# Multilingual Hardcoded Subtitle Extractor

Amarjith V

*Department of Computer Science and Engineering*

*Sree Buddha College of Engineering*

Alappuzha, INDIA

amarvijay2580@gmail.com

Anaswara Anil

*Department of Computer Science and Engineering*

*Sree Buddha College of Engineering*

Pattoor, INDIA

anaswaraanil168@gmail.com

Anju Viswam

*Department of Computer Science and Engineering*

*Sree Buddha College of Engineering*

Pattoor, INDIA

anjuviswam30@gmail.com

Aravind KM

*Department of Computer Science and Engineering*

*Sree Buddha College of Engineering*

Pattoor, INDIA

aravindkm2001@gmail.com

*Abstract-* **The Multilingual Hardcoded Subtitle Extractor revolutionizes the way language barriers are overcome in digital content. Through the utilization of state-of-the-art Optical Character Recognition (OCR) technology, it meticulously extracts text from video frames that contain hardcoded subtitles. This enables a seamless translation experience into the user's preferred language. The entire process is streamlined, resulting in the creation of a SubRip (SRT) file that facilitates global content accessibility and localization efforts. Whether it is individual content creators looking to expand their audience or extensive production teams aiming for international distribution, this extractor offers a versatile toolset to cater to diverse needs and preferences. With its user-friendly interface and automation capabilities, the Multilingual Hardcoded Subtitle Extractor simplifies the complexities of creating multilingual content. It empowers users of all skill levels to effortlessly navigate the translation process, fostering a more inclusive digital environment that celebrates linguistic diversity. By making multilingual subtitles accessible to all, this extractor not only enhances global content accessibility but also promotes cultural exchange and understanding on a global scale. It paves the way for a future where language is no longer a barrier to universal enjoyment.**

*Keywords – OCR, EasyOCR*

## I. INTRODUCTION

In a world increasingly interconnected, multimedia content holds immense potential to educate, entertain, and inspire across geographical and linguistic boundaries. However, a persistent challenge confronts both creators and consumers alike: hardcoded subtitles embedded directly into video frames. These subtitles present significant hurdles for non-native speakers, language learners, and the hearing-impaired community, limiting their ability to fully engage with and comprehend content. Consider the frustration of attempting to follow along with a foreign film or instructional video in an unfamiliar language, hindered by subtitles that cannot be easily translated. This frustration not only limits the global accessibility of countless videos but also impedes cross-cultural communication and inclusivity.

Optical Character Recognition (OCR) emerges as a promising solution to this problem, enabling the extraction of hardcoded subtitles from videos by converting text-containing images into machine-readable formats. This foundational technology, integral to computer vision, addresses various practical challenges such as recognizing handwritten characters, license plates, and street numbers. The application of OCR in video processing involves importing essential libraries like OpenCV, identifying frames containing text, and subsequently extracting subtitles using the EasyOCR API. Additionally, translation services such as Google Translate or the OpenAI API can be utilized to translate extracted subtitles into desired languages, enhancing accessibility for diverse audiences. Key features of this approach include its multilingual extraction capabilities, automation for efficiency, and a user-friendly interface. Its potential applications extend to global content localization, language learning support, and improved accessibility for individuals with hearing impairments, fostering a more inclusive and interconnected global multimedia landscape where language ceases to be a barrier to universal understanding.

## II. LITERATURE REVIEW

In [1] An Adaptive Thresholding Algorithm-Based Optical Character Recognition System for Information Extraction in Complex Images, study focuses on image processing and data segmentation techniques using an adaptive thresholding algorithm. The images used were sourced from the internet in various formats and resolutions. A custom adaptive algorithm was applied to unify the complex backgrounds of the images, using the Gaussian thresholding algorithm. This algorithm dynamically generates the block-size to apply threshing, ensuring that images are processed area-wise. The data was sourced from web repositories, and the algorithm was constantly tweaked to improve performance. The image segmentation method used adaptive thresholding, which converts the image to greyscale and uses a custom adaptive thresholding algorithm to separate essential features based on minimal pixels. The threshed image was then parsed to

Tesseract, a commercial OCR tool, for easy text separation. The Image to Grey Submodule converts the input image to grayscale, while the Set Block Size Submodule determines the pixel usage per area of the image. The Thresh image Submodule returns the segmented image, which is then parsed to the Extract Text module for character extraction.

