2025, 10(40s) e-ISSN: 2468-4376

https://www.jisem-journal.com/

Research Article

Smart Vision: Developing an OCR-Based Speech Synthesis System with LabVIEW

V Venkataramanan^{1*}, Pankaj Mishra², Khushi Khanchandani³, Vijay Kapure⁴, Sarika Dharangaonkar⁵

1,2,3,5</sup>Assistant professor, Department of Information Technology, K J Somaiya School of Engineering, Somaiya Vidyavihar University,

Vidyavihar, Mumbai 400077

4Department of Applied Science and Humanities, Xavier Institute of Engineering, Mahim, Mumbai, India Drishti Parekh⁶, Vishnupriya Singh⁷, KalpShah⁸ 67.8 UG Student, Department of EXTC, D.J. Sanghvi College of Engineering, Vileparle (W), Mumbai-56 Corresponding Author*: reviewer.venkat@gmail.com

ARTICLE INFO

ABSTRACT

Received: 30 Dec 2024 Revised: 05 Feb 2025

Accepted: 25 Feb 2025

In today's world we are surrounded by one thing and that is data. This data needs to be recorded for various purposes and to do this manually is a huge task. Text is one of the forms of data and to record this text in computers is very difficult since typing all the text is time consuming and therefore inefficient. Hence, we can make this process more efficient by using Optical Character Recognition (OCR). OCR can help in conversion of this text data to speech signal. This efficient conversion to speech opens more doors for further applications. One of which is for blind and visually impaired, by which they can easily comprehend the text data without relying on third party for help. This paper aims to show the development of the OCR based speech synthesis system that is cost effective and easy to understand. The OCR converts the image to text which it then converts to speech. After the use of OCR, when the image is converted into speech, speech libraries present, such as in Microsoft SDK, are used for the conversion of text to speech. In this manner, the user, after scanning the image, is able to receive the information in form of speech. OCR training can be done in multiple languages and multiple fonts as well as handwritten font, according to the convenience of the user. Thus, software technology using Lab view is used for easy reception of written data using speech. This system has applications in finance industry, medical, home automation, etc. National Instrument's LabVIEW has been used to develop OCR based speech synthesis system

Keywords: LabVIEW, Natural Language Processing, Optical character recognition, Recognition, Speech, Synthesis

INTRODUCTION

This dream has now become real in machined reading. Text is present at schools in books, in newspapers and mostly in any kind of written pieces. That information helps the mobilisation of their cause but it is information that is not readily accessible to the blind and visually impaired. Therefore, the Optical Character Recognition (OCR) based speech synthesis system will complement the executive and discriminative control of the environment by the visually impaired as an able as a sighted person. Some contributions in this area were made in the past decade.

In fact, most published papers on general background studies have been performed in the domain of text from image or video. For the development, many various types of domain specific OCR applications have been created, such as receipt OCR; invoice OCR, Check OCR, legal billing document OCR, etc. In banking, sector it is widely used for check processing with no operator's involvement. The money is forwarded by reading the writing on the check when necessary. When such recordings are done and instead of writing it as most papers are done it saves the need to go through many papers files which can be quite tiresome for employees, besides, since these are computer database they also save on paper. Use of OCR is surging on multiple folds with the advancement in the technology. These are the purposes for which and owing to the ease provided by the LabVIEW, mentioned software has been used for OCR to Speech synthesis.

2025, 10(40s) e-ISSN: 2468-4376

https://www.jisem-journal.com/

Research Article

Technologies for the depiction of letters and numbers have changed allowing creation of OCR and speech synthesis for the weak-sighted. Some of the pioneering efforts in these categories, for instance the survey on the application of machine learning (ML) for text recognition have enabled the development of sophisticated OCR systems. For example, in [7], the authors propose the microstate structure relating to black hole entropy and compare it with OCR systems, in which each layer of the system adds to the accuracy of the final result.

