# Document Processing By Automatic Color Form Dropout

Ankush D. Kadu[1], Dr. P. R. Deshmukh[2]

[1] ME student, Department of Electronics and Tele-communication, Sipna college of Engineering and Technology, S.G.B. Amravati University, Amravati (Maharashtra State), India.

[2] Professor, Department of Electronics and Tele-communication, Amravati (Maharashtra State),India

---------------------------------------------------------------------***---------------------------------------------------------------------

**Abstract** - *Color dropout refer as dropping out of color form backgrounds from images of completed forms to Obtain color form dropout images retaining only the respondent information. The successful color dropout simplifies the task of extracting textual information from the image for the reader. At least one non-dropout color is selected and transformed to RGB or a Luminance-Chrominance space. Color dropout process includes scanning of image, conversion to RGB or a Luminance-Chrominance space, pixel analysis, dropout threshold detection, finally storing output ie dropout image. Processing may be performed in RGB or a Luminance-Chrominance space, such as YCBCR. Color dropout is obtained by converting pixels that have color within threshold range of dropout colors to white and all others to will keep as it is in RGB or a Luminance-Chrominance space. This approach uses a FPGA platform which lends itself to high speed hardware implementation with low memory requirements. This is done using VHDL coding. Color dropout processing result may be either represented in RGB or YCBCR.*

*Key Words: Color dropout, Color space conversion, FPGA, MATLAB, Threshold detection, VHDL*

## 1. INTRODUCTION

Color forms constitute a large number of documents that are scanned using high speed scanners. In color forms, text that has been entered, while the document background and lines are not of any practical use. When electronically processing a document such as review form or the like having respondent information entered, so there is a need to remove or dropout the background of the document from a scanned image of the document thereby facilitating minimum storage requirement of image referred as color dropout. Representative documents of this type are health forms, insurance forms, business forms, etc. For performing character recognition on these forms, it is desirable to eliminate the color background of the form, and keep only the textual information that is of relevance.

Color dropout is nothing but the image processing function whose objective is to convert the scanned color document to a black and white image where the color form backgrounds are turned to white and the text colors are turned to black There are several advantages to performing color dropout. By removing background and lines significantly file gets compress and reduces the storage requirements for the resulting document files. The main advantages of color dropout during optical character recognition that information to be read separate from the background information, such as line, boxes and other textual instruction and by means of this minimizes line interference with the text characters, and may reduce complications during character recognition. This process provides elimination of all but the desired information. Once this separation is done, the text from the image are extracted and process by an OCR algorithm. There are several advantages of the present apporach including, but not limited to color removal is performed by evaluating local image content without access to the entire image, minimized image processing time, and the color or colors retained represent the aspects of significant interest to the end user, less memory is required than for other techniques, the process does not require buffering the entire image, an operator is not required to set parameters for each image or image type, the invention reduces the information extraction process time. This process results in the elimination of all but the desired information. In this proposed work, no prior knowledge of the form type is assumed.

Color processing has two approaches as RGB space or Luminance/Chrominance color space. Color dropout based on luminance/chrominance processing involves all the steps that are used in RGB processing, as well as one color space transformation from RGB to YCBCR color space. To accomplish this we need to differentiate between the colors of the background and the colors of the entered text so the image is converted from a full-color form to black and white. This approach lends itself to high-speed hardware implementation with low memory requirements, such as an FPGA platform. Idea to develop an algorithm using MATLAB & VHDL programming for Document Processing for Automatic Color Form Dropout on FPGA platform to get impressive speed of operation increased by using hardware instead of software. Developed VHDL coding for distance coding to be tested using a Xilinx FPGA. The core provides an

excellent amount of processing performance given the FPGA space requirements. It also provides excellent system scalability for much greater performance. The method presented in this approach is designed to operate in an automatic environment and is implemented in hardware. Results for color dropout processing either may be represented in RGB or YCBCR spaces are presented.

