Original Article ISSN (Online): 2582-7472

### UNIFIED COMMUNICATION: A SURVEY ON HARMONIZING REGIONAL LANGUAGE DIVERSITY

Thenarasi V<sup>1</sup>, Santhosh Kumar B N<sup>2</sup>, Prakasha Raje Urs M<sup>3</sup>, Rashmi R<sup>4</sup>

- <sup>1</sup> Assistant Professor, Department of Computer Science, Government First Grade College, Siddartha Layout, Mysore, India
- <sup>2</sup> Assistant Professor of Computer Science, Maharani's Science College for Women, Mysore, India
- <sup>3</sup> Assistant Professor of Computer Science, Maharani's Science College for Women, Mysore, India
- <sup>4</sup> Assistant Professor of Physics, Maharani's Science College for Women, Mysore, India





#### DOI

10.29121/shodhkosh.v5.i1.2024.267

**Funding:** This research received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors.

**Copyright:** © 2024 The Author(s). This work is licensed under a Creative Commons Attribution 4.0 International License.

With the license CC-BY, authors retain the copyright, allowing anyone to download, reuse, re-print, modify, distribute, and/or copy their contribution. The work must be properly attributed to its author.



# **ABSTRACT**

This project tackles the complex task of extracting insights from text and images using Optical Character Recognition (OCR). After extracting text, language identification is crucial for a comprehensive multiclass classification approach, especially given the limitations of existing machine translation systems for Indian languages. The paper carefully examines challenges in machine translation, morphological analysis, parsing, word sense disambiguation, and the translation process to enhance the quality of translations. Beyond translation, the project includes automatic text summarization to distill essential content. Through the seamless integration of OCR, language detection, translation, and text summarization, our approach aims to facilitate unified communication by harmonizing diverse voices in multilingual settings.

**Keywords:** Text Extraction, OCR, Language Detection, Translation, Summarization

### 1. INTRODUCTION

Within the ever-evolving scene of advanced domains, where development and innovative progressions persistently shape our intelligent with data, there lies a endless supply of undiscovered potential inside the domain of picture information. Inside these pictures, lies a trove of idle treasures, holding up to be uncovered, each pixel carrying with it a story of semantic abundance that rises above unimportant visual representation. It is inside this complex embroidered artwork of pixels that the quintessence of information dwells, advertising significant bits of knowledge into the exceptionally texture of our advanced presence.

In an time where the storm of data inundates us from all corners, the capacity to extricate meaning from visual information stands as a linchpin in our journey for understanding. Among the bunch instruments at our transfer, robotized methods for extricating content from pictures rise as a guide of brightening, casting light upon the often-obscured pathways that lead to comprehension. This capability isn't just a comfort but a need, serving as a portal to a all encompassing understanding that's significant for a large number of applications, extending from report recovery to the nuanced examination of multilingual literary information.

Setting out on an unfamiliar voyage into the profundities of this advanced ocean, this paper sets cruise to examine the heap challenges related with dialect recognizable proof, machine interpretation, content summarization, and optical character acknowledgment (OCR) inside the setting of interpreting and summarizing territorial dialects. It may be a travel full with deterrents, where the turbulent waters of phonetic differences and semantic subtlety undermine to overwhelm the unwary traveler. However, equipped with the instruments of advanced innovation and the soul of request, we press forward, unfazed by the challenges that lie ahead.

At the heart of our endeavor lies a journey for greatness, a tireless interest of quality, exactness, and contextual relevance in our interpretations and outlines. To attain this elevated objective, we dive profound into the nuances of morphological structures and semantic varieties that recognize one dialect from another. It may be a travel of disclosure, where each etymological subtlety revealed brings us closer to our objective of bridging the crevice between different societies and dialects.

But our journey does not conclusion there. No, it expands past mere translation and summarization to envelop a broader vision of concordant communication among different voices in multilingual situations. It could be a vision of inclusivity, where each voice, in any case of dialect or origin, is given the opportunity to be listened and caught on. To realize this vision, we must not as it were depend on the brute drive of innovation but moreover on the artfulness of human understanding. It is a fragile adjust, where the cold rationale of machines must be tempered by the warmth of human sympathy.