In the recent past, integrating the OCR with the synthesis speech has become of great interest, more so in the area of assistive technology. Such systems are intended to translate between text-based data and people for whom accessing this data with more conventional methods is problematic. In this regard, the value of this technology cannot be overemphasized, given the increasing push toward the inclusion of learners and other individuals with disability, in education, social interaction as well as in everyday life.

These systems help the physically challenged and more specifically the visually impaired to read print materials on their own and thus increase their independence and quality of life. Besides this, the usage of these technologies does not remain limited to day-to-day purposes only; they have implications in fields of finance, health care, automation and many such areas where effectiveness of the data processing techniques and availability of information can play a crucial role. It is seen that the development of such systems in general is a part of several societal advancements in terms of technology that help in easing down the processes and thereby helps in the making of society 'inclusive'.

The integration of ML with OCR has actively improved the speech to text converting systems. Contemporary OCR systems have been made more accurate in identifying texts and other characters through employing improved ML algorithms that enables them to recognize named and new styled fonts together with recognizing text in different languages and even in a noisy background. Such improvements in performance have been made possible mainly by advancements in Deep Learning, where system can learn from large data sets and the system's performance can be made better iteratively. For instance, the convolutional neural networks (CNNs) have proved useful in the improvement of the character recognition whereas the recurrent neural networks (RNNs) have been useful in the enhancement of the context understanding.

In addition, advancement in Assistive Technology like the combination of OCR real-time processes has widened the opportunities for innovation. Whereas before, OCR could only work on batch scans and simple video frames, the current generation of OCR systems can transcribe text in Real Time from webcam feeds almost in real time. This capability is especially helpful in scenarios that include texts on screens or on spaces where there is continuous transformation. For instance, the visually impaired can use these systems to read text on electronic display or billboards encountering when in the public places.

It marks a tremendous influence on the society More often than not the implication of these advance brings long lasting change on the society. This way OCR based systems allow visually impaired people to get a deeper and more independent access to text based information in educational, working or social environments. This technology not only contributes to making the service more independent but also tends to prioritise the principles of inclusive information space. Thus, as the OCR and ML technologies develop further, it is anticipated that their utilization area will also grow in a greater extent and provide solutions that are more elaborate and adaptable for different situations.

LITERATURE SURVEY

Based on the history of OCR system development and utilization studies are discussed. Note that included are commercial systems developments. In relation to R & D two types are distinguished and that is matching of the template type and analysis of the structural type. [4] and [5] indicate that these two approaches are developing into each other and are also becoming integrated. OCR systems are categorized in to three generations and for each generation some of the OCR systems are selected and briefly described. Certain observations are made about recent approaches to OCR – for instance, neural nets and systems of expertise – to underscore some problems. As for any given author's expectations of the future scope and developments, the following points are described. [1]

2025, 10(40s) e-ISSN: 2468-4376

https://www.jisem-journal.com/

Research Article

Sample digits' regions are obtained from the image by labels, edge extraction and also augmentation. Since, the digits inside the images contain slope along with distortion, the digits are identified, after slope and distortion rectification are done. Hough Transformation is used to make correction of the distorted shape while for correction of slope only the numbers enclosed in the skew rectangles are truncated. In these experiments, a testing of a total of 1332 images of signboards was done with 11939 digits. A digit extraction rate of 99. 2 percent with a correct matching rate of 98 percent. 8% was obtained. [2]

That is why, binarization and image augmentation lay the foundation of the proposed methodology, combined with the right approach to the connected constituent analysis. Image binarization successfully carries out interior / open air scene images that contain shadows, unequal light intensity which could be little at times and high signal noise ratio. Final binaries are pre-processed with connected component analysis, based mainly on text areas. The recommended methodology improves the success rates of applying OCR at industrial level. Some experiments done in International Conference on Document Analysis and Recognition. (ICDAR) 2003 Robust Reading Competition seeming to support this approach. [3]

