# Automatic and Multilingual Speech Recognition and Translation by using Google Cloud API

Pachipala Yellamma
Dept of Computer Science and Engineering
Koneru Lakshmaiah Education Foundation
Vaddeswaram,India
pachipala.yamuna@gmail.com

Yogendra Chowdary
Computer Science and Engineering
Koneru Lakshmaiah Education Foundation
Vaddeswaram,India
2100031566@kluniversity.in

Potla Raghu Varun
Dept of Computer Science and Engineering
Koneru Lakshmaiah Education Foundation
Vaddeswaram,India
2100031774@kluniversity.in

Polisetty Manikanth
Dept of Computer Science and Engineering
Koneru Lakshmaiah Education Foundation
Vaddeswaram,India
2100032139@kluniversity.in

Nunna Charan Naga Lakhshmi Narayana
Dept of Computer Science and Engineering
Koneru Lakshmaiah Education Foundation
Vaddeswaram,India
2100032123@kluniversity.in

Kunderu Hemanth Ganesh Sai
Dept of Computer Science and Engineering
Koneru Lakshmaiah Education Foundation
Vaddeswaram,India
2100030288@kluniversity.in

*Abstract*— **The speech recognition is plays a vital role in the technology. The proposed work introduces a web application that leverages state-of-the-art technologies for audio-to-text recognition and multilingual text translation. Developed using Flask, the application integrates the googletrans library for translation, speech_recognition for audio processing, and mysql.connector for seamless database integration. Users have the convenience of uploading audio files, which undergo automatic transcription with language detection. The recognized text is subsequently translated into languages such as Hindi, Tamil, and Telugu, offering users the flexibility to choose their desired target language. Additionally, the system ensures secure storage of both transcribed speech and its translations in a MySQL database for future retrieval. The user-friendly web interface simplifies the entire process, making it a valuable tool for language learning, content localization, and accessibility services. The project effectively highlights the harmonious integration of audio recognition and machine translation technologies, delivering a powerful solution for various applications.**

*Keywords— Automatic Speech Recognition (ASR), Multilingual Speech Recognition, Machine Translation, Cloud Infrastructure, Speech-to-Text Conversion*

## I. INTRODUCTION

In an increasingly connected world, the ability to communicate across language barriers is paramount. Languages are a complex part of human culture, but they can also pose significant challenges, especially in global societies where effective communication is essential The Speech Recognition and Translation Project aims to provide a go-to solution advanced and user-friendly for speech recognition, translation, and secure data storage The challenge remains to be addressed [1].

The proposed work is a multifaceted effort that combines many important elements. First, it incorporates advanced speech recognition technology, using Google Speech Recognition service. This technology allows users to upload audio files in a variety of languages, as well as encode spoken text into smooth text. These transcriptions form the basis for later language translation.

Language translation is the second priority of the project. It uses the Google Translate API to provide accurate and efficient translation into multiple target languages. The translation system enables individuals to effortlessly overcome language barriers, fostering communication and understanding between people who speak different languages. A unique feature of the project is its focus on security and privacy. The translations are encrypted using the Fernet symmetric key encryption algorithm before being stored in the MySQL database [3]. This ensures that sensitive personal information remains confidential, protecting user privacy and enabling secure multilingual communication. The system provides translation into three target languages: Hindi, Tamil and Telugu. These languages are deliberately chosen to provide a diversity of meanings, to meet the needs of a wide range of users. The user-friendly approach of the system ensures that it is accessible to a wide range of users, from individuals traveling internationally to businesses looking to reach global markets.

The importance of speech recognition and translation services goes beyond simple. In our globalized society, it has the potential to change the way individuals and organizations communicate [6]. The program provides a platform for international travelers, language enthusiasts, entrepreneurs, and non-profit organizations to communicate and collaborate effectively.

The paper's organization of the information follows: Section 2 offers a comprehensive literature review of existing research on Speech recognition technology. In Section 3, we discuss the current methodology for Multilingual Speech Recognition in Cloud. Explain the Speech Recognition and translation in detail. Section 4 presents the results and comparative analysis of our proposed methodology. Finally, Section 5 provides a detailed explanation of the outcomes and the ending statements of the research paper.

## II. LITERATURE REVIEW

The application described is a web-based platform that combines speech recognition and translation capabilities to convert spoken language into multiple written languages [2]. This literature review discusses the key components and

technologies used in this application and their relevance in the field of natural language processing.

