

Highlighted Document Image Classification

Yafei Mao^a, Yufang Sur^a, Peter Bauer^b, Todd Harris^b, Mark Shaw^b, Lixia Li^b and Jan Allebach^a

^aSchool of Electrical and Computer Engineering, Purdue University, West Lafayette, IN 47907, U.S.A.;

^bHP Inc., Boise, ID 83714, U.S.A..

Abstract

There are many existing document image classification researches, but most of them are not designed for use in constrained computer resources, like printers, or focused on documents with highlighter pen marks. To enable printers to better discriminate highlighted documents, we designed a set of features in CIE Lch(a^*b^*) space to use along with the support vector machine. The features include two gamut-based features and six low-level color features. By first identifying the highlight pixels, and then computing the distance from the highlight pixels to the boundary of the printer gamut, the gamut-based features can be obtained. The low-level color features are built upon the color distribution information of the image blocks. The best feature subset of the existing and new features is constructed by sequential forward floating selection (SFFS) feature selection. Leave-one-out cross-validation is performed on a dataset with 400 document images to evaluate the effectiveness of the classification model. The cross-validation results indicate significant improvements over the baseline highlighted document classification model.

Introduction

Multifunction printers (MFPs) are popular in home and small office. This is because they offer multiple functions, such as print, copy, and scan, for the price and size of a single device. Apart from the cost and efficiency, the most important factor that the customers care about is image quality. In this regard, different configurations of the scan/copy pipeline are embedded in the device to optimize the image quality of a particular kind of image type, such as text documents, highlighted pages, or photos. For example, the configuration designed for the text mode may increase the contrast and sharpen the edges to get clear text; and the configuration designed for the photo mode may impose a smoothing effect to reduce the noise. A common method to change image quality settings is through manual selection [1]. Users can choose the most appropriate mode from a list of predefined modes according to the content of the document, or adjust each attribute from the submenu. Such a method often requires trained users, which may not always be desirable. Therefore, it is necessary to integrate an automatic document image classification model into the printer firmware.

There are many research papers on document image classification [2–6]. However, they are not all suited for use in entry-level printers due to the memory and computational complexity restrictions of the printer firmware. To avoid such problems, Lu et al [7] developed a low-complexity algorithm using SVMs and several features to classify text, photo, and mixed documents. Xu et al [8] extended the approach by adding more features and two additional classes highlight and faded document. The features in [8] are

- *Luminance and chroma flatness scores* describing the spread of the histograms.

*Research supported by HP Inc., Boise, ID.

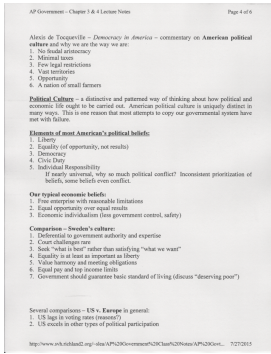
- *Color variability score* indicating color consistency by measuring the height of the histogram bins.
- *Text edge count* based on counting the number of pixels that differ by 100 in gray value on a scale of 0 to 255 from their adjacent pixels.
- *Chroma around text* indicating the distribution of chroma around text edges.
- *Color block ratio* based on counting the number of 32×32 non-overlapping blocks with at least 10% chromatic pixels.
- *White block count* based on counting the number of 32×32 non-overlapping white blocks.

However, [8] is not accurate enough; many highlighted documents are misclassified into the text category and vice versa. Note that highlighter marks on the highlighted documents are made using a highlighter pen after the document has been printed, as shown in Fig. 1(c). Misclassifying these highlighted documents is especially problematic for printers. Most highlighter pens have bright and fluorescent colors [9], so they reflect more light than conventional colors. This reflection will result in unreliable color reproduction by the printer. For example, the scanned or copied highlighters may appear lighter or darker than expected or even change colors, i.e. yellow becomes green. The highlighting may even disappear completely when the scanned document is printed. For these reasons, the need still exists for an improved method for differentiating highlighted documents from other types of documents.

In this paper, we propose two novel gamut-based features and six low-level color features to capture the color specifics of the highlighted regions in the image. These new features are concatenated to the seven features in [8]. The sequential forward floating selection (SFFS) feature selection algorithm [10] is applied to find the best feature subset for our application. The optimum set of features is then used to train a directed acyclic graph support vector machine (DAGSVM) [11] to classify the documents. We work with a specific model of MFP and a specific image processing pipeline equipped with four classes of documents. They are text, photo, highlight, and mixed. Some example images are in Fig. 1. Nevertheless, our work can be easily applied to any MFP by re-measuring the printer gamut. Our cross-validation results show that the new feature subset significantly improved the precision and recall for all document types.

