**Abstract**

Cloud computing and Optical Character Recognition (OCR)
technology has become an extremely attractive area of research over the last few years. Sometimes it is difficult to retrieve text from the image because of different size, style, orientation and complex background of image. There is need to convert paper books and documents into text. OCR is still imperfect as it occasionally mis-recognizes letters and falsely identifies scanned text, leading to misspellings and linguistics errors in the OCR output text. A cloud based Optical Character Recognition Technology was used. This was powered on RunPod cloud in form of a Web Application Programming Interface that load images to an Optical Character Recognition server, process with necessary recognition, export parameters, and obtain the results of the processing. The key idea is to bring together the advantages of the Optical Character Recognition technology and cloud computing in one place in other to enable quicker access and faster turn out time, processing period, and increasing efficiency across the board for application users. The methodology adopted is the object oriented methodology. This was achieved using html, css and javascript programming language.

**General Terms:**

Cloud computing: Cloud computing is on-demand access, via the internet, to computing resources—applications, servers (physical servers and virtual servers), data storage, development tools, networking capabilities, and more—hosted at a remote data centre managed by a cloud services provider (or CSP). The CSP makes these resources available for a monthly subscription fee or bills them according to usage.

**OCR:**
Optical character recognition (OCR) is sometimes referred to as text recognition. An OCR program extracts and repurposes data from scanned documents, camera images and image-only pdfs. OCR software singles out letters on the image, puts them into words and then puts the words into sentences, thus enabling access to and editing of the original content. It also eliminates the need for manual data entry.

**RunPod Cloud:** RunPod provides two cloud computing services: Secure Cloud and Community Cloud. Secure Cloud runs in T4 data centres by our trusted partners. Our close partnership comes with high-reliability w/ redundancy, security, and fast response times to mitigate any downtimes.

**Keywords:** Cloud computing ,Optical character Reader, RunPod cloud, Image, text.

**Introduction**
There is a paradigm shift in the manner people learn, trade, communicate and share knowledge with the drastic introduction of modern computers into every aspect of life The dawn of digital computers has made it unavoidable that everything processed by the digital computer must be processed in digital form. For example, most famous libraries in the world like the Boston public library has over 6 million books, for public consumption, inescapably has to change all its paper-based books into digital documents in order that they could be stored on a storage drive. It can also be projected that every year over 200 million books are

being published, many of which are disseminated and published on papers [6].Thus, for several document-input jobs, Optical Character Recognition is the utmost economical and swift process obtainable. Every year, the innovation liberates storage space once meant for file organisers and boxes loaded with paper records. Before OCR can be utilized, the source material must be scanned with the use of an optical scanner to read the page as bitmaps.(patterns of dots).. The Optical Character Recognition software then processes these scans to differentiate between images and text and determine what letters are represented in the light and dark areas. —Older OCR systems match these images against stored bitmaps based on specific fonts. The hit-or-miss results of such pattern-recognition systems helped establish OCR's reputation for inaccuracy.

**Review of Optical Character**

 Recognition Technologies- Tesseract Tesseract is an optical character recognition engine for

various operating systems. It is free software, released under the Apache License, Version 2.0, and development has been sponsored by Google since 2006. The Tesseract engine was originally developed as proprietary software at Hewlett Packard labs in Bristol, England and Greeley, Colorado between 1985 and 1994, with some more changes made in 1996 to port to Windows, and some migration from C to C++ in 1998. A lot of the code was written in C, and then some more was written in C++. Since then all the code has been converted to at least compile with a C++ compiler.

**GOCR**: GOCR claims it can handle single-column sans-serif fonts of 20–60 pixels in height. It reports trouble with serif fonts, overlapping characters, handwritten text, heterogeneous fonts, noisy images, large angles of skew, and text in anything other than a Latin alphabet. GOCR can also translate barcodes.

