# Document Image Binarization and Segmentation

**[1]Shashidhar Bhat, [2]Prof. Vinod H C, [3]Shaili Srivastava, [4]Shashikumar N, [5]Tulika**
*[2]Assistant Professor*
*Department of ISE, SJBIT, Bengaluru, India*
*Email: [1]shashibhat511@gmail.com, [2]vinod4805@gmail.com, [3]shailisrivastava35@gmail.com,*
*[4]shashikumarn@gmail.com, [5]tulika0615@gmail.com*

## *Abstract*
*Conceptually the Binarization of the chronicled archives is NP-difficult issue since the picture contains commotion, source debasements, and enlightenment. The point of binarization is to locate the best possible picture pixels' limit to enhance the general execution of the framework. This paper presents another half and half meta-heuristic calculation to decide the best edge an incentive for picture archives binarization. The point of Binarization is to locate the correct picture pixels' limit to enhance the general execution of the framework. Record division is a strategy for ripping the archive into unmistakable areas. In this proposed framework at first we displaying Wavelet deterioration and to binarize the record picture, and furthermore utilizes the projection profile to section lines and associated part investigation to fragment the characters. The normal result will be the binarized and fragmented characters, these character can be bolster to OCR for acknowledgement.*

*Keywords: Binarization, wavelet Segmentation, Connected Component Analysis.*

## INTRODUCTION
Over the world, an enormous number of record pictures and generally compositions need to protection and digitization. Inferable from, maturing and physical conditions, these old original copies, and reports are normally seriously debased. Albeit different models and techniques have been produced to accomplish this objective, the binarization of chronicled reports is a basic prerequisite in the preprocessing and upgrade stages for the archive picture investigation and recovery. In the meantime, viable binarization can be utilized to reprocess and upgrade the debased pictures to accomplish better intelligibility. All through the writing, there are numerous propelled binarization techniques, which can be comprehensively classified into two gatherings; worldwide and nearby strategies relying upon how edge esteem is figured for the picture to be divided. In camera captured and scanned documents, where noise, contrast, and brightness vary, segregating pixels as pixels of foreground or background is still a complex and challenging task.

Repo vision is a strategy for ripping the archive into unmistakable districts. A report is a variety of data and a standard mode of conveying data to others. Compatibility of information from reports includes ton of human exertion, time extraordinary and may seriously restrict the use of information systems. So, programmed data compatibility from the record has turned into a major issue. It is being demonstrated that record division will encourage to beat such issues.

## EXISTING SYSTEM
Many common binarization approaches register neighborhood or worldwide limits in view of picture statistics. One drawback of this approach is that edge is invariant to a stage of the pixel. On the other hand, other approaches (e.g. edge detection, Markov irregular field (MRF), connected segments) incorporate solid inclinations about the state of closer view segments.

## PROPOSED SYSTEM

Many significant thresholding strategies have been accounted for achieving picture binarization. The same number of debased archives doesn't have an unmistakable bimodal example, worldwide thresholding is normally not a reasonable approach for the corrupted record binarization. Versatile thresholding, which assesses an adjacent edge for each archive picture pixel, is frequently a superior way to deal with and manage diverse varieties inside debased report pictures. The neighborhood picture differentiate and the nearby picture slope are extremely helpful highlights for sectioning the content from the foundation archive picture in light of the fact that the record message more often than not has certain picture distinction to the adjoining foundation report picture. The proposed approach diminishes the computational complexities of ‚supporting technique after record Binarization and archive division, for example, OCR framework for handling content locales and for preparing non content district.

There are two noteworthy modules in our proposed framework
1. Binarization
2. Segmentation

## RELATED WORK

In literature survey, we found that three binarization classes are Global Binarization, Local Binarization and hybrid binarization. Global approach considers single value of intensity for the separation of text from background. Local threshold use different intensity values depend upon the local information of the different parts of images. Hybrid method use combination of different approaches. Ntirogianniset .al. Propose a joined approach for the binarization, where they utilize Otsu & Niblack methods. They joined the two binarization yields at associated segment level. Worldwide thresholding isn't a decent arrangement for processing printed/manually written report pictures since it can't viably deal with recognizable debasements like black out characters, recolor and drain through which are generally happen in written by hand/printed report pictures. Nearby thresholding additionally not performs well regardless of whether it utilizes a locally fluctuating force esteem that depends either on a chosen property of a pixel from a little window from the image or on measurable assessment from the locally characterized window. There are three principle ways to deal with record segmentation: base up, top-down and hybrid. The process utilizing base up strategies begins from pixel level, pixels are then converged into bigger segments, for example, homogenous square pieces, at that point converged to shape homogeneous districts. Top-down systems continue by beginning from the entire picture and split it repeat safely into various zones until locales in the zone share similar highlights. Top-down systems are efficient for good design organized reports yet frequently bomb in complex design. There are additionally half and half techniques that blend the two beforehand specified techniques. Division utilizing surface investigation falls under the last classification.