Feature Extraction

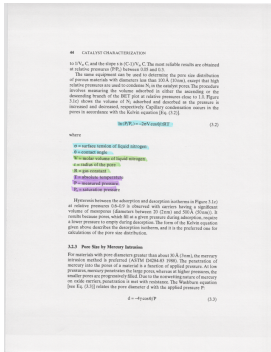
In this section, we will describe the detailed procedure for extracting the new features. Sec. Color Space presents the color space conversion. Sec. Printer Gamut-Based Features demonstrates how to retrieve and characterize the highlighter pixels based on the estimated gamut. Finally, the steps to obtain our new low-level color features are presented in Sec. Low-Level Color Features.



(a)



(b)



(c)



(d)

Figure 1: Example scanned document images: (a) text, (b) photo, (c) highlight, and (d) mixed. Note that in highlight documents, the highlighter marks are drawn manually by the user using a highlighter pen after the page is printed.

Color Space

The choice of the color space is essential for color image analysis [12]. *RGB* is a widely used device-dependent space in imaging devices. However, it is not perceptually uniform. That is to say, the distance between the *RGB* coordinates of two colors is not proportional to the human perception of such a difference [13]. For this reason, several perceptually uniform spaces have been developed, such as CIE $L^*a^*b^*$ and CIE $Lch(a^*b^*)$. In this paper, we will extract color features in the CIE $Lch(a^*b^*)$ color space.

We first convert gamma-corrected *RGB* to linear *RGB* using a 1D gamma uncorrection lookup table (LUT) provided by the organization sponsoring this research. This LUT is applied separately to each of the *R*, *G*, and *B* channels of the MFP scanner output. The conversion is plotted in Fig. 2. Then a 3×3 matrix transformation is applied to the linearized *RGB* values

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} 0.5313 & 0.3519 & 0.1168 \\ 0.2742 & 0.7673 & -0.0415 \\ 0.0051 & 0.0510 & 0.9438 \end{bmatrix} \begin{bmatrix} R_L \\ G_L \\ B_L \end{bmatrix}. \quad (1)$$

This matrix describes the transformation from the device-dependent scanner *RGB* of a particular MFP product to the device-independent CIE *XYZ* color space. It was also provided by the organization sponsoring the research. Then, we convert to CIE $L^*a^*b^*$ through the equations provided in [14]. Finally, the color attributes lightness (L^*), chroma (C^*), and hue (h) can be computed as

$$\begin{aligned} L^* &= L^*, \\ C^*(a^*, b^*) &= \sqrt{a^{*2} + b^{*2}}, \\ h(a^*, b^*) &= \arctan\left(\frac{b^*}{a^*}\right). \end{aligned} \quad (2)$$

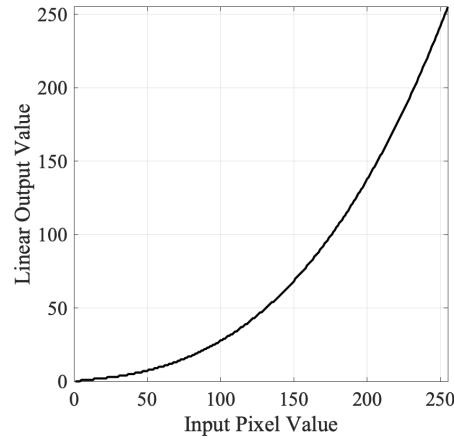


Figure 2: The gamma uncorrection conversion from gamma-corrected *RGB* to linear *RGB*.

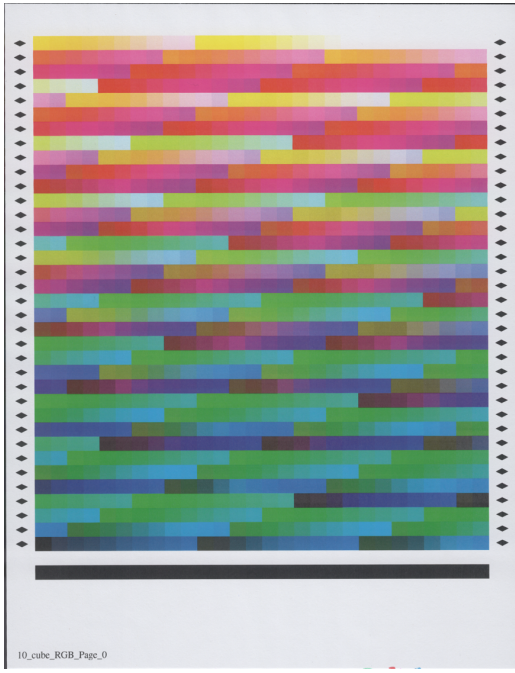
Here, L^* corresponds to how light or dark a color is, C^* represents the color intensity, and h describes the appearance of color – color in its pure form, as in red, green, or blue.