A symbol implies the presence of a fact, circumstance, or feature. Symbols are present at all the places in our life. They help to make our life simpler when we are accustomed to them. A at times they create difficulties for example, symbols or signs of a foreign country may not be understood by a tourist. In this paper issues of automated recognition of sign and its translation are discussed. They put forward a scheme which is capable of image acquisition and detection, which then recognizes the symbols in the acquired image and converts it into the desired language. They use an approach which is user based in system enhancement. This approach takes profits from human intelligence and leverages mental capacity. They are presently working on Mandarin symbol translation. They have currently developed a basic version that recognises Mandarin symbols from the camera and translates it into the desired language of the user. The sign translation, along with spoken language translation, can help tourists to overcome difficulties usually encountered in foreign land. This technology will also help a visually impaired person to increase their environmental alertness. [4]. Text present in a scene or a video footage gives important information about the index and also gives hints for video decoding. In this work, they propose algorithms for tracking and detection of text in the video. Their system uses a scale-space feature extractor whose output is given to an artificial neural network processor to detect text segments. Their text tracking system consists of two sections: a sum of squared difference (SSD) based section to find the starting position of text and a contour-based section to enhance the position. They conducted experiments with a variety of video footages to show that their system can detect and track text efficiently. [5]

A full OCR system for written Bengali, one of the most popular written scripts, is discussed in this paper. This problem isn't easy because (a) there are approximately three hundred simple, modified and composite letter shapes in this written script, (b) the letters of the word are topographically joined and Bengali is a conjugational language. In their method, the candidate image is grabbed by a flatbed scanner which is then subjected to bend correction, textual graphics differentiation, linear segmentation, zonal recognition, word and letter segmentation using common and some novel methods. The simple and modified letters which are approximately 75 in total and occupy round about 96% of the textual body, are recognized with a tree classifier. The composite letters are detected by a feature tree classifier and continued with sample-verification. The recognition is robust and basic wherein processing methods like diminishing, clipping and thinning is evaded. Letter unigram data is employed to increase the efficiency of the classifier. Various investigative methods are also employed to speed up sample verification. An error-correction system, which is based on dictionary, has been used too wherein different dictionaries are checked for parent root recognition and affixes that contain morphological information. [6]

OCR technology is improved by the evolution of high-powered desktop computing which allows for the development of recognition software which are powerful and can detect and read a various handwritten and printed texts. Yet it remains a difficult task to make a system that doesn't falter under difficult conditions and works with very high accuracy rates. The motto of the paper was to introduce the method to the developers and researchers in this field. Since South Indian languages have restricted facility, this field is still sought after and is needed in the society. In this paper they have presented the features of the scripts, methods used in detection techniques and comparison of the results of various optical recognition methods developed for various South Indian scripts. [7].

2025, 10(40s) e-ISSN: 2468-4376

https://www.jisem-journal.com/

Research Article

In this paper the OCR uses smart mobile phones which are based on windows operating system. This paper combines the functionality of a speech synthesizer and an OCR. Their aim was to develop a user-friendly mobile application which performs image to speech conversion using mobile phones only. The input to the OCR is an image, from which text gets recognized and is then converted into speech. This mobile application can be used in finances, judicial processes, home automation, etc. [8]

Extraction of data by just sound is a unique property. Speech is more efficient communication mode than text by itself because the visually impaired and blind can also understand it. In this paper they have tried to develop a more economical and practical OCR system. The speech synthesizer along with the OCR was developed using NI LabVIEW software. [9]

OBJECTIVES

This paper aims to show the development of the OCR based speech synthesis system that is cost effective and easy to understand. The OCR converts the image to text which it then converts to speech. After the use of OCR, when the image is converted into speech, speech libraries present, such as in Microsoft SDK, are used for the conversion of text to speech. In this manner, the user, after scanning the image, is able to receive the information in form of speech. OCR training can be done in multiple languages and multiple fonts as well as handwritten font, according to the convenience of the user.