## PROCESS FLOW

The stream of Binarization work is appeared in following figure (Assumption Input report is skew amended picture).
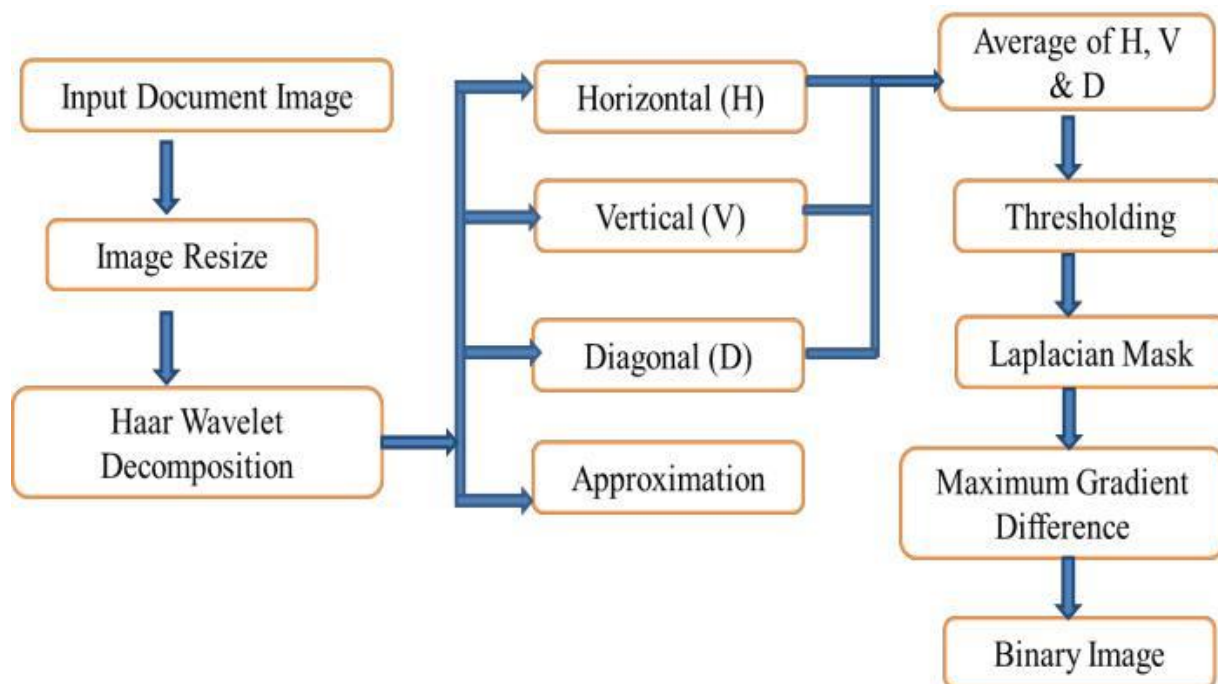
*Fig 1: process flow diagram for image binarization*

The process flow diagram for image segmentation shows the flow of processed binary image to the horizontal, vertical and diagonal projection profile which further undergoes character, word and line segmentation.
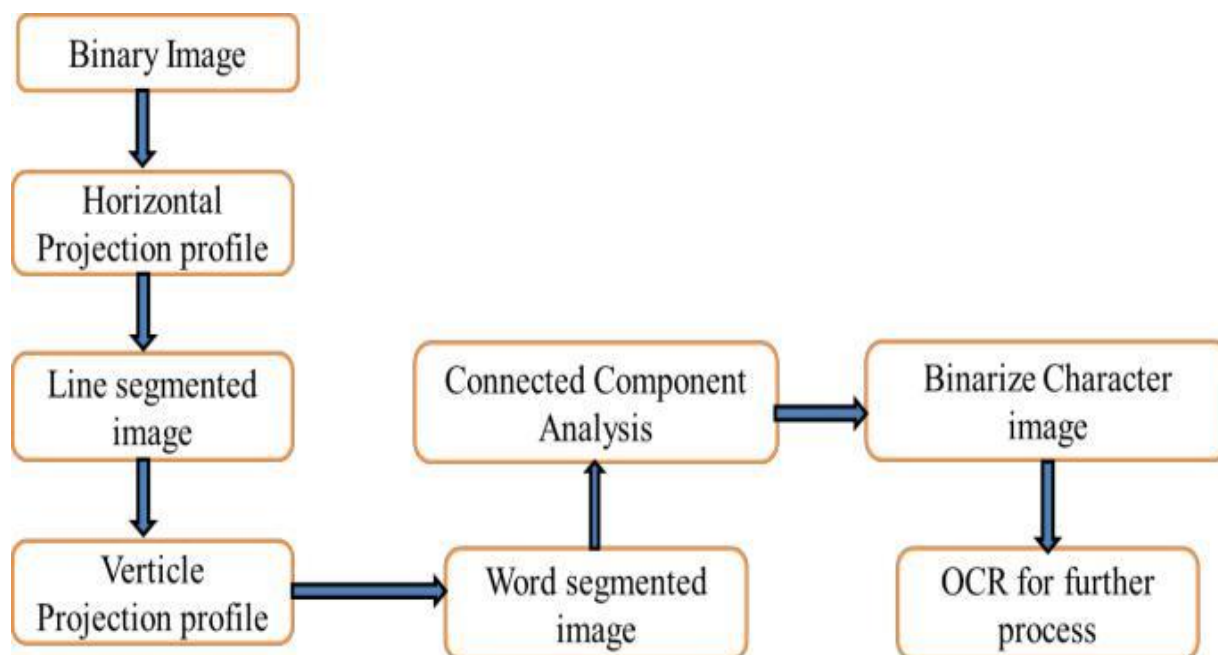