Printer Gamut-Based Features

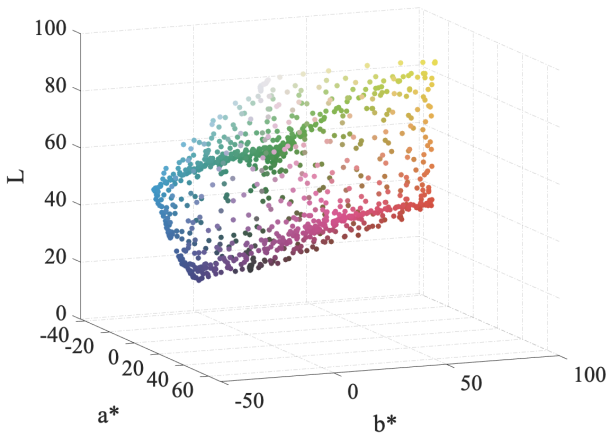
In our discussion here, a gamut is the range of colors that can be reproduced by the printer. In this paper, we only work with printed colors as sensed by a scanner, which has its own gamut of colors that can be uniquely sensed. However, generally, the scanner gamut is larger than that of the printer. So we assume here that the printer gamut is strictly contained within the scanner gamut, and that the scanner gamut contains all highlighter colors.

Given the fact that most colors painted by highlighter pens cannot be accurately reproduced by the printer, we will look at chromatic pixels in the scanned image lying outside of the printer gamut. We follow the procedure described in [14] to estimate the gamut using the test page shown in Fig. 3 (a), which has been printed with our target printer, and then scanned with our target scanner. The resulting gamut is shown in Fig. 3 (b). Inspired by [15], we segment the CIE $Lch(a^*b^*)$ space into 18 non-overlapping 20-degree hue slices. Then the gamut boundary of each hue slice is computed as the convex hull [16] that encompasses all points within the hue slice [17]. We refer to such a convex hull as a gamut hue sector. The vertices of the gamut sectors (convex hulls) are stored in counterclockwise order to facilitate later computations. We will soon compare each image and the gamut at every 20-degree hue slice.

Furthermore, by visual inspection, we note that some highlighter colors, such as yellow and orange, are softer than others, while others, such as magenta and purple, are more visible. From the plots of the four example highlighter patches in Fig. 7, one can see that the range of chroma and lightness varies from one hue slice to another. Thus, we propose to use hue-slice dependent chroma and lightness thresholds to improve the accuracy of highlighter pixel characterization. The thresholds are measured based on a scanned sheet of paper containing assorted colors of highlighter pen marks. It can also be seen from Fig. 7 that many highlight pixels are out of the printer gamut, adding credence to our claim that the highlighted regions often cannot be reliably reproduced. We assume that the chroma and lightness values of a highlight pixel should each fall between two thresholds. From the lightness and chroma 2D histogram in Fig. 4, we note that high chroma/low lightness values could be solid colors that are from a photo or some graphics, and low chroma/high lightness



(a)



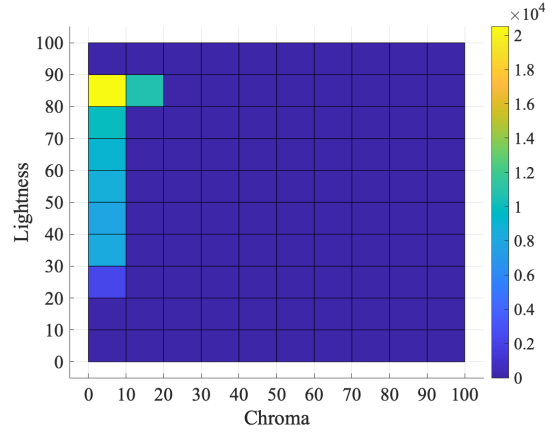
(b)

Figure 3: (a) The printed and scanned test page used to estimate the gamut. (b) The estimated gamut in CIE $L^*a^*b^*$ space. Each dot corresponds to the mean $L^*a^*b^*$ values of a patch.