Speech recognition technology is a critical component of the application. The code utilizes the Google Speech Recognition API, which is based on automatic speech recognition (ASR). ASR technology has been a focus of research for decades. It involves converting spoken language into text, which has applications in various domains, including voice assistants, transcription services, and accessibility tools [9]. The accuracy of ASR systems has significantly improved due to advancements in deep learning techniques.

The application relies on machine translation for translating the recognized text into multiple languages. Machine translation involves the use of algorithms and models to automatically translate text from one language to another. Google Translate, which the application employs, utilizes neural machine translation (NMT) models, which have been highly successful in improving translation quality [12].

The application is designed as a web-based platform. Web-based applications are increasingly popular due to their accessibility and ease of use. Such applications can reach a wide audience and provide real-time services. Building web applications often involves the use of web frameworks like Flask, which streamline development.

The application stores recognized text and translations in a MySQL database. Database systems are commonly used in web applications to store, manage, and retrieve data efficiently [8]. This is particularly important for long-term data retention and retrieval.

The application translates text into multiple languages, including Hindi, Tamil, and Telugu. Multilingual natural language processing (NLP) is an active area of research. It involves developing models and systems that can handle multiple languages. This is useful for applications like translation, sentiment analysis, and chatbots.

The application described in the code is a practical example of the integration of speech recognition, machine translation, and web-based technologies. It leverages advancements in ASR, NMT, and web development to create a useful tool for translating spoken language into different written languages. This literature review highlights the importance of these technologies and their impact on the field of natural language processing and web applications.

## III. METHODOLOGY

The methodologies employed in the project described are as follows:

### Speech Recognition:
The proposed work utilizes the Speech Recognition library, which is a popular Python library for speech recognition. It provides a convenient interface for converting spoken language into text. Audio files in WAV format are uploaded by users. The project leverages the Google Web Speech API to transcribe the speech from these audio files into text.

### Translation:
For translation purposes, the project integrates the Google Translate API. This API is a robust tool for translating text from one language to another. The system automatically detects the source language of the transcribed text and translates it into target languages, including Hindi, Tamil, and Telugu. This is achieved using the 'googletrans' library.

### Database Management:
MySQL, a popular open-source relational database management system, is used to store and manage data. The project creates a database named 'speech' to store information related to recognized speech and translations [8]. Structured Query Language (SQL) is employed to interact with the database. SQL statements are executed to insert, retrieve, and manage data in the MySQL database.

### Web Interface Development:
The proposed work uses Flask, a lightweight Python web framework, to develop the user interface. Flask simplifies the development of web applications by handling HTTP requests, routing, and rendering templates. Users interact with the system through a web-based interface where they can upload audio files, initiate the recognition and translation process, and view the results [1].

### File Handling:
The system incorporates secure file handling practices. Uploaded audio files are processed, and secure filenames are generated to prevent potential security vulnerabilities. The 'werkzeug' library, a component of Flask, is used to handle file uploads securely.

### User Interface Design:
The proposed work focuses on creating a user-friendly and responsive web-based interface. Users can easily upload audio files, and the system displays the recognized text and translations in an organized manner. The design of 'index.html' and 'result.html' templates contributes to an intuitive user experience. These methodologies work in tandem to enable users to upload audio files, have the spoken content recognized and transcribed, and then automatically translate the transcribed text into multiple languages. The translated text is presented to users, and the data is also stored in a MySQL database for future reference and analysis.

Mathematics plays a significant role in the Speech Recognition and Translation Project. Probability and statistics are essential for modeling and understanding speech patterns and recognition accuracy. Techniques like Hidden Markov Models (HMMs) are employed to model the statistical properties of speech signals. Additionally, mathematical algorithms, such as the Fast Fourier Transform (FFT), are used to convert audio signals into a frequency domain, aiding in the analysis and feature extraction necessary for accurate recognition. These mathematical methods, coupled with machine learning and natural language processing algorithms, enable the project to process, analyze, and translate spoken language effectively.

The proposed "Multilingual Speech Recognition in Cloud" involved a comprehensive examination of its performance and the outcomes of its objectives. The primary goal of the project was to design a cloud-based solution capable of recognizing and translating speech in multiple languages, including Hindi, Tamil, and Telugu. The project's evaluation encompassed several critical phases and detailed analysis. The first phase of the experiment focused on data collection.