*Fig 2: process flow diagram for image segmentation*

The proposed approach that uses the computational complexities of supporting procedure after document binarization and document segmentation such as OCR systems for processing text regions and for proposing non-text regions.
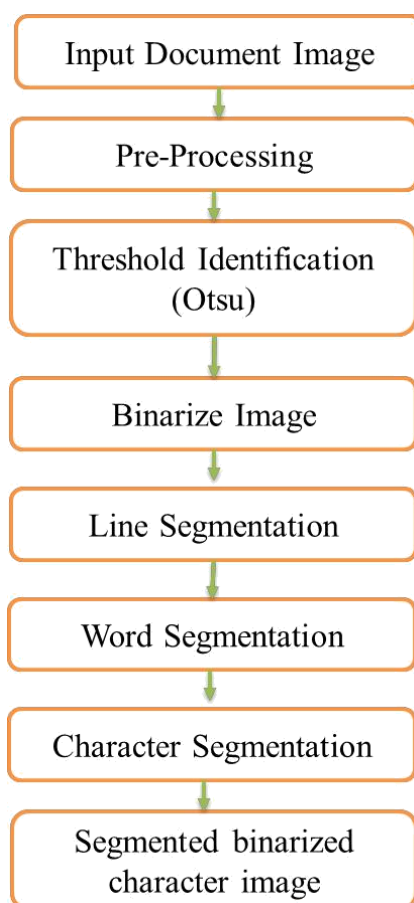
**Fig 3:** *System Architecture*

## METHODOLOGY

**Binarization**-The point of Binarization is to locate the correct picture pixels' limit to enhance the general execution of the framework. Report division is a technique for severing the archive into unmistakable areas. In this proposed framework at first we exhibiting Wavelet deterioration and to binarize the archive picture, and also utilizes the projection profile to portion lines and associated segment examination to section the characters. The normal result will be binarized and divided characters, these characters can be bolster to OCR for recognition. The following calculation is utilized for the same.

**Haarwavelet Decomposition**-In science, the Haar wavelet is a succession of rescaled "square-molded" capacities which together frame a wavelet family or premise. Wavelet examination is like Fourier investigation in that it permits an objective capacity over an interim to be spoken to as far as an orthonormal premise.

Laplacian Operator is likewise a subordinate administrator which is utilized to discover edges in a picture. The real distinction amongst Laplacian and different administrators like Prewitt, Sobel, Robinson and Kirsch is that these all are first request subordinate veils however Laplacian is a moment arrange subsidiary cover. In this veil we have two successive orders, one is Positive Laplacian Operator and other is Negative Laplacian Operator.

Another contrast amongst Laplacian and different administrators is that not at all like different administrators Laplacian didn't take out edges in a specific course

however it takes out edges in following order.

- Inward Edges
- Outward Edges

| 0 | 1 | 0 |
|---|---|---|
| 1 | -4 | 1 |
| 0 | 1 | 0 |

4-connected Laplacian mask

| 1 | 1 | 1 |
|---|---|---|
| 1 | -8 | 1 |
| 1 | 1 | 1 |

8-connected Laplacian mask

**Maximum Gradient Difference:** Obtain the minimum and the maximum gradient values in W over G(x,y) as follows:

**Min(x , y)= min (G($x_i$ ,$y_i$))**

**xi, yi € (x , y)**

**Max (x , y)= max (G($x_i$, $y_i$))**

**xi , yi € W(x, y)**

**GD (x y) = Max(x ,y)-Min(x ,y)**

GD(x,y) indicates the maximum gradient difference, it removes the non-text blocks from the document image.