Hence, our extend isn't just a specialized endeavor but a helpful one, pointed at cultivating a more comprehensive and associated world community. Through the consistent integration of OCR, dialect location, interpretation, and summarization, we endeavor to break down the boundaries that partitioned us and construct bridges of understanding that span the globe. It may be a respectable interest, one which will seem overwhelming in its scope, but one that's all things considered worth seeking after with all the passion and commitment that we are able muster.

Within the conclusion, our objective is simple:

to make a world where communication knows no boundaries, where dialect is now not a boundary but a bridge, interfacing hearts and minds over societies and landmasses. It could be a world where the abundance of human differing qualities is celebrated, not dreaded, and where the trade of thoughts streams openly, unencumbered by phonetic limitations. It could be a world that we may as it were glimpse faintly on the skyline, but one that's worth striving for with each fiber of our being.

# 2. LITERATURE SURVEY

- 1. "An Anatomization of Dialect Location and Interpretation utilizing NLP Procedures" by Bhagyashree P Pujeri and Jagadeesh Sai D. (2020): This paper gives a nitty gritty examination of dialect location and interpretation strategies leveraging Characteristic Dialect Handling (NLP) methods. It investigates different calculations and strategies utilized in these forms, advertising experiences into the fundamental instruments and challenges.
- 2. "Machine Interpretation between Malayalam and English" by Dr. Sreelekha S. (2020): Centering particularly on machine interpretation between Malayalam and English, this paper digs into the complexities of interpreting between these dialects. It addresses the phonetic complexities and social subtleties included within the interpretation prepare, beside headways and challenges in this space.
- 3. "On Extractive and Abstractive Neural Record Summarization with Transformer Dialect Models" by Jonathan Pilault, Raymond Li, Sandeep Subramanian, and Christopher Buddy (2020): This consider examines the utilize of Transformer dialect models for extractive and abstractive neural archive summarization. It investigates the adequacy of these models in creating brief rundowns from huge literary records, highlighting their potential applications and confinements.
- 4. "Making strides English-to-Indian Dialect Neural Machine Interpretation Frameworks" by Akshara Kandimalla, Pintu Lohar, Souvik Kumar Maji, and Andy Way (2022): Centered on upgrading English-to-Indian dialect neural machine interpretation frameworks, this paper proposes novel strategies and techniques to move forward interpretation precision and familiarity. It addresses the special etymological characteristics and challenges related with deciphering between English and Indian dialects.
- 5."Kannada to English Machine Interpretation Utilizing Profound Neural Organize"by Pushpalatha Kadavigere Nagaraj, Kshamitha Shobha Ravikumar, Mydugolam Sreenivas Kasyap, Medhini Hullumakki Srinivas Murthy, and Jithin Paul (2021): This paper presents a profound neural organize approach for interpreting content from Kannada to English. It explores the utilize of profound learning methods to overcome phonetic and auxiliary contrasts between the two dialects, pointing to move forward interpretation quality and precision.

