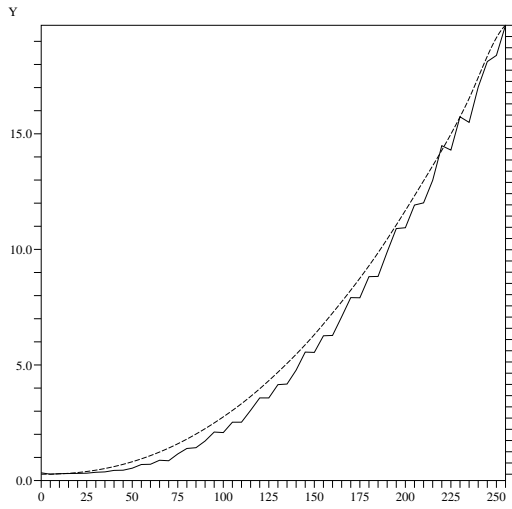
On the other hand, it should be possible to solve the phenomenon of the cross effect between the colour channels, and the method may be extended for a multidisplay system application.

## ACKNOWLEDGMENTS

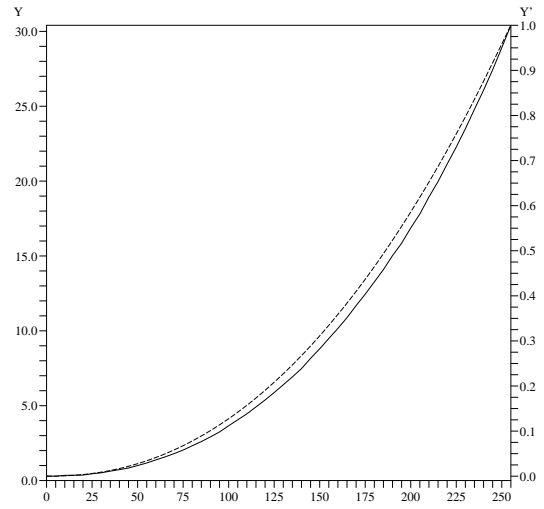
This work was partially supported by the European Union within the scope of the RealReflect project IST-2001-34744, "Realtime Visualization of Complex Reflectance Behaviour in Virtual Prototyping" and by the Spanish Government by project number TIC2001-2416-C03-01.

## REFERENCES

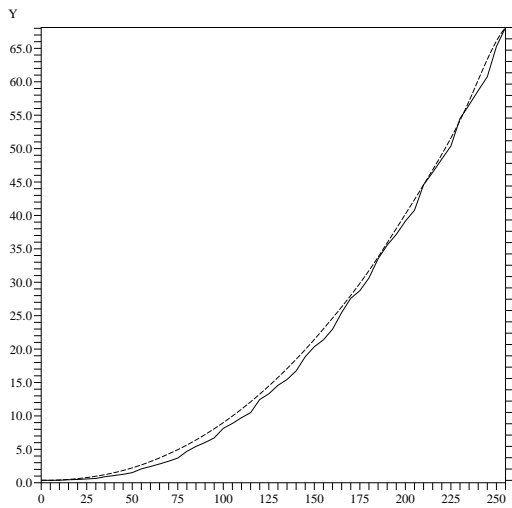
1. A. Neumann, A. Artusi, G. Zotti, L. Neumann, and W. Purgathofer, "An Interactive Perception Based Model for Characterization of Display Devices," in *Proc. SPIE Conf. Color Imaging IX: Processing, Hardcopy, and Applications IX*, pp. –, January 2004.
2. R. S. Berns, R. J. Motta, and M. E. Gorzynski, "CRT colorimetry – Part I: Theory and practice," *Color. Res. Appl.* **8**, pp. 299–314, Oct. 1993.
3. R. S. Berns, R. J. Motta, and M. E. Gorzynski, "CRT colorimetry – Part II: Metrology," *Color. Res. Appl.* **8**, pp. 315–325, Oct. 1993.
4. M. D. Fairchild and D. Wyble, "Colorimetric characterization of the apple studio display (flat panel LCD)," tech. rep., Munsell Color Science Lab., Rochester Institute of Technology, Rochester, NY, July 1998.
5. E. Day, "Colorimetric Characterization of a Computer Controlled (SGI) CRT Display," tech. rep., Munsell Color Science Lab., Rochester Institute of Technology, Rochester, NY, April 2002.
6. N. Kato and T. Deguchi, "Reconsideration of CRT Monitor Characteristics," in *Color Imaging: Device Independent Color, Color Hardcopy, and Graphics Arts III*, pp. 33–40, 1998.