In [2] Based on improved edge detection algorithm for English text extraction and restoration from colour images, Image restoration techniques can be used in image translation systems to improve understanding of image content. By changing the background and adding new text, this technology can enhance the re-use value of images. Traditional edge detection algorithms struggle with image background complexity, making this paper a promising research topic. The paper improves the traditional edge-based detection method by fixing the 5 text parts of a colour image and proposing a new edge detection operator. The proposed algorithm effectively extracts accurate edge information and has strong anti-noise performance, making full use of image information and improving text extraction accuracy.

The FAST technique for corner identification in pictures for text recognition—which may be employed even with noise or blurs—is covered in the work [3] Image processing application in character recognition. For machine vision applications like structure from motion, optical flow computation, object tracking, and 3D scene reconstruction, this approach is crucial. One kind of feature point that is often employed is the corner, which drastically lowers data and offers a trustworthy constraint for calculating picture displacements. With applications in motion detection, video tracking, 3D modelling, image registration, and object identification, corner detection is used in image analysis to extract features. The benefits of this method over other algorithms for character recognition are also covered in the study. In the digital era, optical character recognition (OCR) technology is crucial in character recognition, leading to error-free output and improved image analysis and digitization accuracy.

## III. METHODOLOGY

The Multilingual Hardcoded Subtitle Extractor stands as a robust solution meticulously crafted to tackle the complexities of extracting multilingual subtitles from videos, especially those entrenched as hardcoded text within frames. Operating within a well-defined workflow, the system seamlessly incorporates optical character recognition (OCR) technologies, enabling the precise identification, extraction, and organization of subtitles. The system's operational framework, depicted in Figure 1.1, illustrates its structured approach and underscores its efficiency in handling multilingual subtitle extraction tasks.

### A. Input Processing

The input processing phase serves as the pivotal entry point into the Multilingual Hardcoded Subtitle Extractor system, where a video file containing hardcoded subtitles undergoes meticulous examination. Leveraging the capabilities of OpenCV, an open-source software library renowned for its prowess in computer vision and machine learning, the input processing stage executes a series of image processing operations such as greyscale conversion and scaling. These operations lay the groundwork for preliminary analysis and metadata extraction, deciphering crucial details including the video's format, resolution, and encoding. Moreover, this phase adeptly identifies potential challenges such as variations in frame rates or video quality, crucial factors that may influence subsequent processing steps.

The essence of this phase lies in comprehensively understanding the intricacies of the video's characteristics, thereby facilitating efficient and accurate subtitle extraction throughout the system's workflow. Following the initial analysis, the video undergoes decoding to extract individual frames, priming them for meticulous frame-by-frame analysis in subsequent steps. Notably, this phase may incorporate preprocessing measures aimed at enhancing video quality, standardizing frame rates, or rectifying any anomalies, ensuring a homogenized input conducive to the system's downstream operations. Designed with adaptability in mind, the input processing phase accommodates an array of video formats, encompassing mp4, mkv, AVI, and more. This versatility fortifies the system's robustness, ensuring seamless integration into diverse content creation workflows while upholding the requisite standards of accuracy and efficiency essential for effective subtitle extraction.
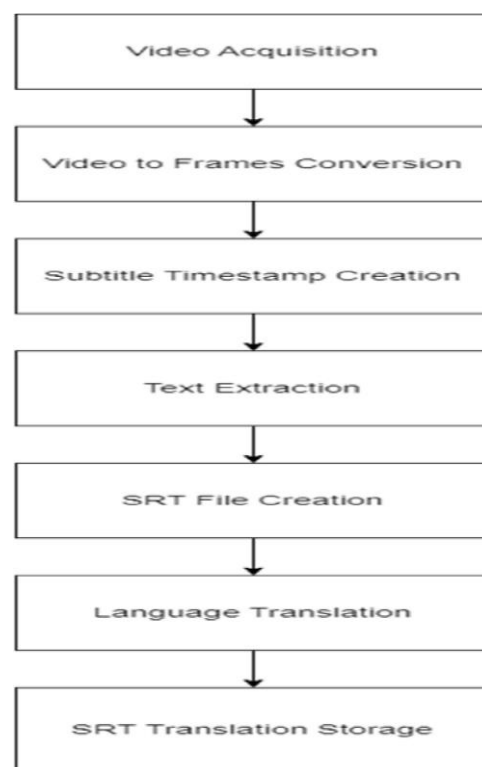