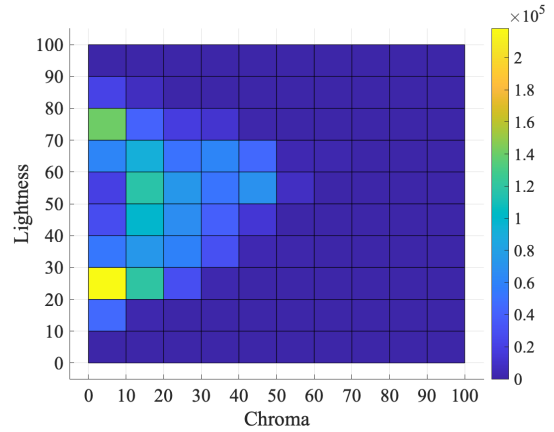
values could be the color of the media. Therefore, to speed up the subsequent computations, such colors can be excluded.

Now, we are going to compare each scanned image and the printer gamut at each hue slice. Specifically, we use the following procedure to extract the highlight pixels of interest. For each 20-degree hue slice s , we first compute a set of image pixels $\mathcal{S}^s = \{(C^*, L^*) \mid C^* \in [C_{LB}^{*s}, C_{UB}^{*s}], L^* \in [L_{LB}^{*s}, L_{UB}^{*s}]\}$, where $C_{LB}^{*s}, C_{UB}^{*s}, L_{LB}^{*s}$, and L_{UB}^{*s} are the lower and upper bounds of the chroma and lightness, respectively, of hue slice s . Note that the set \mathcal{S}^s depends on a specific image and we will repeat this process for each image in the dataset to compute their own feature values. Then, we check if each pixel in \mathcal{S}^s is inside the gamut hue sector or not. Given an edge of the s -th gamut hue sector defined by the vertices $\mathbf{V}_i^s(C_i^{*s}, L_i^{*s})$ and $\mathbf{V}_{i+1}^s(C_{i+1}^{*s}, L_{i+1}^{*s})$, $0 \leq i \leq N-2$, and a pixel $\mathbf{P}_j^s(C_j^{*s}, L_j^{*s}) \in \mathcal{S}^s$, if that pixel lies in the exterior of the gamut hue sector as shown in Fig. 5, then based on [18], $\exists i \in \{0, \dots, N-2\}$ such that

$$(V_{i+1}^s - V_i^s) \times (P_j^s - V_i^s) < 0, \quad (3)$$



(a)



(b)

Figure 4: Lightness and chroma 2D histogram plots of a (a) text image and a (b) photo image. The original images are Fig. 1 (a) and (b).

where \times denotes cross product. Here N is the total number of vertices of the gamut hue sector. Note that the sector is closed since we require that $\mathbf{V}_{N-1}^s(C_{N-1}^{*s}, L_{N-1}^{*s}) = \mathbf{V}_0^s(C_0^{*s}, L_0^{*s})$, and again the vertices are labeled in the counterclockwise orientation.

Now that the highlight pixels of interest are recognized, we can design some features to describe them. We propose to use highlight hue count and maximum highlight strength. The first feature is designed to count the number of highlight colors marked on the document image. The second feature is designed to compute the average distance from each highlight color to the printer gamut boundary for each hue sector. Then, we take as our feature value, the maximum over all the hue sectors of this average distance.

To compute the first feature, we iterate through all hue slices, and count the number of highlight pixels in each slice. If there are a sufficient number of highlight pixels, i.e. at least 1% of the image pixels in the hue slice are highlight pixels, then the hue slice will be counted towards the total number of high-

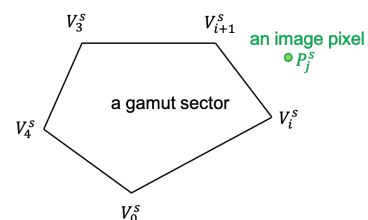


Figure 5: Visual aid for Equation (3).

