

Bangla PDF Speaker : A Complete Computer Application to convert Bangla PDF to Voice

1st Md. Mizanur Rahaman Nayan

EEE department

BUET

Dhaka, Bangladesh

mdmizanur940@gmail.com

S

Abstract—In pdf data is stored as content so we lost our text form which is in Unicode3. [1] .that is why extracting text from pdf directly is not possible. So, for extracting text from a PDting threshold for binarization which increases efficiency highly [2] So, after taking image from pdf page we do some image processing and and make it more useful as an input for pytesseract engine. Then pytesseract engine extract text from the image. We store this text in a tuple .At the time of tuning the program, We have measured the accuracy for the test images and improved the parameter value that we used in image processing . When we found the accuracy is high enough then we used the parameter for general images which are automatically generated from pdf file. After extracting text using OCR ,we normalize the texts. Normalization increases words readability and spoken sentence quality [3].After normalization we send the text to the text to speech engine ,gTTS. This engine convert text to speech that we hear.

Index Terms—component, formatting, style, styling, insert

I. INTRODUCTION

Bangla text extraction and text to voice conversion has become very popular now a days. Many initiative has been taken for Bangla text extraction [4] [5] [6] ,conversion of text to speech [7] and even pdf to voice [8]. But there is no such complete application that can convert Bangla PDF to voice with high accuracy. Extracting feature from different types of pdf file get's error prone if we don't process and normalize the extracted text from the pdf files properly. In this paper , We are presenting Bangla PDF speaker .It is a computer application that can read through a pdf containing a Bangla Writings . Basically, this application is a cascaded version of four big area of research. They are as follows: Image Processing, Image to Bangla Text , Text normalization and finally text to speech .We have focused on image processing part and text normalization part which has made the application more efficient. This application will be very helpful to listen any PDF instead of readings. English pdf speaker is available, but there is no such application that can convert bangla pdf to voice with more than mean 94% accuracy

II. METHODOLOGY

A. Taking Input

Our program will take a pdf file as input and range of page number that the user want to hear. I have used input() function

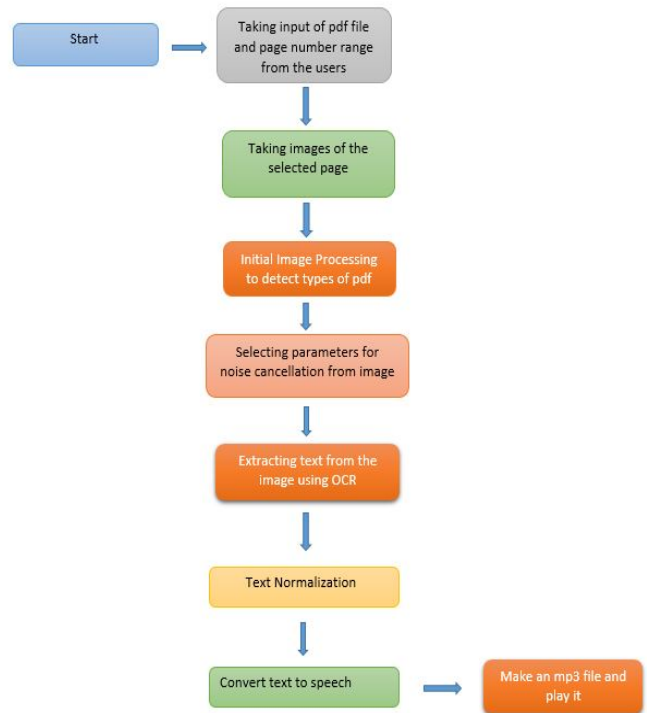


Fig. 1. Flow-Chart of the project

for taking input of the pdf file location as string and range of page number as string. Then I converted page number into integer from string for feeding it in next step.

B. Converting PDF's selected page to images:

I have used a python lib pdf2image that convert pdf pages to images of desired size.