Fig. 1   Flowchart of working model

### B. Frame Detection with EasyOCR

During this pivotal stage, the Multilingual Hardcoded Subtitle Extractor harnesses the formidable capabilities of EasyOCR, an advanced optical character recognition framework renowned for its prowess in identifying text-

containing frames within videos, with a particular emphasis on hardcoded subtitles. Employing cutting-edge deep learning techniques, EasyOCR systematically scrutinizes each frame, showcasing remarkable adeptness in discerning text across a spectrum of fonts and styles. This framework exhibits adaptability to varying text sizes and orientations, thereby demonstrating resilience in the face of the multifaceted presentation of subtitles. Through meticulous character-level and word-level recognition, EasyOCR precisely distinguishes subtitles, adeptly navigating through challenges posed by line breaks or special characters. The process incorporates robust error-handling mechanisms tailored to mitigate obstacles such as low contrast or complex backgrounds, thereby optimizing detection accuracy. The integration of EasyOCR for frame detection augments the overall efficiency and adaptability of the system, endowing it with the capability to proficiently handle diverse linguistic contexts and presentation intricacies prevalent in video subtitles. By leveraging the EasyOCR framework, the system intelligently identifies frames housing text within the video, with an unparalleled focus on robust text recognition. EasyOCR emerges as a cornerstone component, pivotal in the identification and isolation of frames featuring hardcoded subtitles, thus underscoring its indispensable role in the Multilingual Hardcoded Subtitle Extractor's workflow.

## C. Image Export

During the Image Export phase, frames recognized to contain text, particularly hardcoded subtitles, undergo individual extraction as images. Each extracted image is assigned a naming convention incorporating temporal information, denoting the starting and ending times within the video. This meticulous organization ensures chronological order and simplifies temporal synchronization, vital for subsequent processing stages. The exported images serve as the foundation for subsequent optical character recognition (OCR), potentially accompanied by preprocessing steps aimed at enhancing clarity. This phase plays a pivotal role in streamlining the preparation of frames, ensuring they are primed for efficient and accurate multilingual subtitle extraction. By contributing to the overall precision and organization of the system, the Image Export phase reinforces the robustness and efficacy of the subtitle extraction process.

## D. Text Extraction with EasyOCR API

Upon transitioning to the Text Extraction stage, our framework seamlessly integrates the EasyOCR API, an advanced optical character recognition service renowned for its proficiency in identifying text across diverse languages. Employing state-of-the-art algorithms, the API meticulously analyzes each image, adeptly extracting text even amidst variations in fonts, sizes, or orientations. Robust error-handling mechanisms embedded within the API further bolster accuracy, enabling it to tackle challenges posed by distorted or complex text with finesse. The resultant output manifests as a machine-readable representation of the extracted content, primed for organization and synchronization of subtitles. By leveraging the EasyOCR API, our framework ensures the precise extraction of multilingual subtitles, thereby enhancing the overall efficiency of the process. The API's versatility and robust error-handling

capabilities contribute significantly to the framework's efficacy, empowering it to adeptly navigate through diverse linguistic contexts and textual intricacies.

## E. SRT File Generation

In the SRT File Generation phase of the Multilingual Hardcoded Subtitle Extractor, a foundational step is taken to meticulously organize the extracted content and synchronize it with the video timeline. Building upon the machine-readable representation obtained from the Content Extraction stage, this phase involves the creation of a SubRip (SRT) file. Each segment of extracted content is precisely aligned with the corresponding frames' start and end times within the video. This meticulous alignment ensures that the subtitles are accurately synchronized with the video playback, preserving the sequential flow and contextual accuracy of the extracted content. The SRT file serves as a standardized format for storing subtitles, encapsulating not only the textual content but also the timing information dictating when each subtitle segment should appear and disappear during video playback. The SRT file generated in this phase transcends mere subtitles; it represents a structured transcript of the video content, providing content creators with a comprehensive resource for seamlessly integrating multilingual subtitles into their videos. This final output fosters a more inclusive and universally accessible viewing experience, transcending language barriers and enhancing the overall accessibility of multimedia content.