## 2. PROPOSED WORK

Dropping out of color form backgrounds from images of completed forms to obtain color form dropout images retaining only the respondent information is sometimes becomes necessary. While printing image from printer sometimes it gives us some distorted image. Distorted image means like something extra color e.g. Pinkish color or yellowish color, etc. gets added into document. So at this stage    it desirable to remove this unwanted added color. This problem of printer motivates to propose the Color Dropout algorithm to remove this unwanted color and keeping whole thing as it is. Resultant document will be a free of distortion.

## 2.1 Objectives

Our main objective is to proposed Color Dropout Algorithm using MATLAB and FPGA

1. To study document processing by observing individual pixel value.

2. Implementation of Color space conversion.

3. To develop an algorithm for color dropout using MATLAB & VHDL programming.

4. To implement on FPGA platform for image processing engine.

5. To be tested and analyze using FPGA various newer technology.

6. Performance evolution parameters of the proposed algorithm to be observed using MATLAB is as follows:

- Min propagation delay
- Device utilization
- PSNR
- MSE
- Correlation factor.

Color dropout is the image processing function who converts the scanned color document to a binary image where the color form backgrounds are turned to white and the text colors are turned to black.

The document is scanned using high speed scanners. In document image processing there is a need to extract textual information from an image that has color content is useful in the background. The removal of the color content is useful in specific applications, such as forms processing, where the color content on the form used to facilitate data entry adds no value to subsequent data processing. Basic assumption is with ink color i.e., darker colors, such as black & dark blue & lighter colors as the part of document background. Color dropout is the image processing function whose purpose is to convert the scanned color document to a binary image where the form background colors are turned to white and the text colors are turned to black.

### 2.2 Color Dropout Architecture

Color dropout system being implemented by using MATLAB and VHDL programming on FPGA platform. Block diagram of color dropout system shown below
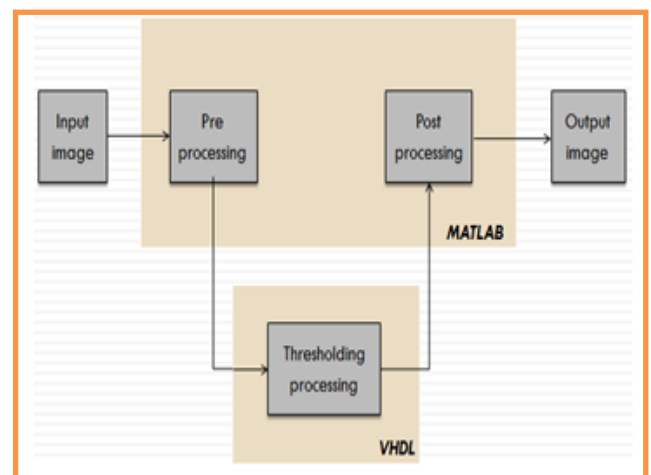


Fig -1 : Block diagram of color dropout system

It involves three main parts as shown below

1) Pre processing

Preprocessing is done in MATAB tool which has more feasible option to visualize image or document. Preprocessing performed in either in RGB space or YCBCR space but in this approach we are performed it using YCBCR and compare their results RGB. Input image given to the preprocessing unit is being resized, converted to YCBCR space, convert this file to text file for further processing. This text file given to threshold processing which is done by using VHDL coding as shown below
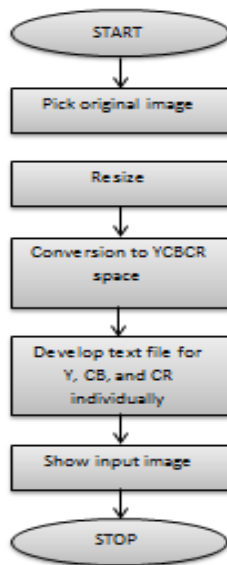
Fig -2: Preprocessing flowchart