METHODS

This section outlines the methodology used in developing the OCR-based speech synthesis system, detailing the process from image conversion to text and its subsequent transformation into speech. The detailed diagram of the Optical character recognition system is shown in Fig.1.

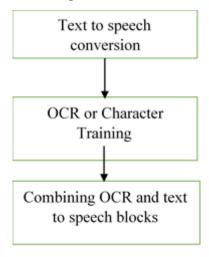


Fig 1. Bolock Diagram of an OCR System [10]

The first block deals with image acquisition and text recognition. This component is best suited for the capturing of the input image which can be acquired from a scanner or even a webcam. After it has been acquired the image goes through a number of processing procedures. It locates areas of interest or areas where there is probably text hence being known as the Region of Interest (ROI) [11] or region of interest. These areas are then taken through binarization and linearization processes which bring the image to black and white together with straightening of the text lines. After that, the text is analytically divided into the particular characters, words, and lines of the material being analysed. The technique for segmenting characters is by using the Vision Assistant tools to map each character to a definite value for slicing, while, for the word segmenting the function of slicing uses a standard measure to determine the correct position of the blanks between the words.

The second block illustrates the management of the important task of transforming the received text to speech. [12-13] This component probably uses a speech library with particular functions adjusted to the needs of the application. Starting with phonetic compilation of the text, the material is divided into phonemes so as to

2025, 10(40s) e-ISSN: 2468-4376

https://www.jisem-journal.com/

Research Article

appropriately pronounce it. Finally, audio processing algorithms produce the beat and voice of the text producing a speaking facility with which the written words are turned into audible vocal expressions.

At the beginning of the System Block Diagram, at the level of analysis, it gives the big picture of the flow. Initially, the input of an image is taken after which the image is subjected to several operations to improve it from the noise [14]. The system then recognizes the depicted text in the processed image and then the character recognition takes place to read the image text into operable form. This recognized text is then output, maybe to a display or a file, then passed through the text to speech synthesis module for conversion.

Using the local variables as means to connect these components of this system is another major consideration in its design as shown in Fig.2. and Fig.3.

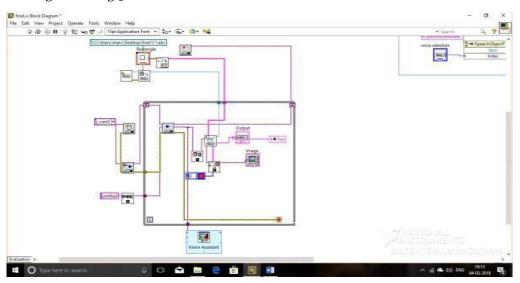


Fig.2: Sub VI for image acquisition and text detection

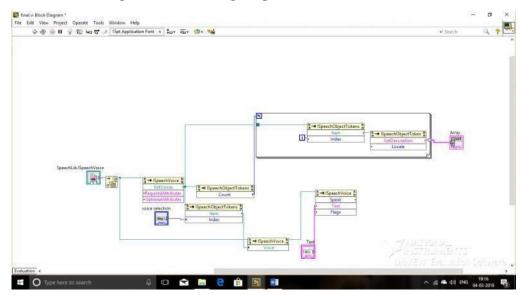


Fig.3: Sub VI for text to speech conversion

In general, there are two working principles for his project, one is related to OCR and other is for Speech Synthesis. For OCR, the input image is received from the scanner or from the webcam. The image is acquired using Image Acquisition (IMAQ) which is a tool in LabVIEW. Then from this image the text gets detected and analysed by respective blocks called ROI (Region of Interest) from the Sub VI [15]. On this detected text binarization and linearization is performed to get plain text. The text is segmented by character, word and line. Character segmentation is done using Vision Assistant by assigning each character in the text its assigned value. Word