## F. Language Translation

During the translation phase of the hardcoded subtitle extractor, users are empowered to select their desired target language from a diverse array of options made available through integrated language libraries or APIs. The translation process is executed utilizing a specific algorithm or tool, prominently featuring the engine employed, such as Google Translate, alongside relevant settings tailored to optimize outcomes. Notably, meticulous attention is devoted to nuances to ensure that the translated subtitles are not only accurate but also culturally sensitive, thus enhancing the overall viewing experience. The formatting of translated subtitles is meticulously managed to uphold consistency with the video, prioritizing readability and precise timing alignment. Rigorous quality checks are conducted, encompassing verification and editing procedures aimed at refining translation precision. Subsequently, the export procedure seamlessly converts the translated subtitles into SRT files, adhering to a predefined structure and encoding standards to ensure compatibility with various video players. Incorporating robust error handling mechanisms, the system is adept at addressing translation issues and promptly providing error messages to assist users throughout the process. Users benefit from a streamlined experience, effortlessly selecting their preferred language from a dropdown menu within the interface. This initiates the conversion process whereby the original subtitle extracted from the video is seamlessly translated into the chosen language and saved as an SRT file. Consequently, users can enjoy content in their language of choice, enriching their viewing experience and fostering inclusivity.

---

### Algorithm

---

1. Load the video using OpenCV

2. Initialize the ROI of the video and convert it into grayscale image using OpenCV.

2.1. Provide required parameters like language for EasyOCR.

3. Iterate through the frames and utilize EasyOCR for text detection in frame.

3.1. Pre-processing involves adjusting contrast, brightness, and other image parameters.

3.2. The CRAFT (Character Region Awareness for Text Detection) algorithm is employed to detect text regions within the pre-processed image.

3.3. A combination of ResNET, LSTM and CTC is used to recognize individual characters and words within the detected text regions.

4. Extract the frames along with its timestamp.

5. The extracted frames are saved as an image along with its timestamp as file name.

6. Text is extracted from the image using EasyOCR and extracted to srt file.

6.1. A Greedy Decoder translates recognized characters into coherent words and sentences.

7. Then the subtitle can be translated to other languages using DeepL or OpenAI API.

8. The translated subtitle is then saved to an srt file.

---

*G. Output Delivery*

In the initial phase, frames containing subtitles are meticulously extracted from the video and saved as images, with their respective start and end times serving as filenames. These images, coupled with their timestamps, constitute the groundwork for subsequent text extraction processes. The culmination of this phase yields a SubRip (SRT) file housing subtitles extracted from the video, encompassing multiple languages present within the content. This SRT file emerges as an invaluable asset for content creators, streamlining the localization of content on a global scale and enhancing accessibility for viewers worldwide. Each subtitle within the SRT file is accompanied by its corresponding timestamp, facilitating seamless synchronization with the video when played on media players like VLC. This SRT file, synonymous with SubRip subtitle files, plays a pivotal role in the translation of subtitles into the user's preferred language. Through an intuitive user interface, individuals can effortlessly select their desired language, triggering the translation of subtitles utilizing a dedicated translation API. The resulting output manifests as a freshly translated SRT file, seamlessly compatible with the original video across various media players, including VLC. Consequently, users are empowered to immerse themselves in their favorite movies or shows in their desired language, thereby enriching their viewing experience and bolstering accessibility.

## IV. MODULES

*A. Easy OCR*

EasyOCR stands out as a Python-based open-source library tailored for Optical Character Recognition (OCR) tasks. Its streamlined functionality facilitates the extraction of text from images, supporting multiple languages and popular formats like JPEG and PNG. Leveraging pre-trained deep learning models, including Transformer-based architectures, EasyOCR excels in achieving high accuracy in text recognition across various sources. Notably, its intuitive API simplifies the integration of OCR functionality into applications, ensuring accessibility for developers. Despite its lightweight nature, EasyOCR delivers competitive performance, proving invaluable for tasks such as document processing, data extraction, and automation, where precise text recognition is paramount. This comprehensive solution serves as a cornerstone in simplifying text extraction from images, offering a versatile toolkit for diverse OCR needs.