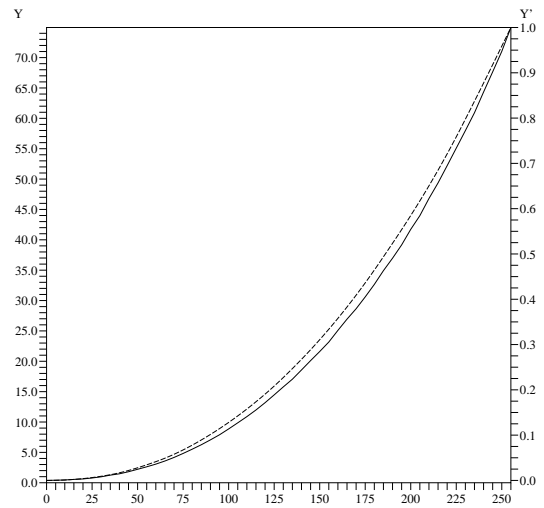
(a) CRT, red channel



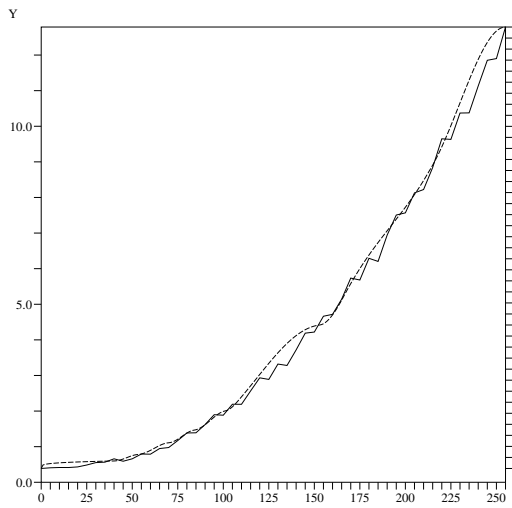
(d) LCD, red channel



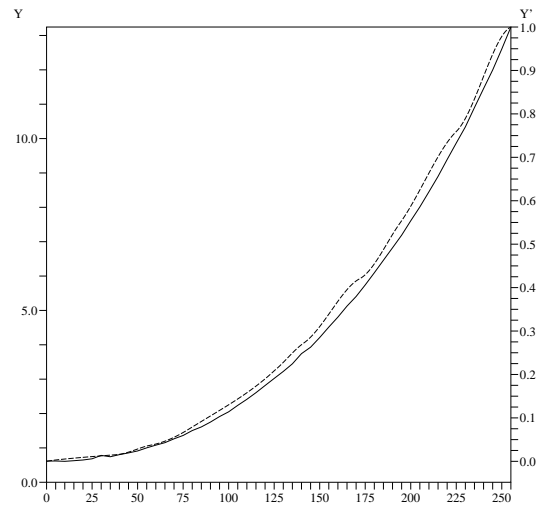
(b) CRT, green channel



(e) LCD, green channel



(c) CRT, blue channel



(f) LCD, blue channel

**Figure 5:** Comparison of TRC acquired with new method (dashed line) and Spectrolino measurements (solid line)

7. D. L. Post and C. S. Calhoun, "An Evaluation of Methods for Producing Desired Colors on CRT Monitors," *Color. Res. Appl.* **14**, Aug. 1989.
8. N. Tamura, N. Tsumura, and Y. Miyake, "Masking Model for Accurate Colorimetric Characterization of LCD," *Journal of the SID* **11**(2), pp. 1–7, 2003.
9. R. Hunt, ed., *Measuring Color*, Ellis Horwood, 2nd ed., 1992.
10. M. D. Fairchild, ed., *Color Appearance Models*, Addison-Wesley, 1998.
11. R. L. D. Valois and K. K. D. Valois, *Spatial Vision*, no. 14 in Oxford Psychology Series, Oxford Univ Press, September 1990.
12. J. Schirillo and S. Shevell, "Brightness Contrast from Inhomogeneous Surrounds," *Vision Research* **36**(12), pp. 1783–1796, 1996.
13. K. Tiippana and R. Nasanen, "Spatial-frequency bandwidth of perceived contrast," *Vision Research* **39**(20), pp. 3399–3403, 1999.
14. A. Watson and C. Ramirez, "A Standard Observer for Spatial Vision." [http://vision.arc.nasa.gov/publications/watson\\_arvo2000.pdf](http://vision.arc.nasa.gov/publications/watson_arvo2000.pdf).
15. R. Fletcher, *Practical Methods of Optimization*, John Wiley & Sons, 2nd ed., 2001.