2025, 10(40s) e-ISSN: 2468-4376

https://www.jisem-journal.com/

Research Article

segmentation is done by changing the minimum spacing distance between two words with a threshold value (for example 30). In Speech Synthesis, phonetic analysis is performed on the text to get the phonetic pronunciation of each word, after this audio processing is done to generate the voice output. Speech Synthesis is done using a speech library available. For interfacing these two concepts a local variable is created in the textbox of speech synthesis section, then this local variable is given to the OCR section which forwards it to Speech Synthesis section. Optical character recognition-based speech synthesis system functions by going through a number of complex processes after the acquisition of an image with the text [16]. Usually, the image is acquired by a scanner or the web camera then with the help of LabVIEW's IMAQ tools the image is processed. The first process of OCR is to locate and extract the image from the image plane, this can be done through Region of Interest identification. This step is important as it begins the text extraction part of the image, where by the focus is on the text and not the image. Next, the ROI obtained is followed by a process called binarization where the image is converted to black and white in order to have a better vision in text extraction. It then undergoes linearization in an attempt to rectify the ink presence of undulating lines by enhancing the recognition of odd characters.

After this, the text is catered further to characters, words and lines. Segmentation of the character is done using the Vision Assistant of the LabVIEW software in which each character is given the appropriate value for recognition. Sentence separation is adjusted using spacing threshold in order to have the system to differentiate between the words in a given sentence. After the text in the input image has been completely segmented and recognized, the system goes to the synthesis of speech. In speech synthesis phase, after recognizing the text, then phonemic analysis is performed using which the recognized text is analysed phonetically which means the conversion of text to the phonemes that are the basic units of sound. This phonetic data is then utilized to produce speech a speech library that already exist like the Microsoft SDK[17-18]. The conversion from text to speech is highly reliant on the OCR as the mistakes in recognizing text would lead to wrong pronunciations or wrong utterances.

RESULTS

In the case of the OCR based speech synthesis system, the performance has evidenced its feasibility applied to the objective of properly transforming the text written on the images to clear and proper speech. The capability of the system was tested through the use of printed text, handwritten text, and different font styles to check on the flexibility of the system. As seen in the functional analysis of the OCR module, the wording recognition did indeed exhibit considerable precision in the identification of characters on images regardless of the text quality as well as text size on the input images. The enhanced prospects of segmenting the text into characters, words, and lines also led the way to precise recognition of the text and subsequent conversion to speech. The practitioners of the research state that the speech synthesis module has achieved a clear and natural speech output, which means the service of the system is perfect for its intended applications as shown in Fig 4

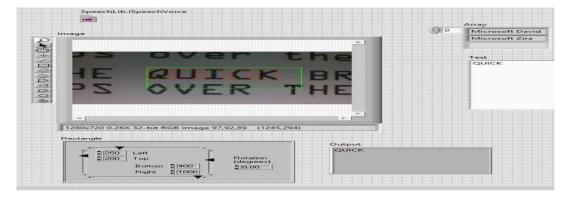


Fig. 4: Text detection front panel

The difficulties were observed during the testing phase of the given system, which improperly recognized, for instance, handwritten text and stylized fonts at times. However, I identified these following challenges which but were overcome through further improvements of the various steps in the image preprocessing techniques, for instance, the thresholds used for binarization and the algorithms used in the identification of Regions of Interest or

2025, 10(40s) e-ISSN: 2468-4376

https://www.jisem-journal.com/

Research Article

ROI. Nevertheless, the performance of the system was quite reasonable, with a rate of conversion from the recognized text to speech outputs that was generally impressive.

However, the fact that the system can cope with different languages and fonts, and is equally well suited to being fed with information from scanners and web cams, suggests that its scope is much wider. They also pointed at the need for integration to the extent that every module; from image acquisition to the actual speech output must operate in harmony for good performance. This system has proved to work as a good proof that LabVIEW could be used for the development of hard-OCR and speech synthesis applications in order to support the future improvements and other more complex applications in diverse areas.