*B. Open CV*

OpenCV, renowned as a versatile open-source tool for computer vision tasks, is aptly dubbed the "Computer Vision Library." Within the realm of the multilingual hardcoded subtitle extractor, OpenCV emerges as an indispensable asset. This library streamlines the video processing pipeline, facilitating the extraction of frames and enhancing text detection capabilities. When integrated with an OCR library like EasyOCR, OpenCV's synergy amplifies the text recognition process, thus fortifying the overall efficiency of the system. The versatility of OpenCV extends to its robust support for various image and video formats, coupled with an extensive array of documentation and vibrant community support. This comprehensive ecosystem proves instrumental in refining subtitle extraction accuracy, particularly in multilingual scenarios, where the flexibility of OpenCV enables the implementation of language-specific filtering techniques. Moreover, the integration of OpenCV offers access to a plethora of image processing techniques, including black and white conversion, cropping, and more, thereby enriching the system's capabilities and bolstering its efficacy in handling diverse video content.

*C. Pysrt*

Pysrt emerges as a versatile Python library meticulously crafted for seamless interaction with SubRip (.srt) subtitle files. Tailored to cater to a myriad of needs, this library furnishes users with a rich array of functionalities, encompassing parsing, creation, and manipulation of subtitles adhering to the SubRip format. With Pysrt at their disposal, users gain the capability to effortlessly read existing subtitle files, extract precise timing information, and dynamically modify subtitles through programmatic interventions. This library further extends its utility by facilitating common subtitle operations such as temporal adjustments, merging, splitting, and beyond, thereby streamlining the entire subtitle handling process. The advent of Pysrt empowers developers to automate tasks associated with subtitle management, rendering it indispensable for projects spanning video processing, content localization, and any endeavor necessitating manipulation of subtitle files. Pysrt epitomizes simplicity and efficiency, offering developers a

convenient interface for both reading and writing SubRip subtitle files, a functionality not inherently supported by Python. Its seamless integration into Python workflows enhances productivity and enables streamlined manipulation of subtitle data, thereby amplifying its utility across diverse applications.

### D. OpenAI API

The OpenAI API with models like GPT 3.5 Turbo can be used for translation of extracted subtitles from the video to a wide range of languages. With different prompts we can get more customised translation according to the type of video or content that we are providing. We can also provide glossary or dictionary for a specific movie or series to get a more accurate result. The downside of using OpenAI API is that it's inconsistent and may not be accurate all the time. Compared to other translation service APIs like Google translate and DeepL its provides a better result overall.

## V. RESULTS & DISCUSSIONS

A graphical user interface (GUI) has been developed to streamline the utilization of the system, enhancing user-friendliness and accessibility. Within this UI, users are provided with fundamental options enabling them to effortlessly select a video file, designate starting and ending times for video processing, and define a Region of Interest (ROI) for targeted analysis. Additionally, the interface offers the flexibility to specify input and output languages, empowering users to tailor the extraction process to their linguistic preferences.
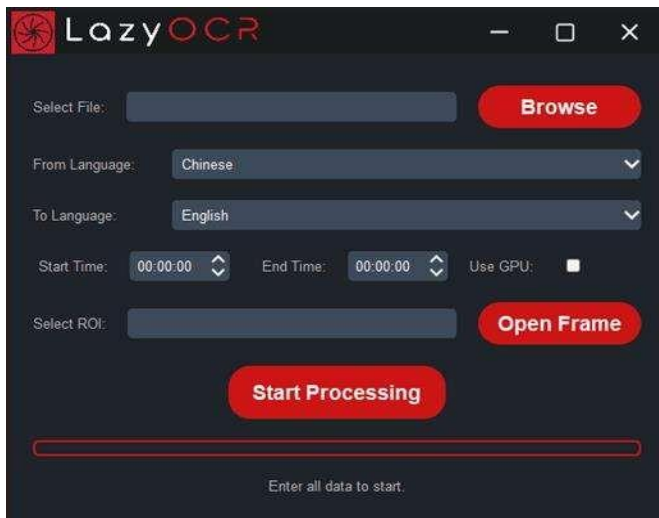


