# **Color Clusterization using Modified HSL Space**

Gabriel Marcu and Satoshi Abe\* Array Corporation, Tokyo, Japan \* Faculty of Computer Science, University of Tokyo, Hongo, Tokyo, Japan

#### Abstract

This paper describes a clusterization procedure based on a set of new transformations for HSL color space. The clusterization is applied for color layer extraction from the large blue print documents overdrawn with color pens and pencils. The process of color layer extraction requires a segmentation and a classification of segmented area to recompose each color layer. Overlapped colors are identified based on extraction of primary and secondary clusters. The new HSL transformations are based on the modification of lightness function in a standard HSL space. The modified HSL color space remains device dependent space, but the RGB regulate domain is not transformed into a HSL regulate domain. The clusterization procedure is based on the Mahalanobis distance in new HSL space. A 3D visualization procedure is used for illustration the efficiency of the clusterization process.

## **Conventional HSV and HSL Color Spaces**

The HSL and HSV color spaces are well known in the computer color processing (i.e.[1]). The HSV color model is described by the following set of relationship (v corresponds to brightness):

$$v = \max (r,g,b);$$
  

$$s = \max (r,g,b) - \min (r,g,b);$$
 (1)  

$$h = hue (r,g,b),$$

where hue() is a hue function (with values in interval (0,1)), defined as follows:

hue(r,g,b) = h / 6, if h > 0, and hue(r,g,b) = (h + 6) / 6, if h < 0, where h is:

$$\begin{array}{ll} h = (g - b)/dh , & \text{if } r = \max (r,g,b), \\ h = 2 + (b - r)/dh , & \text{if } g = \max (r,g,b), \\ h = 4 + (r - g)/dh , & \text{if } b = \max (r,g,b). \end{array}$$

The HSL color model is described by the following set (v corresponds to lightness):

$$v = (\max (r,g,b) + \min (r,g,b))/2;$$
  

$$s = \max (r,g,b) - \min (r,g,b);$$
 (2)  

$$h = hue(r,g,b).$$

The visualization of the HSL space, that corresponds to the relationships (2), is presented in Figure 1. The representation is carried out using the method presented in [2] and extended for color gamuts, but any visualization program can be equally used. The relationships (1) and (2) correspond to the cone or double cone volumes in cylinder coordinates. The lightness and brightness are represented along the vertical axis of the cylindrical reference system. Hue is represented as the circular coordinate and the saturation is the radial coordinate.



Figure 1. HSL color spaces as double cone.

There are also references in which the HSV and HSL spaces are represented as "cylinder" and are defined as follows for HSV and HSL respectively (or an equivalent form if r,g,b values are in other interval than [0,1] ):

$$v = \max(r,g,b);$$
  

$$s = 1 - \min(r,g,b) / \max(r,g,b);$$
 (3)  

$$h = hue(r,g,b),$$

$$v = (max (r,g,b) + min (r,g,b))/2;$$

$$s = \begin{cases} \frac{\max(r, g, b) - \min(r, g, b)}{\max(r, g, b) + \min(r, g, b)} & \text{if } 0. < v < 0.5 \\ \frac{\max(r, g, b) - \min(r, g, b)}{2 - \max(r, g, b) - \min(r, g, b)} & \text{if } 0.5 < v < 1 \end{cases}$$
(4)

h = hue(r,g,b).

A detailed explanation of the cylindrical and cone shapes of these spaces can be found in [2], where the distinction between the definition (1)-(2) and (3)-(4) is analyzed.

Despite the "perceptual" attribute of these spaces, the HSL and HSV spaces defined by the equations (1~4) are strongly not uniform. The non-uniformity is in all three dimensions of the spaces. The most perceptual non-uniformity is in the lightness or brightness. As we can observe from all diagrams, in all these spaces the yellow and the blue have same lightness or brightness even if the B&W images showing only the luminance of the image split these two colors in very light and very dark gray levels. Same observation stands for colors with other hues. This phenomena is not disturbing in a application like color picker, but can conduct to errors in an application like color classification. For such application other color spaces can be used, like CIELAB of CIELUV, but the color transformation required by these spaces are more complicated and time consuming and require color calibrations.

This paper proposes new relationships of the HSL and HSV color spaces based on the replacement of the lightness and brightness definitions, in the relationships (1) and (2) with a new relationship. The new color space is referred as modified HSL (MSHL). In this paper the new space is visualized in three dimensions and its shape is compared with the old regulate HSV and HSL spaces. The MSHL space remains device dependent as the conventional HSV and HSL spaces, but the distribution of color is more suitable for some application like color clusterization.