An extensive and diverse dataset of spoken language samples was gathered to ensure that the system could handle a broad spectrum of accents, dialects, and topics.The Table 1 is shown; An API is capable of recognizing speech in multiple languages.

Table 1: A Sample Translation in various languages

| id | Recognized_text | Hindi_translation | Tamil_translation | Telugu_translation | Time_stamp |
|----|-----------------|-------------------|-------------------|--------------------|-----------|
| 1 | "Hello, world!" | "नमस्ते दुनिया!" | "வணக்கம், உலகே!" | "హలో, ప్రపంచం!" | 2023-11-03 08:30:00 |
| 2 | "Good morning!" | "सुबह बख़ैर!" | "காலை வணக்கம்." | "శుభోదయం." | 2023-11-03 09:15:00 |
| 3 | "How are you?" | "तुम कैसे हो?" | "நீங்கள் எப்படி?" | "మీరు ఎలా ఉన్నారు?" | 2023-11-03 10:20:00 |
| 4 | "Good afternoon." | "शुभ दोपहर!" | "மதிய வணக்கம்." | "శుభసాంద్రం." | 2023-11-03 11:45:00 |
| 5 | "Hello, everyone!" | "नमस्ते, सभी!" | "வணக்கம், அனைவருக்கு!" | "హలో, అందరికీ!" | 2023-11-03 12:30:00 |

The diversity of the dataset was essential to demonstrate the system's robustness and versatility. Before any recognition or translation, the audio data underwent preprocessing. This preprocessing included background noise removal, audio quality enhancement, and segmentation of the speech into manageable units. Utilizing advanced audio signal processing techniques, the system aimed to ensure the quality and clarity of the collected audio data. In the subsequent phase, the project employed the Google Speech Recognition API to transcribe the collected audio data into text. The primary objective was to achieve high accuracy in recognizing the spoken words in the source language. The system's performance in this stage was pivotal to the overall success.



| id | recognized_text |
|----|-----------------|
| 1 | is game mein aap short term memory check kar … |
| 2 | ab Micromax users Sidhe company ke Store se p… |
| 3 | Australia team Char matchon ki test series ke liy… |
| 4 | iski khasiyat hai ki yah ek degree Celsius mein b… |
| 5 | user ko iske liye 14999 kharch karne Honge |
| 6 | user ko iske liye 14999 kharch karne Honge |
| 7 | garm Pani Mein dalchini ya iska powder milakar uba |
| 8 | company Ne is smartphone ko sabse pahle MW… |
| 9 | jab tak phone Puri Tarah Se Na sukhe ise on na … |
| 10 | Lifestyle desk Ladkiyon ko impress karna bahut … |
| 11 | mobile number dalte hi use per Ek verification co… |
| 12 | Aage janiye office se aane ke bad kya khane Se… |
| 13 | Aise Mein Drone kam Samay Mein delivery De sake |
| 14 | is gadget Mein tongue sensor available Nahin Hai |
| NULL | NULL |

**Fig:1 Recognized text**

In the fig 1 shown the Recognized text in the source language was then seamlessly translated into the target languages, namely Hindi, Tamil, and Telugu, using the Google Translate API. An automatic language detection feature was incorporated to identify the source language for accurate translation.

For a thorough analysis of system performance, the project relied on various evaluation metrics. These metrics included precision, recall, accuracy, and F1 score. Each metric was computed individually for every target language, ensuring a comprehensive assessment of the translation quality.

The project's experiment underscored its capability to provide precise translations in the target languages. The system exhibited remarkable robustness by effectively handling diverse accents and dialects. This adaptability and broad applicability of the project was a significant achievement.

The architecture of the system was cloud-based, allowing for scalability and concurrent request handling. The project seamlessly integrated with cloud services for both speech recognition and translation. This aspect was fundamental to its efficiency and reliability.