The effectiveness of the proposed OCR-based speech synthesis system can be also explained by the special preprocessing of input images, which is similar to black hole studies where each stage of the work is significant and important to receive accurate near-horizon geometries. The careful division of the textual information into characters, words, and lines as geometrical features in the higher-dimensional spacetimes highlights the necessity of accurate analysis of the source images for getting high results as shown in Fig 5

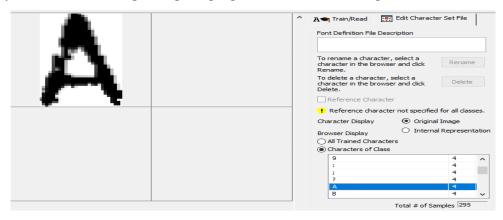


Fig. 5: Optical Character recognition

Real-time examination was also done to determine the effectiveness of the system where the text was directly obtained from the webcam. This element of extra live typing real time text capture means that other factors such as the condition and direction of text images and even the resolution can be affected. However, the system could intimidate the patients and also have occasional wrong interpretations of signs and symptoms, thus revealing its efficiency and versatility in spite of the challenges. On the use of real-time text-to-speech conversion it was found to be useful in that it produced immediate audio feedback and this plays an important role in examples like the assistive technology for the blind. Indeed, the fact that the system can accept complex inputs and maintain good rate of convergence even without compromising much speed adds to the thought that the model can be 'deployed' for use in real life given that consistency is always key in real-world problems.

DISCUSSION

The OCR-based speech synthesis system offers an efficient and accessible solution to convert text into speech, making it easier for visually impaired individuals to access information. By utilizing OCR, the process of converting written data into speech becomes more streamlined. The use of speech libraries, like Microsoft SDK, enhances the system's functionality. Additionally, the ability to train OCR in multiple languages, fonts, and handwritten scripts allows for a highly customizable solution, catering to the user's needs.

CONCLUSION

We have successfully performed the image OCR to speech Synthesis system using LabVIEW 7.1. The image to be synthesized is taken as the input. The various parts of OCR are performed on the image and the image is then converted to a text file. This builds up one VI of the software. The text file is converted into the speech using Microsoft speech Library (5.1). This formulates another VI of the software. There were many obstacles faced during the completion of this project. This is evident from the integration of image processing, character recognition, and text to speech wherein the platform is put into use. Various input formats handled by our system, as well as the

2025, 10(40s) e-ISSN: 2468-4376

https://www.jisem-journal.com/

Research Article

conversion of text to spoken words, can be utilised in various accessibility and automation solutions. The project also focusses on the idea of modularity, when the complex system is built, the problem can be located easier and it is easier to make change to it. In general, this work could contribute to the subsequent investigations in the domain of helping technologies and the automation of the document workflow.