light hues. As for the second feature, we first need to calculate the distance from each highlight pixel to its corresponding gamut hue sector. For the hue sector s , let \mathbf{P}_k^s be the k -th highlight pixel in the gamut sector s . As illustrated in Fig 6, there are three cases depending on the relative location of the pixel with respect to the gamut sector on the lightness-chroma plane. Let $r_{i,k}^s = \frac{(\mathbf{V}_{i+1}^s - \mathbf{V}_i^s) \cdot (\mathbf{P}_k^s - \mathbf{V}_i^s)}{\|\mathbf{V}_{i+1}^s - \mathbf{V}_i^s\|^2}$, and \mathbf{P}_k^s be the projection of \mathbf{P}_k^s to the edge $\mathbf{V}_{i+1}^s - \mathbf{V}_i^s$ of the hue sector s . Then, the distance from \mathbf{P}_k^s to the edge [19] is

$$d_{i,k}^s = \begin{cases} \|\mathbf{P}_k^s - \mathbf{V}_i^s\|, & \text{if } r_{i,k}^s \leq 0 \\ \|\mathbf{P}_k^s - \mathbf{V}_{i+1}^s\|, & \text{if } r_{i,k}^s \geq 1 \\ \|\mathbf{P}_k^s - \mathbf{P}_k^s\|, & \text{otherwise.} \end{cases} \quad (4)$$

Thus, the shortest distance from the highlight pixel to the periphery of the gamut sector s is $d_{\min,k}^s = \min_{i=0,\dots,N-2} d_{i,k}^s$. Next, we compute the average distance \bar{d}^s over all K^s outlier pixels for the gamut hue sector s according to $\bar{d}^s = \frac{1}{K^s} \sum_{k=0}^{K^s-1} d_{\min,k}^s$. Finally, the second gamut-based feature can be computed as

$$\bar{d}_{\max} = \max_{s=0,\dots,17} \bar{d}^s. \quad (5)$$

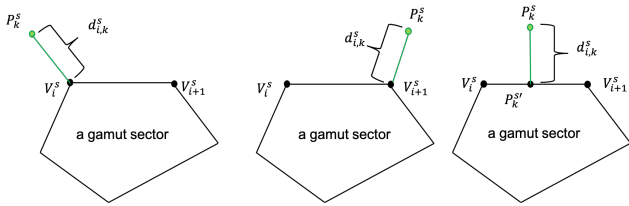


Figure 6: Visual aid for Equation (4). The three cases left to right correspond to the equations from top to bottom in Equation (4), respectively.

Low-Level Color Features

We use two properties of the highlighted regions in order to design the low-level color features. First, highlighter marks are typically bright and relatively translucent colorants drawn on a light-colored background [9]. So the average value of the lightness and chroma in a highlighted image block should be higher than those of a non-highlighted block. Second, within a small region, i.e. 32×32 pixels in a document scanned at 75 dpi or $0.427 \text{ in} \times 0.427 \text{ in}$, there is usually a single highlighter color, and the fluctuations in chroma and lightness should therefore be smaller than those in the mixed and photo images that contain various colors. On the basis of these two properties, we developed six color-moment features [20], namely minimum block mean, maximum block standard deviation, and the minimum block unnormalized skewness of the lightness and chroma channels to describe the characteristics of the color distribution of the highlighted image blocks.

To be specific, we partition the query image into 32×32 non-overlapping blocks, and compute the color moments for both the L^* and C^* channels within each block. Let the pixel value (L^* or C^*) of the i -th channel at the j -th image pixel in block l be $I_{i,j}$ and the number of image pixels in the block be Q , then the block

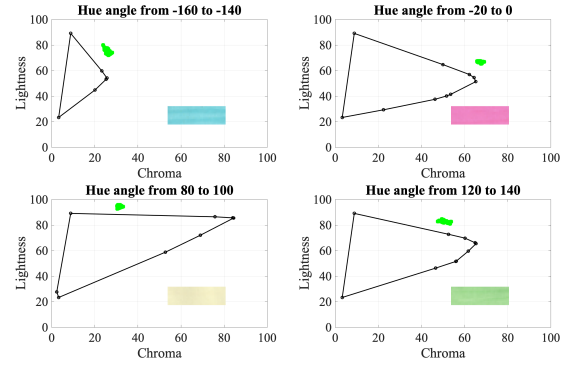


Figure 7: Four highlighter patches with their corresponding gamut hue sectors. The black convex hull indicates the gamut hue sector within the hue range. The green dots are the (C^*, L^*) coordinates of the highlight pixels within the hue range.