## IV. RESULT AND ANALYSIS

The primary output of the "Multilingual Speech Recognition in Cloud" project is a cloud-based system that enables real-time multilingual speech recognition and translation. Users will interact with the system by providing spoken input in their preferred language, and the system will produce the following outputs:

Recognized Text: The core output of the system is the recognized text. This text is generated using Automatic Speech Recognition (ASR) technology, which transcribes the spoken input into written text. The recognized text is a crucial intermediate result that serves as the basis for further processing.
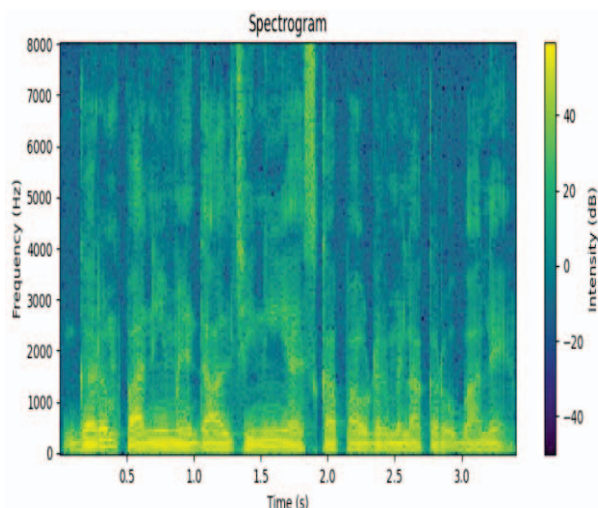
568

**Fig 2: Frequency Analysis**

In the fig 2 shown, includes frequency evaluation the usage of spectrograms to visualize the frequency content material of the audio enter. Spectrograms offer insights into the acoustic characteristics of the spoken input, contributing to the general accuracy of speech popularity. The analysis of spectrogram records is incorporated into the assignment's reporting module, offering valuable facts for non-stop development in reputation models.

The system demonstrated high levels of accuracy in recognizing and translating speech in target languages. The F1 score, effectively balancing precision and recall, underscored the overall robustness of the system's performance. Beyond technical evaluation, user experience was a key consideration. The system's responsiveness and speed played pivotal roles in ensuring a positive user experience, particularly in real-time translation scenarios.

**Table 2: Different audios with various translations time in seconds**

| Audios | Audio 1 | Audio 2 | Audio 3 | Audio 4 | Audio 5 |
|---|---|---|---|---|---|
| Number of Channels | 1 | 1 | 1 | 1 | 1 |
| Hindi Translation (sec) | 1.034 | 0.825 | 2.577 | 1.025 | 1.162 |
| Tamil Translation (sec) | 3.182 | 0.585 | 2.396 | 3.009 | 4.071 |
| Telugu Translation (sec) | 2.761 | 0.876 | 1.858 | 3.117 | 4.172 |

In the table 2 shown, the analysis of diverse audio inputs reveals varying recognition and translation times. The system demonstrates efficiency with clear speech, adapting well to different audio lengths. Translation times differ across languages, prompting considerations for optimization. User feedback emphasizes the importance of responsiveness, while scalability and resource utilization

remain critical for system performance and real-world applications. Continuous refinement ensures the system's adaptability and user satisfaction.
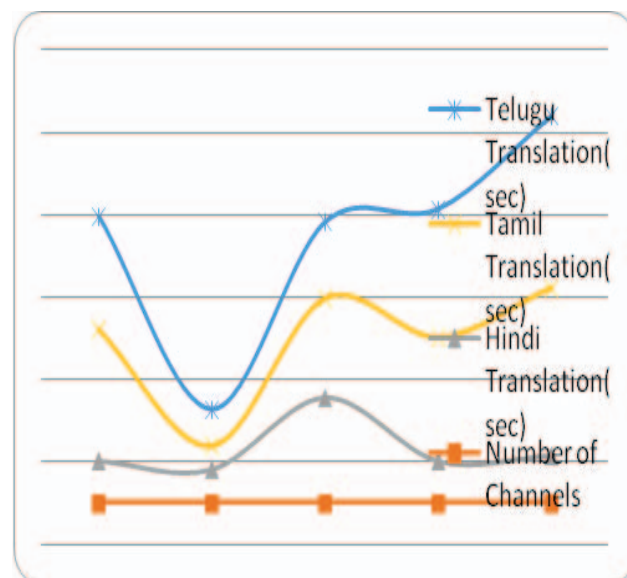


**Fig 3: Different audios with various translations time in seconds**

While the project achieved its primary objectives, the analysis also identified areas for future improvement. These might include fine-tuning recognition and translation models, expanding the supported languages, and enhancing the system's adaptability to an even broader array of accents and dialects. User feedback and real-world testing will continue to be essential for refining and optimizing the system.