- 6."An Proficient Long Short-Term Memory Demonstrate for Computerized Cross-Language Summarization" by Y. C. A. Padmanabha Reddy, Shyam Sunder Reddy Kasireddy, Nageswara Rao Sirisala, Ramu Kuchipudi, and Purnachand Kollapudi (2023):Centering on computerized cross-language summarization, this paper presents an productive long short-term memory (LSTM) show for summarizing multilingual literary information. It addresses the challenges of extricating important data from different dialect sources and proposes inventive arrangements to improve summarization exactness.
- 7. "Content Summarization Utilizing NLP" by ChetanaVaragantham, J.SrinijaReddy, UdayYelleni, MadhumithaKotha, and Dr P.VenkateswaraRao: This paper investigates content summarization methods utilizing Normal Dialect Preparing (NLP) strategies. It explores different approaches to naturally produce brief outlines from expansive literary reports, counting extractive and abstractive summarization methods. The consider assesses the viability of these strategies in condensing data whereas protecting the fundamental meaning and setting of the initial content.
- 8. "Programmed Kannada Text Extraction from Camera Captured Images" by Vipin Gupta, G.N. Rathna, and K.R. Ramakrishnan: Centered on Kannada content extraction from camera-captured pictures, this paper presents computerized methods to extricate printed substance from pictures captured by cameras or versatile gadgets. It talks about picture handling calculations and Optical Character Acknowledgment (OCR) strategies custom-made for extricating Kannada content, tending to challenges such as shifting textual styles, sizes, and foundations.
- 9. "Picture Division and Content Extraction:Application to the Extraction of Printed Data in Scene Pictures" by Danial Md Nor, Rosli Omar, M. Zarar M.Jenu, and Jean-Marc Ogier: This paper explores picture division and content extraction methods for extricating literary data from scene pictures. It investigates strategies to recognize and separate content locales inside complex picture foundations, empowering the extraction of significant literary substance for encourage handling or investigation.
- 10. "Script Recognizable proof and Dialect Discovery of 12 Indian Dialects utilizing DWT and Format Coordinating of As often as possible Happening Character(s)" by Jeelen Kumar Sarungbam, Bhupendra Kumar, and Ankur Choudhary: Centering on script distinguishing proof and language detection of 12 Indian dialects, this paper proposes strategies based on Discrete Wavelet Transform (DWT) and template matching strategies. It addresses the challenges of recognizing dialects and scripts in multilingual situations, pointing to create strong arrangements for dialect distinguishing proof errands.
- 11. "Telugu to English Interpretation utilizing Coordinate Machine Interpretation Approach" by T. Venkateswara Prasad and G. Mayil Muthukumaran: This paper presents a coordinate machine interpretation approach for interpreting content from Telugu to English. It examines the application of rule-based and factual machine interpretation procedures to encourage coordinate interpretation between the two dialects, tending to phonetic and basic contrasts to make strides interpretation precision and familiarity.
- 12."Tamil OCR Change from Computerized Composing Cushion Acknowledgment Exactness Makes strides through Adjusted Profound Learning Designs" by V. Jayanthi and S. Thenmalar: Centered on Tamil Optical Character Acknowledgment (OCR) transformation from advanced composing cushion acknowledgment, this paper investigates adjusted profound learning structures to improve OCR exactness. It examines strategies to make strides the acknowledgment of Tamil characters from advanced composing cushions, pointing to create more strong OCR frameworks for Tamil script recognition.

### 3. METHODOLOGY

Envisioned against the backdrop of this linguistic tapestry, our project unfurls with four intertwined goals, each a crucial chapter in the narrative of empowering multilingual communication:

1. Unveiling Linguistic Narratives from Visual Canvases:

The first brushstroke of our endeavour involves liberating textual gems embedded within images. Text extraction from textual images emerges as a pivotal process, setting the stage for subsequent linguistic exploration. The characters confined to visual spaces become the raw material for unravelling linguistic nuances.

2. Pioneering Linguistic Recognition:

In the symphony of languages, the extracted text takes centre stage, playing a dual role. Beyond being mere characters on a digital canvas, this text becomes the linchpin for language recognition. Unveiling the essence of languages is paramount in applications like machine translation, information retrieval, and summarization, where the first step is understanding the language before unravelling its content.

3.Bridging Linguistic Divides through Translation:

Acknowledging the global thirst for multilingual communication, our project extends its arms across linguistic boundaries. The extracted text, diverse in its regional origins, undergoes a transformative journey, crossing linguistic frontiers to emerge unified in a common language—English. This translation metamorphosis serves as a bridge, ensuring accessibility and comprehension across the diverse linguistic kaleidoscope.