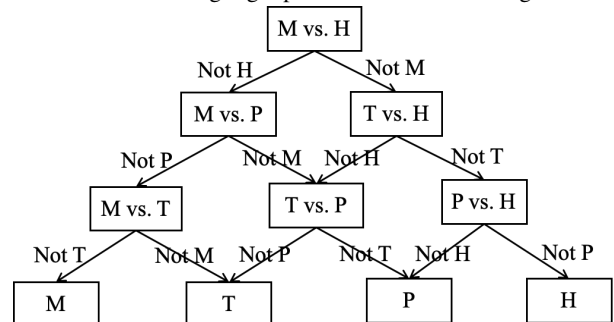


Figure 8: Illustration of the tree structure of the DAGSVM model. M = mixed, T = text, P = photo, and H = highlight.

color moments of channel i are defined as:

$$E_{i,l} = \frac{1}{Q} \sum_{j=0}^{Q-1} I_{i,j}$$

$$\sigma_{i,l} = \sqrt{\frac{1}{Q} \sum_{j=0}^{Q-1} (I_{i,j} - E_{i,l})^2} \quad (6)$$

$$s_{i,l} = \sqrt[3]{\frac{1}{Q} \sum_{j=0}^{Q-1} (I_{i,j} - E_{i,l})^3}$$

Finally, we compute the minimum and maximum color moment values across all blocks. These features are simple, yet effective. Intuitively, they summarize the chroma and lightness characteristics of the most prominent block in the image, which will be the highlighted block, if there is any.

Classification Model

Along the lines of [8], we use a DAGSVM model [11] to solve the multi-class classification problem. The DAGSVM model that we use has a tree structure, as shown in Fig. 8. It consists six 1-vs.-1 SVMs, one for each pair of the four classes. At the root level, the classifier decides if the image is in the mixed or highlight class. If it does not belong to the highlight class, then we go to the left child. If it does not belong to the mixed class, then we go to the right child. This procedure is repeated until the final decision is reached. The radial basis function (RBF) kernel is used for all the SVMs.

In our application, different misclassifications are weighted differently. For example, it is more problematic to misclassify text as photo than to misclassify text as highlight. If the text

		Classifier Output			
		M	T	P	H
Ground Truth	M	0	3	5	4
	T	3	0	10	2
	P	3	10	0	15
	H	10	10	10	0

(a)

		Classifier Output			
		M	T	P	H
Ground Truth	M	114	4	6	5
	T	1	80	0	3
	P	4	0	96	0
	H	5	1	0	81

(b)

Table 1: (a) The error weight matrix W . It shows how different classification results weight differently towards the total cost. (b) The leave-one-out confusion matrix U . It summarizes the performance of our classification model. In both tables, M = mixed, T = text, P = photo, and H = highlight.

image is processed through the photo mode configuration, which has a smoothing effect, the text strokes will look too blurry. However, if the text image is processed using the highlight mode configuration, which will move the colors inside the printer gamut, the reproduction will not be negatively impacted. Therefore, in the training process for the DAGSVM model, our goal is to minimize the weighted error

$$\mathcal{E} = \sum_{i,j} W_{i,j} U_{i,j}, \quad (7)$$

where $W_{i,j}$ is the weight of classifying the i -th class as the j -th class and $U_{i,j}$ is the number of images in the i -th class being classified as the j -th class. The matrix W is presented in Table 1 (a). The weights were chosen by engineers working for the organization sponsoring this research.

Experimental Results

Our dataset consists of the images in [8]. The images were labelled by engineers working for the organization sponsoring this research. There are in total 400 images, including 129 images in the mixed class, 84 images in the text class, 100 images in the photo class, and 87 images in the highlight class. The image contents include book and magazine pages, posters, portrait and natural pictures, handwritten notes, lecture slides, and application forms. Each image has a size of 825×638 pixels and a resolution of 75 dpi.

We evaluate the performance of the model using leave-one-out cross-validation (LOOCV). LOOCV repeatedly splits the data points into a training set containing all but one sample point, and a validation set containing only that remaining sample point. It provides a confusion matrix that we can use to compute the weighted error based on Equation (7). We employ SFFS [10] to select the best feature subset which has the minimal LOOCV weighted error. In the SFFS process, we start from an empty feature set and add one of the non-used features to the set to train the model. Then the cost function is evaluated on the validation dataset. The one feature that gives us the lowest cost will be

included. After each inclusion, a number of exclusions will be performed to the current feature set if the cost can be further decreased. This process is iterated a number of times until there is no further decrease in the cost. Our final feature subset, selected by SFFS, contains 12 features, with 6 from [8] (all but the color variability score) and 6 new ones (2 gamut-based features and the first 4 color-moment features). The optimal LOOCV confusion matrix is shown in Table 1 (b).