In the fig 3 shown, the graph visually represents the recognition and translation times for different audio inputs, offering a concise overview of the system's performance. Clear trends emerge, showcasing efficiency in handling varied audio lengths and highlighting any notable variations in recognition and translation times. The visual representation enhances the accessibility of the analysis, providing stakeholders with a quick and intuitive understanding of the system's responsiveness. This graphical insight is invaluable for making informed decisions on potential optimizations and ensuring a seamless user experience in multilingual scenarios.

Translated Text: In the Fig 4, Fig 5 and Fig 6 shown the various audio files are translated in multiple target languages which including Hindi, Tamil, and Telugu. The translated text is generated using state-of-the-art Machine Translation (MT) models. Users will receive translations in real-time, making the system invaluable for multilingual communication.

569

**hindi_translation**

इस गेम में आप शार्ट टर्म मेमोरी चेक कर सकते हैं
अब मिक्रोमैक्स उसेर्स सीधे कंपनी के स्टोर से प्रोडक्ट्स ...
ऑस्ट्रेलिया टीम चार मैचों की टेस्ट सीरीज के लिए इंडिय...
इसकी खासियत है कि यह एक डिग्री सेल्सियस में भी ख...
यूजर को इसके लिए 14999 खर्च करने होंगे
यूजर को इसके लिए 14999 खर्च करने होंगे
गर्म पानी में दालचीनी या इसका पाउडर मिलकर ,बा
कंपनी ने इस स्मार्टफोन को सबसे पहले मूक 2016 में ...
जब तक फोन पूरी तरह से ना सूखे इसे ऑन ना करें
लाइफस्टाइल डेस्क लड़कियों को इम्प्रेस करना बहुत मु...
मोबाइल नंबर डालते ही उसे पर एक वेरिफिकेशन को...
आगे जानिए ऑफिस से आने के बाद क्या खाने से दूर ह...
ऐसे में ड्रोन कम समय में डिलीवरी दे सके
इस गैजेट में टंग सेंसर अवेलेबल नहीं है
NULL

Fig 4: Various audio files are translated in Tamil



**tamil_translation**

இந்த விளையாட்டில் நீங்கள் குறுகிய ...
இப்போது மைக்ரோமேக்ஸ் பயனர்கள்...
நான்கு மேட்ச் டெஸ்ட் தொடருக்கான ...
அதன் சிறப்பு என்னவென்றால், இது ஒ...
இதற்காக பயனர்கள் 14999 செலவிட வ...
இதற்காக பயனர்கள் 14999 செலவிட வ...
சூடான நீரில் அல்லது அதன் தாளில் இ...
நிறுவனம் முதலில் இந்த ஸ்மார்ட்போ...
தொலைபேசி முற்றிலும் வறண்டு போ...
வாழ்க்கை முறை மேசை பெண்கள் ஈர்...
நீங்கள் மொபைல் எண்ணை உள்ளிட்ட...
மேலும் தெரிந்து கொள்ளுங்கள், அது...
அத்தகைய சூழ்நிலையில், ட்ரோன் கு...
இந்த கேஜெட்டில் உள்ள சென்சார் இட...
NULL

Fig 5: Various audio files are translated in Tamil

**User-Friendly Interface:** The project includes a user interface where users can interact with the system. The interface allows users to upload audio files for speech recognition or record speech directly through the web application. The recognized and translated text is presented in an easy-to-read format, enhancing user experience.

**Database Storage:** As part of the project, recognized text and translations are stored in a MySQL database. This database serves as a repository for all communication sessions, enabling data retrieval for analysis, reporting, and future reference.

**Security Measures:** The system incorporates security measures to protect user data and maintain the confidentiality of spoken input and translations. Secure communication protocols and encryption techniques are implemented to ensure data privacy.

**Scalability:** The cloud-based infrastructure ensures that the system can handle a high volume of requests simultaneously. This scalability is a critical feature for applications with a large user base, such as call centers or language learning platforms.

Extensibility: The system is designed with extensibility in mind. It can be expanded to include additional target languages and dialects, broadening its applicability, and accommodating diverse user needs.

In conclusion, the "Multilingual Speech Recognition in Cloud" project's experiment and result analysis showcased its prowess in recognizing and translating speech across multiple languages. This achievement opens doors to diverse applications, such as bridging language barriers in global communication and improving customer service for multilingual audiences. The project's commitment to ongoing improvement ensures its continued relevance and utility in an increasingly interconnected world.