C. Image Processing

Image processing is the most important part of our program since efficiency of the program depends largely on the quality of processed image. If we don't give properly filtered and noise cancelled image we won't get good accuracy. We will explain how accuracy falls if we don't use image processing properly .We will complete image processing in two steps.

First step is to detect the quality of image which measures mean of standard deviation of each pixel value and mean value of pixels . In second step, depending on the value of mean standard deviation and mean pixel value we will set parameters for next level image processing where we will grey-scale and binarize our image finally and will send to pytesseract engine to extract text from the image .We will set a threshold to binarize . Thresholding is a common image processing operation applied to gray-scale images to obtain binary or multilevel images.O’Gorman, Lawrence show how thresholding helps to increase image readability [9]

1) Image Pre-processing to detect type of noise in image:

We convert images in grey-scale. After that we calculate amount of noise by calculating standard deviation and mean of pixel .We found if mean standard deviation of pixel values is very high then we can identify the image of low noise since we expect image only contain bangla text. So, by measuring mean standard deviation we can measure the quality of noise of an image . We classify different type of pdf file depending on the mean standard deviation value and mean of pixel value. Where high standard deviation with high mean means high quality of image hence good pdf and with low mean standard deviation means low quality of image hence low quality of pdf .We will now explain how by measuring mean and standard deviation of pixel value we can measure image noise quantity.



Fig. 2. Image Pre-processing

To do this at first we take image as an array and then we flatten it to give it input to the numpy.std () function. Then we use numpy.std() and numpy.mean() function to measure standard deviation and mean of pixel values.We observed for several picture but here we are summarizing by presenting few image and their mean and standard deviation showing the noise measurement.(Here noise is considered as very close pixel value to the pixel value of writings)

আমাদের দেশের নাম বাংলাদেশ। এখানে আসে পাখি বাকেরাকো। পাখির
কলকাকলিতে মুখরিত থাকে সবসময়। আমার সেনার বাংলা আমি তোমায়
ভালবাসি। আমাদের দেশের নাম বাংলাদেশ। এখানে আসে পাখি বাকেরাকো।
পাখির কলকাকলিতে মুখরিত থাকে সবসময়। আমার সেনার বাংলা আমি
তোমায় ভালবাসি।

Fig. 3. Image with very low noise

average	float64	1	240.2959423517041
im	Image	(951, 340)	<JpegImageFile @ 0x29AAD8E9FD0> ...
im_arr	uint8	(340, 951, 3)	[[[255 255 255] [255 255 255] ...
im_flatten	uint8	(970020,)	[255 255 255 ... 255 255 255]
standard_deviation	float64	1	53.15424899218377

Fig. 4. Standard deviation and mean of Image in figure 2

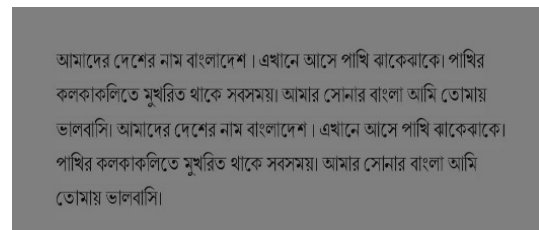


Fig. 5. Image with medium noise

average	float64	1	121.2318348580141
im	Image	(954, 352)	<JpegImageFile @ 0x29AAD8E9630> ...
im_arr	uint8	(352, 954, 3)	[[[127 127 127] [127 127 127] ...
im_flatten	uint8	(1007424,)	[127 127 127 ... 127 127 127]
standard_deviation	float64	1	23.825888805954975