- Color space conversion

RGB color space is the most widely used, but it is device dependent and color differences are not perceptually the same throughout the space that's the reason processing in RGB space is not preferred. In this approach, RGB image data is converted into YCBCR color space because YCBCR is more uniform color space as compared to others. It is possible to transform the RGB values to one of the Luminance/Chrominance color spaces, such as CIE Lab. Here we use the YCBCR color space which consists of Luminance Y, Blue Chrominance CB, and Red Chrominance Cr .Even though YCBCR is not perfectly uniform, it has much better characteristics than RGB. Typically, the scanner color output is in red, green, blue (RGB) form. Better results can be obtained by using one of the Luminance Chrominance color spaces, such as CIELUV, CIELAB, or YCBCR. In this unit MATLAB helps to separate out this Red, Green and Blue pixel value for getting more efficiency.

E.g. In image shown below,

R=224, G= 157, B=175; after conversion it will be

Y =170, CB=126, CR=155;

Range for RGB image:

 R=224-255, G=148-158 AND B=170-177;

Range for YCBCR image: Y=170-172, CB=128 AND CR=:167-174

According to above example we can conclude that YCBCR image provides the compress image and also it is more uniform that of RGB image as ranges of R, G and B varies widely compared to YCBCR.

Here we need to know threshold range pixel value range of dropout by observing pixel value of dropout color is determined by using MATLAB tool which needed for threshold processing. We can analyze Y, CB and CR range of particular pixel for color dropout process which may be considered as threshold value used for comparing with other pixel.
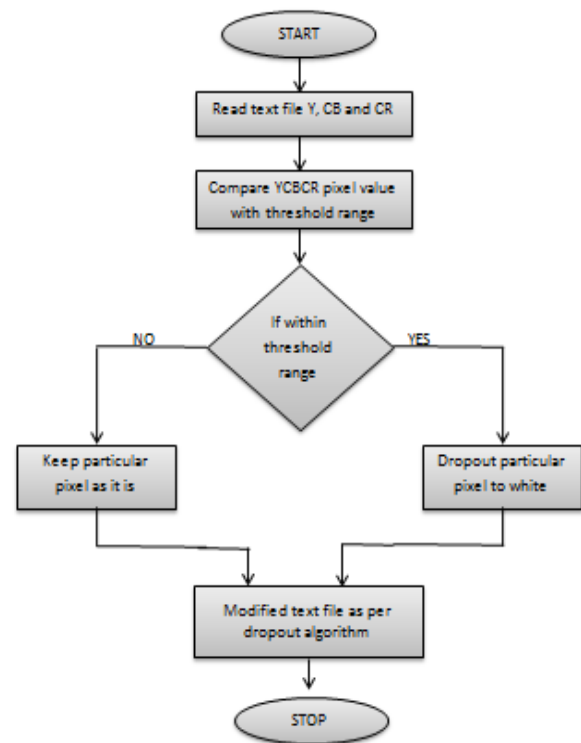
2) Threshold processing



Fig -3: Thresholding flowchart

Threshold processing is the unit where actual color dropout algorithm is implemented because we have developed VHDL coding for color dropout. Process executed in modelsim simulator. Process done at this stage is in bit by bit fashion.

In threshold processing, each image pixel is converted to white if the pixel value is within threshold range otherwise it is keep it as it is and we get modified text file. Now this modified text files given to post processing unit.