Table 2 summarizes the results according to different document image types in terms of precision and recall, as well as the overall accuracy and weighted error. Here, precision and recall are defined [21] as

$$\text{Precision} = \frac{TP}{TP + FP} \quad \text{Recall} = \frac{TP}{TP + FN}, \quad (8)$$

where TP, FP, and FN represent true positive, false positive, and false negative, respectively. It can be computed that the decrease in the weighted error is 61% and the increase in the accuracy is 10.8%. For the highlight class, the precision rises 11.5% and the recall rises 13%. These new results are relative to those reported in [8].

	Precision (%)				Recall (%)				Accuracy (%)	\mathcal{E}
	M	T	P	H	M	T	P	H		
[8]	80.6	76.2	89.0	81.6	83.9	81.0	84.0	78.0	82.0	3.6
Ours	88.4	95.2	96.0	93.1	91.9	94.1	94.1	91.0	92.8	1.4

Table 2: Precision and recall (Equation (8)) for different document types and the overall accuracy and weighted error (Equation (7)). M = mixed, T = text, P = photo, and H = highlight.

Conclusion

A set of highlighter features is proposed. We utilize the characteristics of the highlighter colors and their distribution to describe the highlighted document images. The identification performance of the highlight class was evaluated by means of precision and recall. The overall performance of the model was measured by accuracy and the weighted cost. The newly added features significantly enhance the performance and substantially decrease the cost.

References

- [1] X. Dong, K.-L. Hua, P. Majewicz, G. McNutt, C. A. Bouman, J. P. Allebach, and I. Pollak, "Document page classification algorithms in low-end copy pipeline," *Journal of Electronic Imaging*, vol. 17, no. 4, pp. 043011, 2008.
- [2] H. Cheng and C. A. Bouman, "Document compression using rate-distortion optimized segmentation," *Journal of Electronic Imaging*, vol. 10, no. 2, pp. 460–475, 2001.
- [3] R. L. de Queiroz, "Compression of compound documents," in *Proceedings 1999 International Conference on Image Processing (Cat. 99CH36348)*. IEEE, 1999, vol. 1, pp. 209–213.
- [4] S. J. Simske and S. C. Baggs, "Digital capture for automated scanner workflows," in *Proceedings of the 2004 ACM Symposium on Document Engineering*, 2004, pp. 171–177.
- [5] W. Wang, I. Pollak, T.-S. Wong, C. A. Bouman, M. P. Harper, and J. M. Siskind, "Hierarchical stochastic image grammars for classification and segmentation," *IEEE Transactions on Image Processing*, vol. 15, no. 10, pp. 3033–3052, 2006.
- [6] A. Das, S. Roy, U. Bhattacharya, and S. K. Parui, "Document image classification with intra-domain transfer learning and stacked generalization of deep convolutional neural