## 4.Distillation of Information through Summarization:

Navigating the labyrinth of unstructured data, our project embraces the challenge of text summarization. Here lies the quest to distil the essence, to craft concise summaries that not only expedite research but also reduce reading time. Text summarization, a frontier in machine learning and natural language processing, becomes a tool for enhancing overall information accessibility. Incorporating diverse voices from newspapers, medical records, legal documents, and reports, our Automated Text Summarization (ATS) system embarks on a dual journey—exploring extraction-based summarization for key points and abstraction-based summarization for refined narratives. This endeavour seeks not only to harmonize voices but to create a symphony, resonating in the intricate dance of language recognition, translation, and text summarization. In doing so, we strive to foster unified communication in the kaleidoscopic world of multilingual textual information within images.

### 4. IMPLEMENTATION

### 1. GUI Design using QT Designer:

The GUI layout is designed using QT Designer, a drag-and-drop tool for creating Qt-based interfaces. The .ui file generated from QT Designer defines the layout and components of the interface.

# 2. Integration with PyQt5:

The QT Designer-generated UI file is integrated with the PyQt5 library, a set of Python bindings for the Qt application framework. This integration allows the application logic to interact with the UI components defined in the .ui file.

# 3. Text Extraction from Images:

The application utilizes the pytesseract library for OCR, extracting text from images. OpenCV (cv2) is used for image processing tasks such as grayscale conversion and Gaussian blurring to enhance OCR accuracy.

#### 4. Translation of Extracted Text:

The Googletrans library is employed for translating the extracted text into the desired language. Users can select the source and target languages using dropdown menus in the interface.

### 5. Summarization of Translated Text:

The Sumy library is used for summarizing the translated text into shorter, more concise versions. Summarization is performed based on word frequencies and sentence scores calculated from the translated text.

#### 6. Error Handling:

Robust error handling mechanisms are implemented to handle exceptions and errors encountered during text extraction, translation, and summarization processes. This ensures the reliability and stability of the application.

### 7. Event Handling and User Interaction:

PyQt5's event handling mechanisms are utilized to capture user interactions with the interface components. For example, mouse events are used to define regions of interest for text extraction from images.

#### 8. Optimization and Refinement:

Continuous optimization and refinement are performed to improve the accuracy, performance, and efficiency of text extraction, translation, and summarization tasks. This may involve adjusting parameters, fine-tuning algorithms, and enhancing error handling mechanisms.

### 9. Testing and Validation:

The implemented features are tested and validated to ensure they meet the requirements and specifications of the application. This involves both unit testing of individual components and integration testing of the entire system.

# 10. Documentation and Deployment:

Documentation is prepared to provide guidelines for users and developers on how to use and extend the application. Once validated, the application can be deployed for use by end-users, either as a standalone executable or as a packaged Python application.

Overall, the implementation process involves designing the GUI, integrating with external libraries for OCR, translation, and summarization, implementing error handling mechanisms, optimizing performance, testing, and finally deploying the application for use.

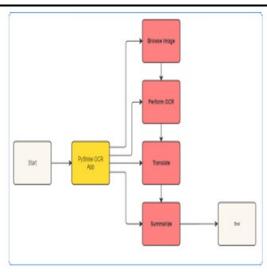


FIGURE 4: USE CASE DIAGRAM

#### LIBRARIES USED:

#### PyQt5:

PyQt5 is used for creating the graphical user interface (GUI) of the application. In the project, PyQt5 is utilized to design and manage the main window, buttons, labels, combo boxes, and other UI components. Signals and slots mechanism is also used to connect GUI events to corresponding functions.

# pytesseract:

pytesseract is employed to perform optical character recognition (OCR) on images. It extracts text from images by leveraging the Tesseract-OCR Engine developed by Google.In the project, pytesseract is used to convert selected image regions (cropped portions) into text. This text extraction functionality is triggered when the user selects a region of an image using the GUI.

### cv2 (OpenCV):

OpenCV (cv2) is utilized for various image processing tasks, including reading images, converting color spaces, applying filters, and cropping images. In the project, cv2 is used to read images from files, convert image color spaces (e.g., from BGR to RGB), and crop selected regions of images based on user input.

#### os:

The os module provides functions for interacting with the operating system, such as accessing files and directories. In the project, the os module is used to manipulate file paths and access files required for language data.