Fig. 6. Standard deviation and mean of Image in figure 4

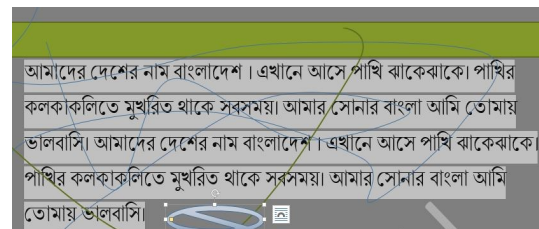


Fig. 7. Image with high noise

average	float64	1	149.56073189509723
im	Image	(833, 342)	<JpegImageFile @ 0x29AA96574E0> ...
im_arr	uint8	(342, 833, 3)	[[[127 127 127] [127 127 127] ...
im_flatten	uint8	(854658,)	[127 127 127 ... 127 127 127]
standard_deviation	float64	1	51.09475578476383

Fig. 8. Standard deviation and mean of Image in figure 6

From our observation on more than 50 different types of noisy and clean image we classified them as follows on basis of standard deviation and mean value of pixels:

Standard Deviation	Mean	Type
High(> 40)	High(> 200)	Very small noise(Standard pdf)
High(> 40)	Low(< 170)	Very noisy
Low(< 40)	High(> 200)	Noisy
Low(< 40)	Low(< 170)	Medium Noisy

TABLE I

CLASSIFICATION OF IMAGE WITH RESPECT TO NOISE

Now, depending on the type described above program will set binarization threshold value automatically. It's very important to select threshold value properly. Otherwise, accuracy will be very poor [9] [10] Threshold value differ from image to image which depends on type of noise. We have found optimum threshold value for different types of image e.g. One value for very small noisy image and another value for other noisy image.so that our program can extract text from different types of pdf files with very high accuracy .If we use same threshold value for different images of different type accuracy will fall . We'll discuss it in result and analysis section again.

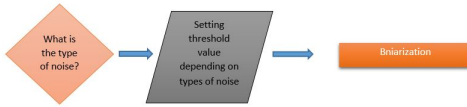


Fig. 9. Image post processing

2) *Extracting text from image:* We have used Pytesseract engine to perform this task. Pytesseract Engine: Pytesseract or Python-tesseract is a wrapper for Google's Tesseract-OCR engine. It extract text from an image. We use data for Bengali language to extract Bengali text from the image. Pytesseract engine can extract data from a noise free image with 100% accuracy. But if it's given noisy image then it's accuracy falls tremendously. Let's test for an image: Following image is directly given to the engine and after processing , again the processed image has sent to the engine. We observed the following result:

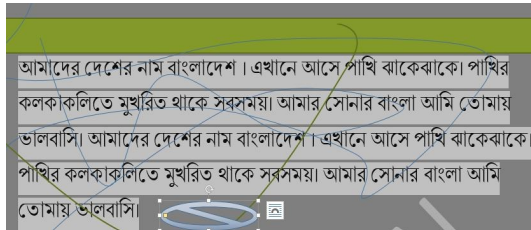


Fig. 10. Input Image

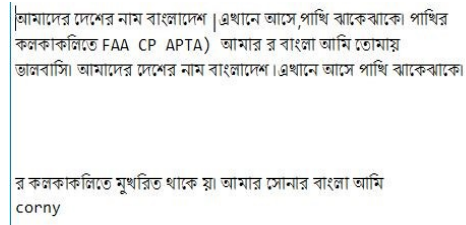


Fig. 11. Output of pytesseract without processing

Observation: Indentation has break down and failure of text extractions is 9 words.

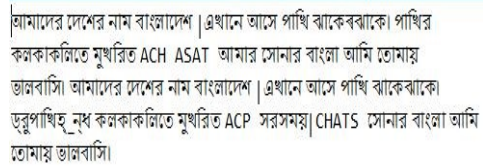


Fig. 12. Output of pytesseract with processing

Observation: By applying , preprocessing indentation has been preserved and number of failure of text extraction is 4 words. That is error words amount has reduced by less than half.