A disadvantage of this space as well as the cone HSL and HSV color spaces is that the regulate cube domain is not converted in a regulate cylinder domain, and the variation of saturation is dependent on the variation of brightness or lightness. This may be important for a color picker where the user want to see an uniform domain for all dimensions of all the available colors. In the MSHL space, the saturation range depends not only on the lightness but also on the hue. This is important in case reversal color transformation are designed, requiring more cautions concerning the definition of the input values (l, s, h) for transformation to RGB. However if the (l,s,h) vector is resulted from direct transformation, the domain of input values is automatically verified.

#### **Modified HSL Color Space**

The HSV and HSL definition can be modified for more perceptual lightness and brightness. We proposed to replace the v relationship in  $(1)\sim(4)$  equations with the following relationship:

$$v = 0.3 r + .59 g + .11 b$$
,

the well known equation that define the luminance in television industry. The new form of HSL color space is:

$$v = 0.3 r + .59 g + .11 b,$$
  

$$s = 1 - \min(r,g,b) / \max(r,g,b) ;$$
(5)  

$$h = hue(r,g,b),$$

for "cylinder" volume and

$$v = 0.3 r + .59 g + .11 b,$$
  

$$s = \max (r,g,b) - \min (r,g,b) ; \qquad (6)$$
  

$$h = hue(r,g,b),$$

for "cone" volume.

The 3D representation of these equations is presented in Figure 2. In the next section, the new HSL color space is used for an application of color classification.



Figure 2. The 3D representation of MHSL "cone" space

#### **Color Classification in Blueprint Documents**

The clusterization procedure presented here was designed for a particular class of color images, the blue print copies of schematic diagrams. On these documents, additional diagrams are hand drawn with color pens and pencils. A small sample of a blue print image is presented inFigure 3. These documents are scanned, then processed for extraction of color layers. Next, the color layers are processed for symbol recognition and coded to store the information in CAD format. The images to be processed for color classification are usually A1 format size and are scanned with 16 dots/mm and 8 bits/color components resulting in about 390 MBytes of raw data (about  $13600 \times 9500$  pixels). For color classification, the processing time was limited to a maximum of 20 min. Due to the time requirement, the clustering methods based on processing of both the color information and the spatial information of the pixels, as proposed in [3],[4],[5],[6], cannot be used. For example, for a 390MB file image, the algorithm proposed by Hedley and Yan,<sup>3</sup> the fastest from the algorithms previously mentioned requires about 2 hours.

The clusterization method presented here uses only the color information of the original image in the form of the colormap and the histogram table. The spatial information is neglected due to the time processing limitation. As result of the clusterization procedure applied on the colormap of the original image, a classification table is derived for each color layer. Each binary file of a color layer results by passing the original image through the LUT corresponding to the classification table. A spatial filtering process is then applied to the binary output image of each color layer, to filter the transition colors and to remove partially the errors of color classification that may appear when the pixel spatial information is ignored. Reference [7] offers more details about analysis and processing method of blue print documents, briefly presented here.



Figure 3. Sample of a blueprint document

The representation in MHSL space of the 3D color histogram of the original image is given in Figure 4. For the conventional HSL and HSV color spaces the lightness or brightness information cannot be used as a feature to discriminate the colors due to the fact that for same saturation and hue, colors that naturally on a black and white image can be easily separated, in the color image have same lightness or brightness. Additionally, the HSV color spaces disperse the distribution of colors for shadow (both cylinder and cone) and for highlight (cylinder) colors that conduct to poor results of color classification in that regions. These is the reason to modify only the HSL color space for color classification and not HSV color space, even if the lightness equation can be equally applied. The definition of saturation of HSV color space conducts to poor performance of color classification.



Figure 4. 3D histogram of image from Figure 2 in MHSV space.

For clusterization of a colormap based on histogram information, an agglomerative clusterization procedure<sup>2</sup> was implemented based on the nearest neighbor criteria. The clusterization procedure described here starts assigning the histogram colors in different clusters. The number of colors selected for the colormap is not critical and all images represented with a number of colors between 150 and 250 preserve the appearance of image from Figure 3, that uses 200 colors. The initial number of clusters is equal to the number of colors in the colormap, each cluster containing a single color. The clusterization is performed iteratively. For each iteration step, the distances between all pair of clusters are computed. The pair of clusters corresponding to the minimum distance are merged into a single cluster and the current cluster number is decreased by 1. The algorithm stops at a threshold number of clusters. The distance between clusters is:

$$Dclust = D - R1 - R2, \tag{7}$$

where D, R1 and R2 are elements of the 3D histogram representation. D represents the euclidean distance between the centers of the clusters computed in color space reference system used for visualization and R1 and R2 represent clusters radius. If the cluster contains a single color, C1, the center of cluster is the center of histogram ball (x1,y1,z1) in the color space reference system and its radius is the ball radius. If the cluster is formed merging 2 or more different colors, the Mahalanobis<sup>8</sup> distance is used to define the cluster radius. In this case, the cluster radius is the Mahalanobis distance, Dm, such that, 80% of the elements of the cluster have the distance to its center smaller than Dm. The distance of one of cluster elements, **v**, to its center is

$$d = ((v - m) \cdot C^{-1} \cdot (v - m)^{t})^{1/2},$$

where **m** and **C** represent the mean vector and covariance matrix of the cluster,  $(\mathbf{v}-\mathbf{m})^t$  is the transpose of matrix  $(\mathbf{v}-\mathbf{m})$ , and  $\mathbf{C}^{-1}$  is the inverse of covariance matrix **C**. The adjustment of the clusterization parameters is discussed in details in [2].

Since the background color results in not useful clusters (at least one, usually more than one, due to the non-uniformity of background color of paper), apriori information about colors to be classified is used to specify the clusters of interest. The pen and pencil colors to be classified are known before clusterization process and this information is used to specify that the useful clusters are placed in certain regions of the color space. The cluster regions are specified as boundary limits of the vector components of color space reference system. A cluster is validated if 80% of its elements are placed inside the region apriori decided for its color. Only the validated clusters are considered in order to stop the clusterization algorithm. This procedure eliminates the clusters corresponding to the background color.

Here are emphasized the benefits of using the modified HSL color space. The definition of the apriori cluster information during the calibration procedure required by the classification process is much more natural, and the lightness information of the color can be used effectively as a feature to classify the color layers.

Figure 5a presents the result of clusterization procedure for 7 clusters identified based on principal clusters only visualized in the MHSL color space. In Figure 9b~d, same image is processed for extraction of clusters using additional information of the secondary clusters. A secondary cluster can be classified in two different classes, as they physically result by overlapping two colors. As a rule, the secondary clusters are placed always between the primary clusters. Hence, the tolerance in classification of secondary clusters is not critical due to the fact that all colors between a primary and a secondary cluster lies in the same color class. This feature enable the use of the device dependent HSL color space instead of more complicated but device independent CIELAB or CIELUV color spaces. The modification of the lightness enables us to exploit the new feature for better color classification and to avoid the use of more complex color transformation. Figure 6a,b illustrate the results of clusterization process applied on the sample image given in Figure 3 when the color layers are extracted based only on primary clusters (a) and based on primary and the secondary (b). The effect of using the secondary clusters as well as the modified HSL color space conduct to a improved results of color classification with significant consequences for the next process of symbol and line recognition.





(b) red layer (primary & sec. clusters

(d) blueprint layer (primary & sec. clusters)





(a) - blueprint layer without secondary clusters
 (b) - blueprint layer with secondary clusters
 Figure 6. Blueprint layer extracted without (a) overlapped colors and with (b) overlapped colors.

### Conclusions

This paper introduced a modified HSL color transformation that enables a more accurate color classification that the conventional HSL and HSV color transformations. Specifically, the modified HSL color transformation enables the use of lightness feature to classify the colors with better performance, while keeping the form of color transformation reasonable simple. A 3D visualization procedure was used to compare the solid volumes and the classification results for the conventional and the new introduced color spaces.

#### References

 J. D. Foley, A. vanDam, S. K. Feiner, J. F. Highes, Computer Graphics: Principle and Practice, Addison Wesley, p.590 (1993).

- G. Marcu, S. Abe, Three dimensional histogram visualization and applications, *J. of Electronic Imaging*, V4, N4, p.330 (1995).
- 3. M. Hedley, H. Yan, "Segmentation of color images using spatial and color space information", *J. Electronic Imaging*, **1**(4), p.374 (1992).
- 4. S. Tominaga, "Color classification of natural images", *Color Research and Applications*, **17**(4), p.230 (1992).
- 5. J. Liu, Y. H. Yang, "Multiresolution color image segmentation", *IEEE Trans. on PAMI*, **16**(7), p.689 (1994).
- 6. J. Wu, H. Yan, A. N. Chalmers, "Color image segmentation using fuzzy clustering and supervised learning", *J. Electronic Imaging*, **3**(4), p.397 (1994).
- G. Marcu, S. Abe, "Analysis of large documents for color classification", J. Color Science Association of Japan, 3(4), p.397 (1995).
- J. M. Jolion, P. Meer, S. Batauche, "Robust clustering with applications in computer vision", *IEEE Trans. on PAMI*, 13(8),p.791 (1991).