Fig. 2    User interface

Now it will start extracting the frames with subtitles from the selected video as images along with respective timestamps and move the extracted frames to a new file with starting and ending time stamps as the filename.
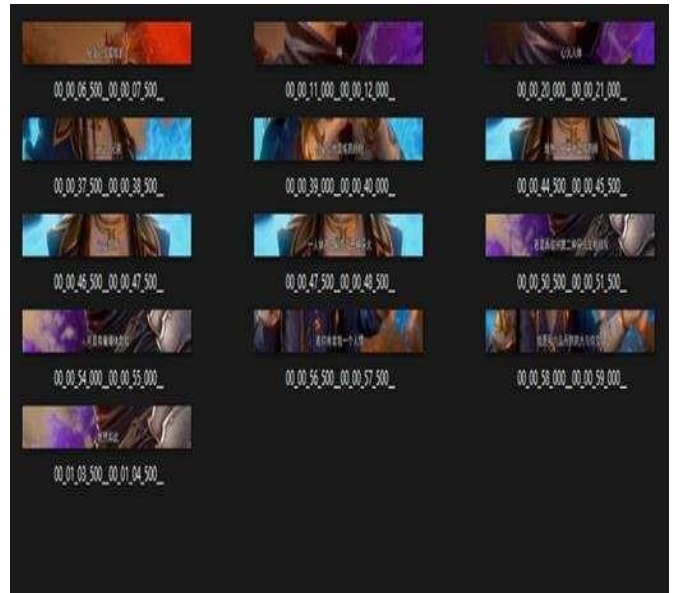


Fig. 3    Frame extraction with timestamps

The images shown in Fig 3 along with their timestamps can be used to extract the text. The extracted text is moved to a comprehensive SRT file. Then this SRT file is used for translating the subtitles into the language which the user have selected. The final output is a comprehensive SRT file containing the multilingual subtitles extracted from the video.
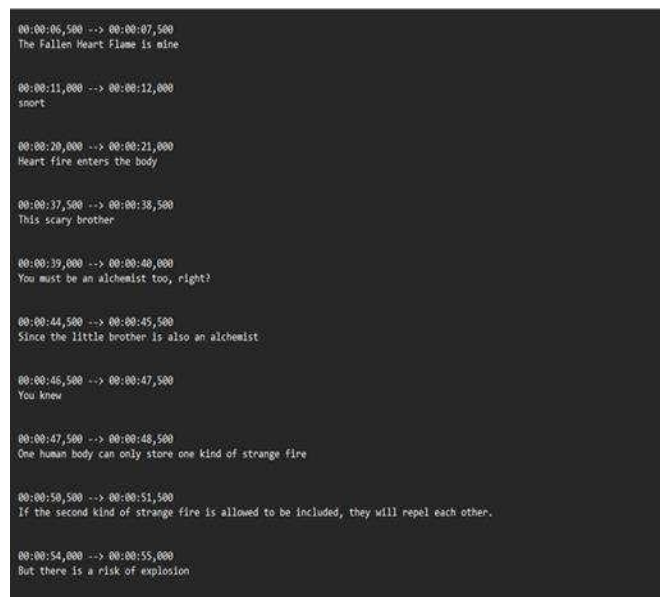


Fig. 4    SRT File with translated subtitles.

The SRT file Fig 4 contains all the subtitle translations. This file can be attached to the video files in any media player like VLC. This newly obtained SRT file can be added to the original video in any media player like VLC thus allowing the user to enjoy his favourite movie or show in his desired language. This newly file can be added to original video in a media player.

Fig. 5      Subtitles translated to English.



Fig. 6      Subtitle translated to Spanish

## VI.  EXPERIMENTAL RESULTS

Our multilingual hardcoded subtitle extractor stands out from existing systems by providing a seamless and user-friendly solution that requires minimal human intervention during text extraction. Unlike other systems, which often rely on manual intervention for accurate extraction and have complex user interfaces, our system simplifies the process, making it more user-friendly and accessible. Additionally, our system minimizes the risk of system crashes, especially on lower-end systems, ensuring a reliable and robust solution for users with different hardware configurations. Importantly, our approach combines text extraction and translation into a single service, eliminating the need for users to navigate between multiple platforms or services. This integration improves efficiency and reduces complexity, ultimately offering a more streamlined and efficient experience for extracting and translating subtitles from videos. By utilizing EasyOCR's advanced capabilities supporting over 80 languages, our system ensures comprehensive language coverage, enhancing accuracy and enabling seamless translation. Integration of OpenAI's GPT-3.5 Turbo further boosts translation quality and speed, delivering contextually relevant translations and significantly reducing processing time for faster results on large datasets. The combination of EasyOCR and GPT-3.5 Turbo in our user-friendly interface offers a solution for extracting and translating subtitles with better accuracy, efficiency, and language support.