#### SVS:

The sys module provides access to system-specific parameters and functions related to the Python interpreter. In the project, sys is used to manage command-line arguments and control the execution of the application.

### PIL (Python Imaging Library):

PIL is utilized for handling images and performing various image processing operations. In the project, PIL is used to create image objects from numpy arrays (as returned by OpenCV), which can then be displayed in the PyQt5 GUI.

# glob:

The glob module is used for searching file pathnames matching a specified pattern. In the project, glob is used to find language data files (\*.traineddata) within a specified directory.

### googletrans:

googletrans is a Python wrapper around the Google Translate API, allowing translation between different languages.In the project, googletrans is used to translate text from one language to another based on user input.

#### textblob:

textblob is a library for natural language processing (NLP), providing tools for common text processing tasks such as tokenization, part-of-speech tagging, and sentiment analysis. In the project, textblob is not explicitly used in the provided code snippet. It might be intended for additional text processing or analysis tasks.

### sumy:

sumy is a library for automatic text summarization, which summarizes longer texts into shorter versions while retaining the main points. In the project, sumy is used to generate summaries of translated text. It extracts key sentences from translated text to provide a concise summary.

### re (Regular Expressions):

The re module provides support for regular expressions, which are patterns used for searching and manipulating text. In the project, re is used to perform pattern matching and text processing tasks, such as splitting text into sentences for summarization.

#### **OUTPUT:**

# **User Interface (UI) Design:**

The project's user interface (UI) is a testament to thoughtful design and intuitive interaction. Developed using Qt Designer, a robust tool for creating graphical interfaces in Qt applications, the UI provides users with a seamless experience while interacting with the translation model. It encompasses a variety of elements including buttons, labels, text fields, and dropdown menus, strategically arranged to optimize user experience. Through Qt Designer's intuitive interface, the UI achieves clarity and efficiency, guiding users through the translation process with ease. By prioritizing usability and accessibility, the UI ensures that users can navigate effortlessly and access the translation model's functionalities without complexity.

#### **Core Functionalities:**

At the heart of the paper lie with essential functionalities that enable accurate and efficient text processing and translation tasks. These core functionalities are meticulously implemented to ensure reliability and effectiveness in handling various aspects of text extraction, translation, and summarization. The project leverages pytesseract for its OCR tool implementation, enabling precise text extraction from images. OpenCV (cv2) is utilized for image preprocessing tasks such as grayscale conversion and Gaussian blurring, enhancing the quality of images prior to text extraction. Translation capabilities are facilitated through integration with the Googletrans library, allowing for seamless translation of extracted text into English or other specified languages. Additionally, text summarization functionalities are achieved using Sumy, enabling the generation of concise summaries from translated text. These core functionalities collectively form the backbone of the translation model, ensuring its efficacy in meeting user needs and expectations.

#### **Accomplishments:**

The paper accomplishments extend beyond the mere implementation of functionalities, reflecting a commitment to excellence and innovation. Noteworthy achievements include the development of a functional interface that enhances user interaction and accessibility. Through meticulous error handling mechanisms, the project ensures robustness and reliability, effectively managing exceptions and errors that may arise during text processing tasks. Moreover, the project demonstrates cross-platform compatibility, enabling deployment across various operating systems without compromising functionality or performance. Optimized performance further enhances efficiency in text processing tasks, ensuring timely and accurate translation and summarization. Integration with external libraries and tools enriches the project's functionality and features, enhancing the overall user experience. These accomplishments underscore the project's success in delivering a robust machine translation model that meets the needs of users seeking seamless text processing capabilities

### 5. CONCLUSION

The endeavor to harmonize diverse voices through the translation and summarization of regional languages for unified communication is underscored by formidable challenges. The scarcity of resources, absence of annotated corpora, and script variations in Indian languages pose significant hurdles for effective language processing. Cultural nuances further complicates the translation efforts, demanding a nuanced understanding beyond literal meanings. The dearth of comprehensive tools for regional languages amplifies these challenges, especially in addressing idiomatic expressions and complex sentence structures. Overcoming linguistic diversity in India requires high translation accuracy, coherent summarization, and efficient handling of various scripts. Beyond textual challenges, the extraction of text from images introduces complexities in accuracy, script diversity, multilingual recognition, contextual understanding, and resource-intensive image annotation. Tackling these challenges necessitates not only robust machine learning models but also comprehensive tools and frameworks that can adapt to the linguistic intricacies of India. The integration of image text

extraction, translation and summarization in a unified system must prioritize efficiency, scalability, and adaptability to foster meaningful communication across diverse linguistic landscapes.