- networks,” in 2018 24th International Conference on Pattern Recognition (ICPR), 2018, pp. 3180–3185.
- [7] C. Lu, J. Wagner, B. Pitta, D. Larson, and J. P. Allebach, “SVM-based automatic scanned image classification with quick decision capability,” in *Color Imaging XIX: Displaying, Processing, Hardcopy, and Applications*. International Society for Optics and Photonics, 2014, vol. 9015, p. 90150G.
- [8] S. Xu, C. Lu, M. Shaw, P. Bauer, and J. P. Allebach, “Page classification for print imaging pipeline,” in *Color Imaging XXII: Displaying, Processing, Hardcopy, and Applications*. Society for Imaging Science and Technology, 2017, vol. 2017, pp. 137–142.
- [9] C. Schmid, J. L. Stoffel, and B. Sperry, “Ink compositions for use in highlighter markers and associated methods,” Aug. 30 2011, US Patent 8,007,096.
- [10] P. Pudil, J. Novovičová, and J. Kittler, “Floating search methods in feature selection,” *Pattern Recognition Letters*, vol. 15, no. 11, pp. 1119–1125, 1994.
- [11] J. C. Platt, N. Cristianini, and J. Shawe-Taylor, “Large margin DAGs for multiclass classification,” in *Proceedings of the 12th International Conference on Neural Information Processing Systems*. 1999, p. 547–553, MIT Press.
- [12] G. Paschos, “Perceptually uniform color spaces for color texture analysis: an empirical evaluation,” *IEEE Transactions on Image Processing*, vol. 10, no. 6, pp. 932–937, 2001.
- [13] G. Wyszecki and W.S. Stiles, *Color Science: Concepts and Methods, Quantitative Data and Formulae*, New York, USA: Wiley, 1982.
- [14] S. A. Gindi, “Color characterization and modeling of a scanner,” Masters thesis, Dept. Elect. Comput. Eng., Purdue Univ., West Lafayette, IN, USA, July 2008.
- [15] M. Shaw, “Gamut estimation using 2D surface splines,” in *Color Imaging XI: Processing, Hardcopy, and Applications*. International Society for Optics and Photonics, 2006, vol. 6058, p. 605807.
- [16] R. L. Graham, “An efficient algorithm for determining the convex hull of a finite planar set,” *Info. Pro. Lett.*, vol. 1, pp. 132–133, 1972.
- [17] W. Kress and M. Stevens, “Derivation of 3-dimensional gamut descriptors for graphic arts output devices,” in *Proceedings of the Technical Association of the Graphic Arts*, 1994, pp. 199–199.
- [18] J. F. Whitney and H. M. Whitney, “The right-hand rule,” in *A Handbook of Mathematical Methods and Problem-Solving Tools for Introductory Physics*, 2053-2571, pp. 7–1 to 7–3. Morgan & Claypool Publishers, 2016.
- [19] B. Kolman and D. Hill, *Elementary Linear Algebra with Applications*, Boston MA: Pearson, 2007.
- [20] M. A. Stricker and M. Orengo, “Similarity of color images,” in *Storage and Retrieval for Image and Video Databases III*. International Society for Optics and Photonics, 1995, vol. 2420, pp. 381–392.
- [21] D. L. Olson and D. Delen, *Advanced Data Mining Techniques*, Berlin Heidelberg: Springer, 2008.

Author Biography

Yafei Mao is a PhD student of Electrical and Computer Engineering (ECE) at Purdue University. She received her BS in ECE from Purdue University in 2016. Her PhD research has focused on image processing, halftoning, and computer vision.

Yufang Sun received her BS in Electrical Engineering from the University of Jilin from China (2004). She is currently a PhD student, working as image processing and data analysis research assistant with Prof. Jan Allebach, in the School of Electrical and Computer Engineering at Purdue University. Her research interests are in image information embedding, decoding error analysis, etc. She has been working on the projects of circular coding and stegaframe detection, both sponsored by HP Labs.

Peter Bauer received his Diplom-Ingenieur (FH) in Computer Science from the University of Applied Sciences in Rosenheim Germany. He worked in Hewlett Packard Research Laboratories in Bristol for 4 years. For the last 23 years he has worked for HP Inc. in product development focused on image processing. He is an image pipeline architect for embedded systems, PC and cloud applications.

Todd Harris received his BS in computer engineering from the University of Idaho (1991). Since then he has worked at Halliburton, focusing on nuclear power plant control system data validation, display and archival, and at HP Inc., focusing on scan image processing and compression methods.

Mark Shaw is an Imaging Distinguished Technologist and Strategist at HP Inc. Mark received his B.S. in Graphic Media Studies in the UK, M.S. degree in Color Science, MCSL at RIT, and his PhD in ECE, at Purdue University in the School of Electrical and Computer Engineering. His work includes video coding, color modelling and reproduction, gamut mapping, color management and image understanding through ML. He is on the Industrial Advisory Board at RIT and Boise State.

Lixia Li has been an imaging and optical engineer in HP for 10+ years, and has worked on printers, scanners and VR headsets. She has previous industrial experience in optical lens design, testing, machine vision. Lisa had a PhD degree in Optical Science and Engineering with focus on micro-nano optics design.

Jan P. Allebach is Hewlett-Packard Distinguished Professor of Electrical and Computer Engineering at Purdue University. Allebach is a Fellow of IEEE, IS&T, and SPIE. He was named SPIE IS&T Electronic Imaging Scientist of the Year, received IS&T Honorary Membership, the IEEE Daniel E. Noble Award, the OSA IS&T Edwin Land Medal, and the IS&T Johann Gutenberg Prize. He was elected to Membership in the National Academy of Engineering, and Fellowship in the National Academy of Inventors.