EasyOCR distinguishes itself from other OCR solutions such as Paddle OCR and MMOCR by offering unmatched accuracy

and fast processing on GPU. It excels in extracting text from various languages and font styles, ensuring dependable results even with intricate visual content. Its efficient performance reduces GPU processing time, enhancing productivity and workflow efficiency. Moreover, EasyOCR's extensive language support covers more than 80 languages, making it the preferred option for efficient and precise text extraction in diverse applications and industries.

### COMPARISON

| | Paddle OCR | MM OCR | EasyOCR |
|---|---|---|---|
| Word(s) Accuracy | | | |
| Correct | 1141 | 270 | 1503 |
| False | 1185 | 2056 | 823 |
| | | | |
| Exact match % | 49.05 | 11.61 | 64.62 |
| Levenshtein Avg | 4 | 6.8 | 2.6 |
| Levenshtein % | 67.70 | 6.70 | 83.80 |
| Performance GPU(sec) | 2.40 | 1.14 | 0.79 |
| Performance CPU(sec) | 3.22 | 1.85 | 37.76 |

Fig. 7     Analysis

## VII.  CONCLUSION & FUTURE SCOPE

In conclusion, this strategic approach not only serves as a testament to the current strides made in multilingual support but also sets the stage for a future where language barriers within entertainment are effectively dismantled. As the system advances, prioritizing the expansion of language support, integration of advanced language models, and optimization for diverse devices emerges as critical milestones toward establishing a more inclusive and streamlined platform. By democratizing content access across linguistic divides, this methodology aligns seamlessly with the overarching objective of nurturing global connectivity and eradicating linguistic barriers within the entertainment sphere, thereby fostering a more interconnected and culturally diverse digital landscape.

Looking forward, several promising avenues for the project's evolution come to the fore. Primarily, the inclusion of additional languages and dialects promises to refine the system, catering to an even broader spectrum of users. Integration of state-of-the-art language models like ChatGPT for translation services holds immense potential in enhancing linguistic capabilities, ensuring more precise and contextually rich translations. Furthermore, optimizing the system to deliver superior performance on low-end devices paves the way for enhanced accessibility across a diverse range of platforms, furthering inclusivity and usability objectives.

## *REFERENCES*

[1]     Akinbade, D., Ogunde, A. O., Odim, M. O. & Oguntunde, B. O. (2020). An Adaptive Thresholding Algorithm-Based Optical Character Recognition System for Information Extraction in Complex Images. Journal of Computer Science, 16(6),784-801.

https://doi.org/10.3844/jcssp.2020.784.801

[2] J. Xu, W. Ding and H. Zhao, "Based on Improved Edge Detection Algorithm for English Text    Extraction and Restoration From Color Images," in IEEE Sensors Journal, vol. 20, no. 20, pp. 11951-11958,15Oct.15,2020 doi:10.1109/JSEN.2020.2964939

[3] Image processing application in character recognition Pradeep K. Nandaa,          Laxmi          Goswami. https://doi.org/10.1016/j.matpr.2021.03.697

[4] Bhansali, M. and P. Kumar, 2013. An alternative method for facilitating cheque clearance using smart phones application. Int. J. Applic. Innov. Eng. Manage., 2: 211-217.

[5] Bishop, C.M., 2006. Pattern Recognition and Machine Learning. 1st Edn., Springer, New York, ISBN-10: 0387310738.

[6] Chaudhuri, A., K. Mandaviya, P. Badelia and  S.K. Ghosh, 2017. Optical Character Recognition Systems for Different Languages with Soft Computing. 1st Edn., Springer International Publishing, ISBN-13: 9783319502519.

[7] Das, R.L., B.K. Prasad and G. Sanyal, 2012. HMM based offline handwritten writer independent English character recognition using global and local feature extraction. Int. J. Comput. Applic., 46: 45-50. DOI: 10.5120/6948-9428.

[8] Deb, K., A. Pratap, S. Agarwal and T.A.M.T. Meyarivan, 2002. A fast and elitist multiobjective genetic algorithm: NSGA-II. IEEE Trans. Evolut. Comput., 6: 182-197. DOI: 10.1109/4235.996017

[9] Ding, J., G. Zhao and F. Xu, 2018. Research on video text recognition technology based on OCR. Proceedings of the 10th International Conference on Measuring Technology and Mechatronics Automation, Feb. 10-11, IEEE Xplore Press, Changsha, China, pp: 457-462. DOI: 10.1109/ICMTMA.2018.00117