REFRENCES

- [1] S. Mori, C. Y. Suen, and K. Yamamoto, "Historical review of OCR research and development," Proceedings of the IEEE, vol. 80, no. 7, pp. 1029–1058, Jul. 1992, doi: 10.1109/5.156468.
- [2] T. Yamaguchi, Y. Nakano, M. Maruyama, H. Miyao, and T. Hananoi, "Digit classification on signboards for telephone number recognition," in Seventh International Conference on Document Analysis and Recognition, 2003. Proceedings., vol. 1, pp. 359–363, 2003, doi: 10.1109/icdar.2003.1227689.
- [3] T. Yamaguchi and M. Maruyama, "Character extraction from natural scene images by hierarchical classifiers," in Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004., vol. 2, pp. 687–690, 2004, doi: 10.1109/icpr.2004.1334352.
- [4] J. Yang, J. Gao, Y. Zhang, X. Chen, and A. Waibel, "An automatic sign recognition and translation system," in Proceedings of the 2001 Workshop on Perceptive User Interfaces, pp. 1–8, Nov. 2001, doi: 10.1145/971478.971490.
- [5] H. Li, D. Doermann, and O. Kia, "Automatic text detection and tracking in digital video," IEEE Transactions on Image Processing, vol. 9, no. 1, pp. 147–156, 2000, doi: 10.1109/83.817607.
- [6] B. B. Chaudhuri and U. Pal, "A complete printed Bangla OCR system," Pattern Recognition, vol. 31, no. 5, pp. 531–549, Mar. 1998, doi: 10.1016/s0031-3203(97)00078-2.
- [7] S. M. Azizul Hakim and Asaduzzaman, "Handwritten Bangla numeral and basic character recognition using deep convolutional neural network," in 2019 International Conference on Electrical, Computer and Communication Engineering (ECCE), pp. 1–6, Feb. 2019, doi: 10.1109/ecace.2019.8679243.
- [8] T. Sato, T. Kanade, E. K. Hughes, M. A. Smith, and S. Satoh, "Video OCR: indexing digital news libraries by recognition of superimposed captions," Multimedia Systems, vol. 7, no. 5, pp. 385–395, Sep. 1999, doi: 10.1007/s005300050140.
- [9] R. S. Kunte and R. D. S. Samuel, "A bilingual machine-interface OCR for printed Kannada and English text employing wavelet features," in 10th International Conference on Information Technology (ICIT 2007), Dec. 2007, doi: 10.1109/icoit.2007.4418296.
- [10]T. Dutoit, "High quality text-to-speech synthesis: a comparison of four candidate algorithms," in Proceedings of ICASSP '94. IEEE International Conference on Acoustics, Speech and Signal Processing, vol. i, pp. I/565-I/568, 1994, doi: 10.1109/icassp.1994.389231.
- [11] "A Real Time Text Detection & Recognition System to Assist Visually Impaired," International Journal of Science and Research (IJSR), vol. 6, no. 7, pp. 2112–2115, Jul. 2017, doi: 10.21275/art20175527.
- [12]J. D. Sanghavi, A. M. Shah, S. S. Rane, and V. Venkataramanan, "Smart traffic density management system using image processing," in Proceedings of International Conference on Wireless Communication, pp. 301–312, 2018, doi: 10.1007/978-981-10-8339-6_33.
- [13]M. A. Rahiman and M. S. Rajasree, "A detailed study and analysis of OCR research in South Indian scripts," in 2009 International Conference on Advances in Recent Technologies in Communication and Computing, vol. 84, pp. 31–38, 2009, doi: 10.1109/artcom.2009.45.
- [14]S. K. Singla and R. K. Yadav, "Optical character recognition based speech synthesis system using LabVIEW," Journal of Applied Research and Technology, vol. 12, no. 5, pp. 919–926, Oct. 2014, doi: 10.1016/s1665-6423(14)70598-x.
- [15] V. Venkataramanan, G. Kavitha, M. R. Joel, and J. Lenin, "Forest fire detection and temperature monitoring alert using IoT and machine learning algorithm," in 2023 5th International Conference on Smart Systems and Inventive Technology (ICSSIT), pp. 1150–1156, Jan. 2023, doi: 10.1109/icssit55814.2023.10061086.
- [16] A. Sharma, A. Srivastava, and A. Vashishth, "An assistive reading system for visually impaired using OCR and TTS," International Journal of Computer Applications, vol. 95, no. 2, pp. 13–18, Jun. 2014, doi: 10.5120/16566-6231.

2025, 10(40s) e-ISSN: 2468-4376

https://www.jisem-journal.com/ Research Article

[17] A. Sharma, A. Srivastava, and A. Vashishth, "An assistive reading system for visually impaired using OCR and TTS," International Journal of Computer Applications, vol. 95, no. 2, pp. 13–18, Jun. 2014, doi: 10.5120/16566-6231.

[18]T. Khete and A. Bakshi, "Autonomous assistance system for visually impaired using Tesseract OCR & gTTS," Journal of Physics: Conference Series, vol. 2327, no. 1, p. 012065, Aug. 2022, doi: 10.1088/1742-6596/2327/1/012065.