   3)  Post processing

Post processing is done in MATAB tool which has more feasible option to visualize image or document. Result may be displayed either in YCBCR space or RGB. In this post processing unit ,output image file is read out by using MATLAB where this text file again converted to image, reshaped , converted to RGB which to visible to human being
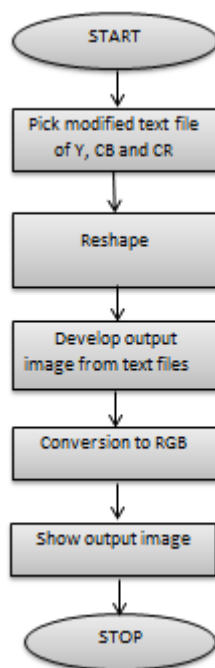


Fig -4:  Post processing flowchart

## 3. RESULT

The system used in this project is evaluated and tested by using two processing method that is either in RGB or YCBCR space. We concentrated on RGB and YCBCR images while comparing these processing techniques, we get know that YCBCR processing has better performance parameters than that of RGB. We have given an emphasis on the color processing as well as data extracting means processing text as well. Comparing this also will get better parameter values for text extraction.

We have developed a conclusion on the basis of parameter explained as shown below and evaluated different parameters for different resolution, different image format on two different category types of processing technique which include another two process technique as follows

- YCBCR processing

  1.  YCBCR color processing

  2.  YCBCR text processing

- RGB processing

  1.  RGB color processing
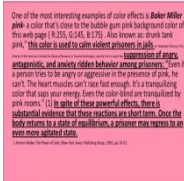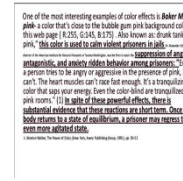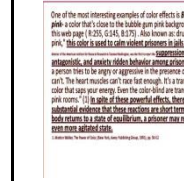
  2.  RGB text processing

| INPUT IMAGE | YCBCR COLOR PROCESS Y=150-190, CB=110-130, CR=150-190 | RGB COLOR PROCESS R=215-255, G=130-160, B=160-180 |
|---|---|---|
| The Power of Colors | The Power of Colors | The Power of Colors |
| ORIGINAL SIZE=147kb | COMPRESSED SIZE=69kb | COMPRESSED SIZE=72k |
| PSNR | PSNR=15.524 | PSNR=15.5775 |
| CORRELATION FACTOR | CF=0.7099 | CF=0.7185 |

Table-1:  Comparison between YCBCR and RGB color processing

In general expected PSNR should a higher for better quality compression and other purpose but in this case it is different approach. As we are dropping some color from image file, our output image must be a different than input image. Hence on this basis we can say that if the output image is equivalent to input, there will be no process being done .In our case error will produce when color gets remove, as error increases, PSNR should be small.

While observing for correlation factor, it also same case as that of PSNR in which it compare the two images, provides the information that how much our output image changes with respect input image. If correlation factor is 1 then we can say that input image exactly resembles to output image but in our case correlation factor must not be 1.
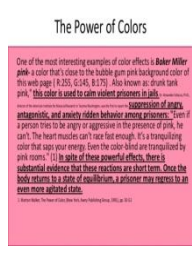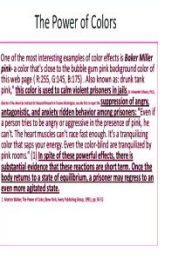
| INPUT IMAGE | YCBCR TEXT PROCESS | RGB TEXT PROCESS |
|---|---|---|
| | Y=0-150 CB=128, CR=128 | R=0-160, G=0-140, B=0-140 |
|  |  |  |
| ORIGINAL SIZE=147kb | COMPRESSED SIZE=68kb | COMPRESSED SIZE=71kb |
| PSNR | PSNR=15.5246 | PSNR=15.5775 |
| CORRELATION FACTOR | CF=0.6771 | CF=0.7019 |

Table-2: Comparison between YCBCR and RGB text processing

According to developed data for different format in table 3, we can conclude that it has better compression and reduces storage requirement for YCBCR text processing rather than that of RGB processing. Next to this YCBCR text processing has better compression ratio. In terms of image format better compression is achieved through tiff image format.

| IMAGE TYPE(kb) | YCBCR PROCESSING | | RGB PROCESSING | |
|---|---|---|---|---|
| | COLOR | TEXT | COLOR | TEXT |
| JPG(24 ) | 24 | 24 | 24 | 24 |
| PNG(125) | 66 | 28 | 84 | 43 |
| BMP(769 ) | 769 | 769 | 769 | 769 |
| TIFF(719) | 119 | 62 | 148 | 102 |