3) *Text Normalization:* Text Normalization is basic text processing before sending text to a text to speech engine. Without text normalization speech form will hamper a lot (e.g if there is a date write on as 12/12/12 then text to speech engine will say 'twelve twelve twelve ' in 'Twelve December Two thousand twelve') [3]. So make speeches more realistic we have done text normalization to the text that we get from pytesseract engine.To make paper independent, normalization process has described as follows:

4) *Converting Normalized text to Voice:* We used Google text to speech engine for converting text to voice. It take input as text format and then produce an mp3 file.There are several speech synthesis system like festival synthesis system [11] and google. But we have used google text to speech synthesis system due to its efficiency and pronunciation is much better than others.

III. EXPERIMENTAL RESULT & DISCUSSION

A. Checking Image Quality generated from PDF

We have used pdf2image python library's pdf2image() function to generate image from pdf of specific pages.We will see some of the image that we have generated from pdf in this section .

Following two images are of randomly selected two pdf page that we converted to image from pdf format.We generated images of size (800,1200) of each.

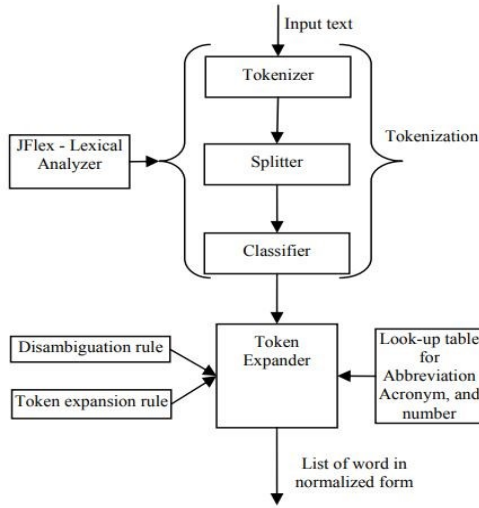


Fig. 13. Text Normalization process [3]

তত্ব
আল-কিনায়া ওয়ান নিহায়া

সাথে বিয়ে শাদী করবে না। তাদের সাথে বোকা-কোলা লেনদেন করবে না এবং তাদের আশ্রয় দেবে না।-অর্থী যতক্ষণ না তারা রাসুলুল্লাহ (সা)-কে তাদের হাতে তুলে দেয়। এ সময় নবী করীম (সা) আরো বলেছিলেন: لا يكثر المسلم الكفار ولا الكفار المسلم. মুসলমান কাকিদেরের দিরাঙ্গ পাশে না এবং কাকির ও মুসলমানের ওয়ায়িহ হবে না। (মধ্যবর্তী) রাবী মুহেবী (র) বলেন, এ হাদীসের ধারক (খিফ) শব্দের অর্থ হল উপভাষ্য। বুখারী মুসলিম (র) এ হাদীস আবদুর রায্বাক (র) সূত্রে উদ্ধৃত করেছেন।

বুখারী ও আহমদ (র)-এর এ হাদীস দুটি থেকে প্রতীয়মান হয় যে, নবী করীম (সা) মুহাসসায়ে অবস্থানের ইচ্ছা করেছিলেন একটি মুখ্য উদ্দেশ্য বুখারীশী কাকিররা রাসুল (সা)-বিরুদ্ধে জোটবদ্ধ হয়েছিল এবং রাসুলুল্লাহ (সা)-কে তাদের হাতে সোপান করার দাবীতে বনু হাশিম ও বনু মুজাশিহের বিরুদ্ধে মুক্তি পত্র সাক্ষর করে (কা'বা ঘরে তুলিয়ে-রেখে)। ছিল তার বার্ষিক পরিণতি ও তাতে তাদের হেয়তা ও পরাজয়ের প্রাণির স্মৃতি বস্তু (যেমন সর্পিষ্ট অধ্যায় বর্ণিত হয়েছে)। ওখানে অবতরণের সিদ্ধান্ত নিয়েছিলেন। অনুরূপ মক্কা-বিজয়ের সময়ও সেখানে অবতরণ করেছিলেন। এ দৃষ্টিকোণের বিচারে সেখানে অবতরণ একটি সুন্নাত ও কাকিতক বিষয় সাব্যস্ত হবে। এক-এটি আশিমাগণের এ বিষয় সম্পর্কিত দুই অভিমতের একটি।