**OCRopus**; OCRopus is a free document analysis and optical character recognition (OCR) system released under the Apache License, Version 2.0 with a very modular design through the use of plugins. These plugins allow OCRopus to swap out components easily. OCRopus is currently developed under the lead of Thomas Breuel from the German Research Centre for Artificial Intelligence in Kaiserslautern, Germany and is sponsored by Google. OCRopus is developed for Linux; however, users have reported success with OCRopus on Mac OS X and an application called TakOCR has been developed that installs OCRopus on Mac OS X and provides a simple droplet interface.
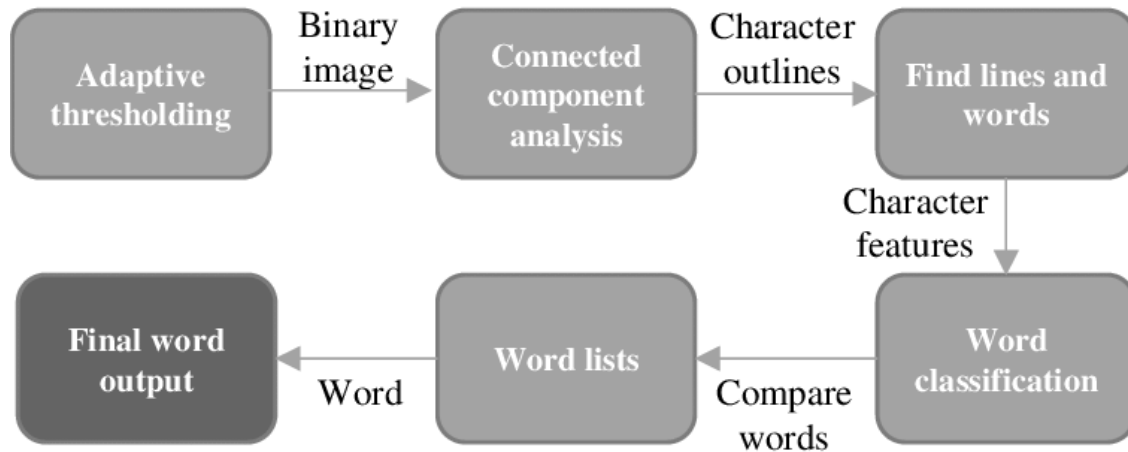
## THE EXISTING SYSTEM

Over the years Optical Character Recognition technology has improved speedily. However there are limitations with respect to the source materials and character formatting. Most documents formatting are lost during text scanning, except for paragraph marks and tab stops.Sometimes, the output from a finished text scan will be a single column editable computer

file. In most cases this computer file will at all times require.proofreading and spellchecking and in addition to reformatting to the required final layout. Using one of the open source and free OCR software for instance, Tesseract OCR is an elegant engine with various layers. It works in step by step manner. The first step in Named as adaptive thresholding, and converts the image into binary images. Second step is to do the connected component analysis of the image, which does the task of extracting character outlines.

After this the outlines extracted from image are converted into Blobs (Binary Long Objects). It is then organized as lines and regions and further analysis is for some fixed areas. After extraction the extracted components are chopped into words and delimited with spaces. Recognition in text then starts which is a two pass process.

Tesseract was originally designed to recognize English text only and there are limits as to the range of languages it can handle. It can only handle left-to-right languages. While you can get something out with a right-to-left language, the output file will be ordered as if the text were left-to-right. Top-to-bottom languages will currently be hopeless. It is unlikely to be able to handle connected scripts like Arabic. It will take some specialized algorithms to handle this case, and right now it doesn't have them. Moreover, it is likely to be so slow with large character set languages (like Chinese) that it is probably not going to be useful. There also still needs to be some code changes to accommodate languages with more than 256 characters.