Table-3: Observing size of image for different image format

## 4. CONCLUSION

As per our objectives using MATLAB & VHDL programming on FPGA platform get impressive speed of operation increased. We also show that proposed method provides better PSNR value, which gives more perfect output image. Elimination of color means it significantly reduces the memory requirements for the resulting document files which is the dropout image. By Providing different method of processing i.e. for color and text adds the flexibility and becomes the more Generalize as per requirement. High performance can easily be achieved by simply using a newer technology FPGA. It significantly reduces the storage requirements for the resulting document files which is the dropout image, reduces process time and improves image transmission time. The VHDL synthesis tools provided the speed analysis for and device utilization summary and other component.

## 5. FUTURE SCOPE

Future work will involve approach and to concentrate on all parts of image rather than text and background and also processing in other uniform color spaces with some supervise learning. Automatic detection of background threshold range without any human intervention for any case not only printer or scanner specific. We can work on color dropout for security purpose while dropping out the color in terms secure data, this color contains secure transformation, this secure data then automatically reform at receiver side.

## REFERANCE

[1]    B. Yu and A. Jain, "A Generic System for Form Dropout," IEEE Trans. PAMI, 1998

[2]    J. Mao and K. Mohiuddin, "Form Dropout using Distance Transformation," Proc. ICASSP'95, 1995, pp. 328-331.

[3]    P. Rudak, "Automatic Detection and Selection of a dropout color using zone calibration in conjunction with optical character recognition of preprinted forms," US Patent 5014329, 1991.

[4]    Y. Murai and T. Amagai, "Image processing apparatus with function of extracting visual information from region printed in dropout color on sheet," US Patent 5,664,031, 1997.

[5]    Shima, Y. ; Ohya, H. ; Yasuda, M. " A form dropout method based on line-elimination and image-subtraction", Eighth International Conference IEEE, 2005

[6]    Smith,E.H.B. ; Goyal,S. ; Scott,R. ; Lopresti,D. "Evaluation of voting with Form Dropout Techniques for Ballot Vote Counting",in

Document Analysis and Recognition (ICDAR), International Conference on IEEE ,2011, pp-473-477

[7] Shuli Sun ; Lihua Xie ; Wendong Xiao ; Nan Xiao, " Optimal Filtering for Systems With Multiple Packet Dropouts" , IEEE Trans, 2008,Pp: 695 - 699 .

[8] A.Savakis and J.Madigan," Automatic Color Form Dropout using Luminance/ Chrominance Space Processing," U.S. Patent Number 6035058, 2000.

[9] Gonzalez, R. C., Woods R. E. 2003, Digital Image Processing, Pearson Education.

[10] Li-jun Zhang , Li-Xin Yang , Li-Dong Guo and Jun Li "Optimal Estimation for Multiple Packet Dropouts Systems Based on Measurement Predictor" Sensors Journal,IEEE,2011,Pp:1943-1950.

[11] Jung Uk Cho , Seung Hun Jin 1, Xuan Dai Pham , Dongkyun Kim , and Jae Wook Jeon "A Real-Time Color Feature Tracking System Using Color Histograms"in IEEE International Conference on Control, Automation and Systems ,pp-1163-1167,2007.

[12] A. M. Sapkal , Mousami Munot ,Dr. M. A. Joshi , "R',G'B' to Y'CbCr Color Space Conversion Using FPGA",IEEE.

Ankush D. Kadu,
ME student, Department of Electronics and Tele-communication, Sipna college of Engineering and Technology, S.G.B. Amravati University, Amravati (Maharashtra State), India

Dr. P. R. Deshmukh,
Professor, Department of Electronics and Tele-communication, Amravati (Maharashtra State),India

BIOGRAPHIES