পঞ্চমতঃ, বুখারী (র) বলেন, আবু মু'আযম (র) আইশা (রা) সূত্রে তিনি বলেন, তা তো ছিল একটি (নাযার) মানসিক যেখানে নবী করীম (সা) অবতরণ করতেন শুধু তার (পরবর্তী) সফর ও প্রস্থানের সুবিধার্থে অর্থাৎ আবজাহ (মুহাসসায়ে) থেকে। মুসলিম (র) এ হাদীস রিওয়াত করেছেন। হিশাম (র) থেকে এ সমস্ (আবু আবু দাউদ (র)-এর রিওয়াত আহমদ ইবন হাফল (র) আইশা (রা) হতে রাসুলুল্লাহ (সা) মুহাসসায়ে অবস্থান করেছিলেন শুধু তার প্রস্থান সহজ হওয়ার উদ্দেশ্যে তা সুন্নাত নয়। সুতরাং যার ইচ্ছা সেখানে অবতরণ করবে, যার ইচ্ছা অবতরণ করবে না। বুখারী (র)-বলেছেন, আলী ইবন আবদুল্লাহ (র) সুফিয়ান ইবন আকাস (রা) হতে তিনি বলেন, (মুহাসসায়ে অবতরণ অবস্থান) মন্বিল যেখানে (খাজিক জায়েই) নবী করীম (সা) অবস্থান করেছিলেন। মুসলিম (র) এ হাদীসটি রিওয়াত করেছেন। আবু বকর ইবন আবু শায়বা (র) প্রমুখ সূত্রে (সুফিয়ান ইবন উয়ায়না) হতে এ সমস্। আবু দাউদ (র)-বলেন, আহমদ ইবন হাফল (র) উছমান ইবন আবু শায়বা (শব্দ অর্থ) ও মুলাদাদ (র) (অর্থ-জাফা) সুফিয়ান ইবন উয়ায়না (র) হতে, তিনি বলেন, আবু রাফি (রা) বলেছেন, তিনি অর্থঃ রাসুলুল্লাহ (সা) আমাকে সেখানে অবতরণ করার হুকুম করেন নি। তবে সেখানে তারু তৈরী করা হলে তিনি সেখানে অবতরণ করতেন। মুলাদাদ (র) বলেছেন। আবু রাফি (রা) ছিলেন নবী করীম (সা)-এর বোকা বহরের শিখার। উছমান (রা) বলেছেন। (সেখানে) অর্থঃ আবজাহে মুসলিম (র) এ হাদীসটি রিওয়াত করেছেন- সুতরাং। মুহাম্মদ ইবন হারব ও আবু বকর (র) সুফিয়ান ইবন উয়ায়না (র) হতে এ সমস্।

এ আসোচনা লক্ষ্য হল যিনি থেকে প্রস্থান কালে নবী করীম (সা)-এর মুহাসসায়ে অবতরণের ব্যাপারে এরা সকলেই একমত পোষণ করেছেন। তবে তা উম্মিহ হওয়ার ব্যাপারে তাঁদের মাত্র মত পার্থক্য সেবা দিয়েছে। কারো কারো মতে তিনি সেখানে কোন বিশেষ

Fig. 14. image of pdf page no 336

B. Efficiency Measurement

[h!] We have shown in Methodology section, how image processing increased overall efficiency of our computer application. Now we will measure overall accuracy of our application. We will use following formula for measuring accuracy of our application.

$$= \frac{\text{number of accurately detected character}}{\text{Total number of characters}} * 100\%$$

let's test for last two images that we have generated in last subsection.

for first page:

ইমাম আহমদ (র) বলেন, মুহ ইবন মায়দুন (র) ইবন উমর (রা) হতে এমর্যে বর্ণনা করেন যে, রাসুলুল্লাহ (সা) এবং আবু বকর, উমর ও উসমান (রা) মুহাসসায়ে অবতরণ করেছেন। ইমাম আহমদ (র)-এর মুসনায়ে আবদুল্লাহ আল আমরী দাফি (র) সনদের হাদীস রূপে আমি এভাবে অধ্যয়ন করেছি। তিরমিযী (র) এ হাদীসটি রিওয়াত করেছেন, ইসহাক ইবন মনসুর (র) হতে। ইবন মাছা (র) রিওয়াত করেছেন, মুহাম্মদ ইবন ইয়াহয়া (র) হতে- (উভয় সনদ....) দাফি-ইবন উমর (রা) হতে, তিনি বলেন, রাসুলুল্লাহ (সা), আবু বকর, উমর ও উছমান (রা) আবজাহে অবতরণ করতেন। তিরমিযী (র) বলেছেন, এ প্রসঙ্গে আইশা, আবু রাফি ও ইবন আকাস (রা) হতে রিওয়াত রয়েছে। ইবন উমর (রা) হতে আবু হাদীসটি হাদান-গারীব একক সূত্রে উত্তম।

কেননা, শুধু আবদুর রায্বাক- উবায়দুল্লাহ ইবন উমর (রা) সনদের হাদীসে এটিই সাথে আমাদের পরিচিতি। মুসলিম (র) এ হাদীসটি রিওয়াত করেছেন। মুহাম্মদ ইবন মিহ্রান আল রাযী (র) ইবন উমর (রা) সূত্রে এমর্যে যে, রাসুলুল্লাহ (সা) এবং আবু বকর ও উমর (রা) আবজাহে অবতরণ করতেন। মুসলিম (র) সাব্বুর ইবন হুওয়ায়রিয়া দাফি ইবন উমর (রা) হতেও রিওয়াত করেছেন যে, তিনি ইবন উমর মুহাসসায়ে অবতরণ করতেন এবং দাফির দিবসের যুহর সাতাত হাসবার (মুহাসসায়ে) আদায় করতেন। দাফি (র) বলেন, রাসুলুল্লাহ (সা) এবং তাঁর পরবর্তী খলীফাগণ মুহাসসায়ে অবতরণ করেছেন। ইমাম আহমদ (র) আরো বলেন, ইউনুস (র)....ইবন উমর (রা) হতে এ মর্যে, রাসুলুল্লাহ (সা) যুহর, আশর, মাগরিব ও ইশার সাতাত ব্যাভাষ্য (আবজাহে) আদায় করতেন, তারপর কিছুক্ষণ তরে থাকতেন। তারপর মক্কার প্রবেশ করতেন এবং ব্যাভাষ্যে তারপর বসতেন। আহমদ (র) আহমদ (র)....ইবন উমর (রা) হতেও রিওয়াত করেছেন। এ রিওয়াতের শেষ অংশে অধিক রয়েছে ইবন উমর (রা) ও তা ফরতেন। আবু দাউদ (র) ও আহমদ ইবন হাফল (র) হতে অনুরূপ রিওয়াত করেছেন।

বুখারী (র) বলেন, হুযায়না (রা)....(মুহেবী)....আবু হুযায়না (রা) হতে তিনি বলেন, যিনার দায় নিবসের দশ তারিখ-এর পরের দিন রাসুলুল্লাহ (সা) বলতেন-
نحن نزلون عدا بخيف لبي كذابة حيث تاملوا على الكفر