In Figure 1, the first part is when attempt to recognize each word is made. Each satisfactory word is accepted and second pass is started to gather remaining words. This brings in the role of adaptive classifier. The adaptive classifier then will classify text in more accurate manner. The adaptive classifier needs to be trained beforehand to work accurately. When the classifier receives some data, it has to resolve the issues and assign the proper place of the text.
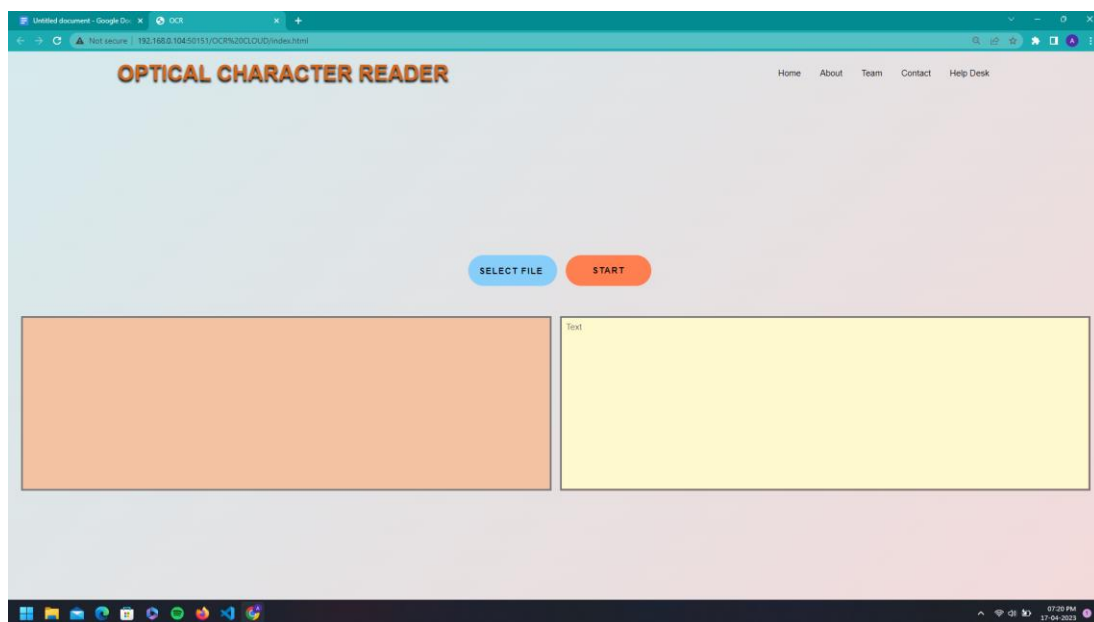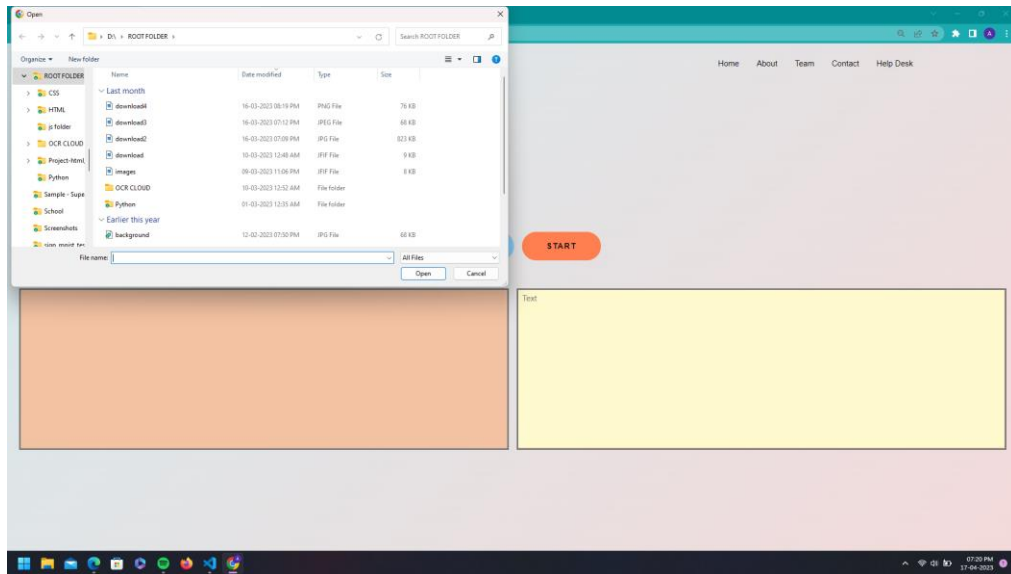
**System Architecture of the Existing SystemTesseract OCR Architecture**
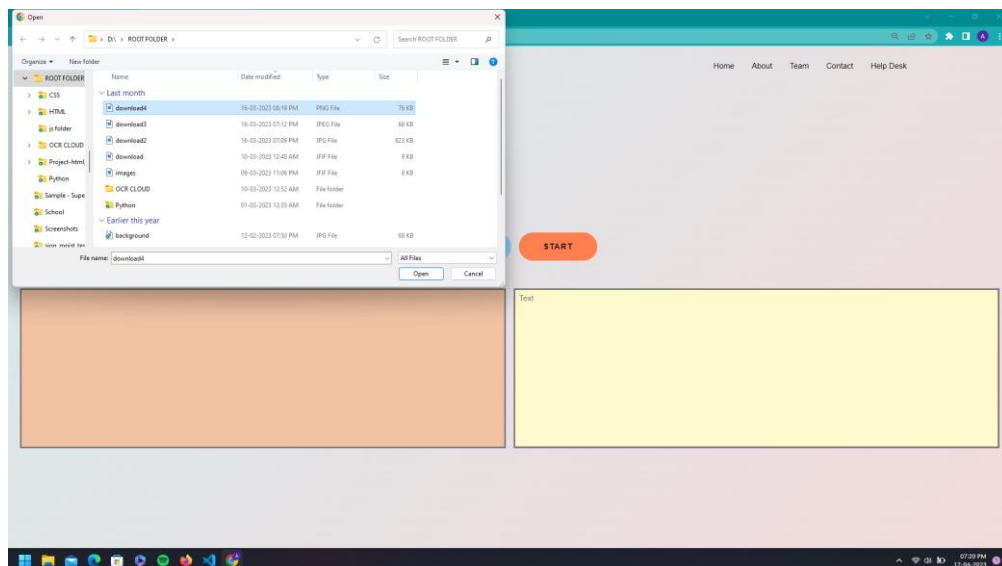
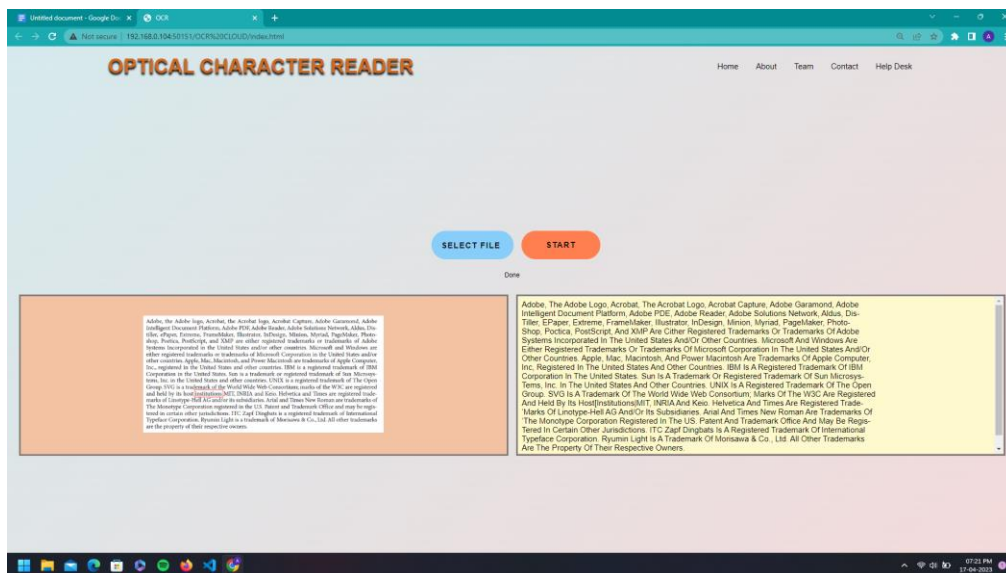## IMPLEMENTATION AND RESULT DISCUSSIONS

Part1:
The system was implemented and tested, the JavaScript language uses a runtime environment. Also Tesseract uses OCR Development library. The user interface forms where users, Allow selecting a file to identify, the file's jpeg or png name, Open File , and start.

**Part 2:** After that the correct image file is identified for the optical character reader. For example the image file can be in jpeg or png.

## CONCLUSION:

Cloud computing and OCR technology can greatly enhance document management and workflow automation. By combining cloud computing and OCR, users can achieve faster turnaround times, increased productivity, and cost savings. Cloud-based OCR solutions automate document processing, reduce manual data entry, and improve accuracy. This integration represents a powerful solution for businesses looking to optimize their document management processes and improve their bottom line.