Otsu Function-The Otsu work is utilized to decide the best edge which amplifies the difference between the diverse picture gatherings. It is characterized as

$$( ) = 0( 0 - )2 + 1( 1 - )2 (2)$$

$$0 = \sum 0 -1 =0 ,1 = \sum 1 -1 = = h \sum h -1 =0$$

Where and $h$ speak to the mean force of and recurrence of the dark level. The best limit is dictated by utilizing: $* = argmax( ( ))$

**Segmentation**- In PC vision, picture segmentation is the way toward apportioning a computerized picture sections (sets of pixels, otherwise called super-pixels). The objective of segmentation is to disentangle and adjust the portrayal of a picture into something that is more critical and less demanding to investigate. Picture segmentation is normally used to find items and limits (lines, bends, and so on.) in pictures. In addition, picture segmentation is the way toward appointing a name to each pixel exhibit in a picture to such an extent that pixels with the basic name share specific attributes fragments (sets of pixels, otherwise called super-pixels). The objective of segmentation is to rearrange and additionally change the portrayal of a picture into something that is more important and less demanding to examine. Picture segmentation is commonly used to make sense of articles and limits (lines, curves, and so on.) in pictures. All the more absolutely, picture segmentation is the way toward allotting a name to each pixel in a picture to such an extent that pixels with a similar name share certain qualities.

**RESULTS**
The introduced method is used for both printed and handwritten manuscripts. The input document image awsscanned, binarized into grey scale image and while went under the process of recursive image segmentation (line segmentation, word segmentation, character segmentation).

The resulting image was fed into OCR (Optical Character recognition) and a processed image was obtained with 99% data loses removed.

## CONCLUSION

This paper presents a document image Binarization method that is compatible to different types of document degradations. The proposed technique is straightforward and powerful, just couple of parameters are included. We have exhibited an approach for document image processing by making use of bianrization and segmentation. The document input image (handwritten and printed) is first converted into a binary image using Haar wavelet decomposition and the processed image then undergoes the segmentation process and fed to OCR. The final image is obtained enhanced effectiveness with sharpened edges and reduced data loss.

## FUTURE WORK

Learning on the extent of the substance and the fitting square size is one of the critical components for effective partition. Future work will include finding a versatile technique in light of the setting to set suitable piece sizes to build the precision of the outcomes. One of the rest of the issues is the outskirts of every segment, since each square has been dealt with as a unit which may bring about a piece containing at least two different substance and this more often than not happens around the boondocks of the segments. Post-handling for finding better outskirts is required for the expansion in exactness of the division. Since each piece has been dealt with as a unit which may bring about a square containing at least two different substances. Furthermore, this generally happens around the outskirts of the segments. Post-preparing for finding better wildernesses is required for the expansion in precision of the division.

## ACKNOWLEDGEMENT

## REFERENCES

1. Mohammed Mudhsh, Shengwu Xiong "Hybrid swarm Optimisation for Document Image Binarisation based on Otsu Function ", CASA 2017.
2. Jino PJ, "Combined Approach for Binarisation of offline Handwritten Documents", IEEE 2017.
3. Zhongi Wang, Jin Zhang, Jing Huang,"Multi-granularity Hierarchial topic based Segmentation of structured, digital library resources", The electronic library, 2017.
4. M W Lin, J R Tapamo, B Ndovie,"a Texture based method for Document Segmentation and Classification", University Kwa-Zulu-Natal, 2017
5. Jino P J, Kannan Balakrishnan, "Combined Approach for Binarization of Offline Handwritten Documents" 4th International Conference on Electronics and Communication System (ICECS), 2017.
6. Michele Alberti, Manuel Bouillon, Rolf Ingold, Marcus Liwicki, "Open Evaluation Tool for Layout Analysis of Document Images", arXiv:1712.01656v1 [cs.CV] 23 Nov 2017.
7. Christoph Wick and Frank Puppe, "Fully Convolutional Networks for Page Segmentation of Historical Document Images", arXiv:1711.07695v1 [cs.CV] 21 Nov 2017.
8. Priyadharshini N, Vijaya MS, "Genetic Programming for Document Segmentation and Region Classification Using Discipulus", (IJARAI) International Journal of Advanced Research in Artificial Intelligence, Vol. 2, No. 2, 2013.

9. J. Sauvola and M. Pietika inen, "Adaptive document image binarization," *Pattern recognition*, vol. 33, no. 2, pp. 225–236, 2000.

10. N. Otsu, "A threshold selection method from gray-level histograms," *Automat- ica*, vol. 11, no. 285-296, pp. 23–27, 1975.

11. Most of our works have been referred from google.co.in