আমরা উগারী দিন বনু কিনানা উপত্যকায় অবতরণ করব, যেখানে তারা মুহেবীতে জট থাকার মুক্তিতে অশীকারাবদ্ধ হয়েছিল। তিনি এ গিঠিত্য বলতে মুহাসসায়ে বুদ্ধিবিহীন। (পূর্ণ হাদীস) মুসলিম (র) হাদীসটি অনুরূপ রিওয়াত করেছেন। মুহাম্মদ ইবন হারব (র) (আবু হাযী) সূত্রে ইমাম আহমদ (র) বলেন, আবদুর রায্বাক (র) উসমান ইবন যাদম (রা) হতে। তিনি বলেন, আমি বললাম, ইয়া রাসুলুল্লাহ (সা) আশা কী কীভাবে অবস্থান যেমন? এটি তার হজ্জের সময়ের কথা। তিনি বলেন, لم نزلنا على مكة. আলী (ইবন আবু তালী) কি আর আমাদের কোন খবরটি রয়েছে? তারপর বসতেন, আশা কী কীভাবে ইমাম আহমদ (র) বনু কিনানা উপত্যকায় অর্থঃ মুহাসসায়ে-অবস্থান যে যেখানে তারা কুরআনীদের সাথে কুম্বীতে ও অশীকারাবদ্ধ হয়েছিল। ঘটনাটি হল যে, বনু কিনানা বনু হাশেমীর বিরুদ্ধে আশা কুরআনীদের সাথে এ ব্যাপারে শপথ মুক্ত আঁতাত ও মুক্তি করল যে তারা হাদীসীদের

Fig. 15. image of pdf page no 337

$$= \frac{1132}{1189} * 100\% = 95.21\%$$

for second page:

$$= \frac{1378}{1451} * 100\% = 94.902\%$$

from above result we observed that accuracy is almost 95%

IV. CONCLUSION

In this project , we tried to develop a computer application that can read pdf file. We tried to increase the accuracy of the system by implementing processing method where image taken from the pdf is classified depending on the value of standard deviation and mean of the pixels . The method automatically select parameters(e.g threshold for binarization) depending on the value of SD and Mean. We have also focused on normalization process to increase text readability. Our application is ready and it Is very efficient on reading any bangla pdf .

REFERENCES

- [1] Deepak Massand. System and method for reflowing content in a structured portable document format (pdf) file, September 30 2010. US Patent App. 12/413,486.
- [2] B Fadiora, F Wada, and OB Longe. Combining optical character recognition (ocr) and edge detection techniques to filter image-based spam. *African Journal of Computing & ICT January*, 5(1):59–68, 2012.
- [3] Firoj Alam, SM Habib, and Mumit Khan. Text normalization system for bangla. Technical report, BRAC University, 2008.
- [4] Umapada Pal, Partha Pratim Roy, Nilamadhava Tripathy, and Josep Lladós. Multi-oriented bangla and devnagari text recognition. *Pattern Recognition*, 43(12):4124–4136, 2010.
- [5] Supriya Kurlekar. Reading device for blind people using python, ocr and gttts.

- [6] Umapada Pal and BB Chaudhuri. Ocr in bangla: an indo-bangladeshi language. In *Proceedings of the 12th IAPR International Conference on Pattern Recognition, Vol. 3-Conference C: Signal Processing (Cat. No. 94CH3440-5)*, volume 2, pages 269–273. IEEE, 1994.
- [7] Shaheena Sultana, MAH Akhand, Prodip Kumer Das, and MM Hafizur Rahman. Bangla speech-to-text conversion using sapi. In *2012 International Conference on Computer and Communication Engineering (ICCCCE)*, pages 385–390. IEEE, 2012.
- [8] Md Rafiqul Islam, Ram Shanker Saha, and Ashif Rubayat Hossain. Automatic reading from bangla pdf document using rule based concatenative synthesis. In *2009 International Conference on Signal Processing Systems*, pages 521–525. IEEE, 2009.
- [9] Lawrence O’Gorman. Binarization and multithresholding of document images using connectivity. *CVGIP: Graphical Models and Image Processing*, 56(6):494–506, 1994.
- [10] Nobuyuki Otsu. A threshold selection method from gray-level histograms. *IEEE transactions on systems, man, and cybernetics*, 9(1):62–66, 1979.
- [11] Firoj Alam, Promila Kanti Nath, and Mumit Khan. Text to speech for bangla language using festival. Technical report, BRAC University, 2007.