# **CONFLICT OF INTERESTS**

None

# **ACKNOWLEDGEMENTS**

None

### REFERENCES

- Bhagyashree P Pujeri, Jagadeesh Sai D. (2020). "An Anatomization of Language Detection and Translation using NLP Techniques" .International Journal of Innovative Technology and Exploring Engineering (IJITEE). Volume-10 Issue-2,pp.69-77
- Dr. Sreelekha S. (2020). "Machine Translation between Malayalam and English". Linguistics Journal. Volume 14 Issue 2,pp. 7-30
- Jonathan Pilault, Raymond Li, Sandeep Subramanian and Christopher Pal. (2020). "On Extractive and Abstractive Neural Document Summarization with Transformer Language Models". Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing. pp. 9308–9319
- Akshara Kandimalla 1, Pintu Lohar 2, Souvik Kumar Maji and Andy Way. (2022). "Improving English-to-Indian Language Neural Machine Translation Systems. Information 13,245.
- Sangeetha G, Prathusha Laxmi B, and Vijayaraja V "A Survey On Web-Based Intelligent Chat Bot", MDPI, 2018
- Pushpalatha Kadavigere Nagaraj , Kshamitha Shobha Ravikumar, Mydugolam Sreenivas Kasyap, Medhini Hullumakki Srinivas Murthy, Jithin Paul. "Kannada to English Machine Translation Using Deep Neural Network". Ingénierie des Systèmes d'Information Vol. 26. No. 1, February, 2021, pp. 123-127.
- Y. C. A. Padmanabha Reddy, Shyam Sunder Reddy Kasireddy, Nageswara Rao Sirisala, Ramu Kuchipudi and Purnachand Kollapudi." An Efficient Long Short-Term Memory Model for Digital Cross-Language Summarization". CMC, 2023,vol.74, no.3.
- ChetanaVaragantham, J.SrinijaReddy, UdayYelleni, MadhumithaKotha, Dr P.VenkateswaraRao." Text summarization using nlp". Journal of Emerging technologies and Innovative Research(JETIR).
- Vipin Gupta, G.N.Rathna, K.R.Ramakrishnan." Automatic kannada text extraction from camera captured images".
- Danial Md Nor, Rosli Omar, M. Zarar M.Jenu, and Jean-Marc Ogier." Image segmentation and text extraction: Application to the extraction of textual information in scene images". International Seminar on Application of Science Mathematics 2011.
- Jeelen Kumar Sarungbam, Bhupendra Kumar, Ankur Choudhary. "Script Identification and Language Detection of 12 Indian Languages using DWT and Template Matching of Frequently Occurring Character(s)". International Conference- Confluence The Next Generation Information Technology Summit.Walaa Hassan, Shereen elBohy, Min Rafik, Ahmed Ashraf, Sherif Gorgui, Michael Emil, Karim Ali "An Interactive Chatbot for College Enquiry", 2023
- T. Venkateswara Prasad, G. Mayil Muthukumaran." Telugu to English Translation using Direct Machine Translation Approach". International Journal of Science and Engineering Investigations. vol. 2, issue 12, January 2013.pp. 25-32.
- Suyash Awasthi, Anupriya Purwar, Dhananjay Batra, Prof. Prakash Devale, "Web Based College Chatbot SDABot", 2021. V. Jayanthi and S. Thenmalar. "Tamil OCR Conversion from Digital Writing Pad Recognition Accuracy Improves through Modified Deep Learning Architectures". Hindawi Journal of Sensors.