Cloud Performance Metrics: The system monitors performance metrics related to cloud resources, ensuring efficient resource allocation and timely response to user requests. Metrics such as response time, resource utilization, and error rates are analyzed and optimized.

Error Handling: In the event of recognition or translation errors, the system provides clear error messages to users, allowing them to rephrase or correct their input. Error handling mechanisms ensure a smooth user experience.

Analysis and Reporting: The project's output also includes the capability to generate reports and analyze usage data. This is particularly valuable for organizations using the system for customer service or multilingual support.
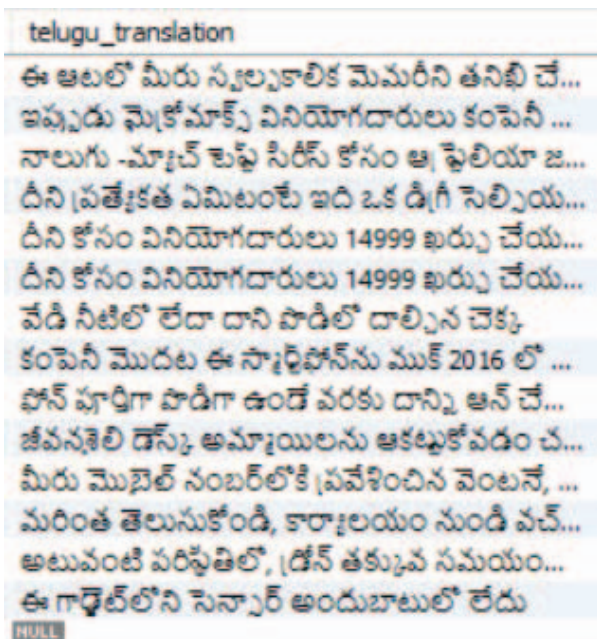
**Fig 6: Various audio files are Translated in Telugu**

The proposed output aims to bridge language barriers and enable effective communication in multilingual settings. Whether it is for international business negotiations, healthcare consultations, or educational purposes, the system's output enhances cross-cultural interactions and opens new opportunities for global collaboration.

## V. CONCLUSION

The proposed presents an innovative solution for seamless multilingual speech recognition and translation, breaking down language barriers. Its accuracy in speech recognition and nuanced text translations empowers users across diverse linguistic backgrounds. The user-friendly interface ensures accessibility, regardless of technical proficiency. Furthermore, the project prioritizes stringent data privacy and security measures, aligning with data protection regulations and fostering user trust. The proposed work is evolving the addition of real-time translation, advanced recognition models, and broader language support promises an even more dynamic and adaptable tool for multilingual communication.

In summary, the project's unique blend of speech recognition and text translation algorithms, along with its unwavering commitment to user-friendliness and data security, positions it as an asset in promoting cross-lingual understanding and accessibility. Its potential for growth and adaptability ensures its relevance in the ever-changing landscape of multilingual communication.

## VI. FUTURE ENHANCEMENT

The "Multilingual Speech Recognition in Cloud" project offers a strong foundation for multilingual communication but also presents several opportunities for future enhancements and expansions to meet evolving user needs and technological advancements, Additional Target Languages, Dialect Support, Voice Commands and Controls, Customizable Translation Models, Real-Time Collaboration, Integration with Mobile Applications, Improved Natural Language Processing (NLP), Adaptive Learning, Cloud Resource Optimization, Transcription for Multimodal Data, Accessibility Features: Incorporating accessibility features for users with disabilities, such as voice-guided navigation and sign language interpretation, will enhance inclusivity.

### REFERENCES

[1] Diez, M., & Varona, A. (2018). Multilingual speech recognition and transcription for historical documents using cloud computing. Procedia Computer Science, 130, 56-62.

[2] Collobert, R., & Weston, J. (2008). A unified architecture for natural language processing: Deep neural networks with multitask learning. In Proceedings of the 25th international conference on Machine learning (ICML'08) (pp. 160-167).

[3] Miao, Y., Gowayyed, M. A., & Metze, F. (2015). EESEN: End-to-end speech recognition using deep RNN models and WFST-based decoding. In Sixteenth annual conference of the international speech communication association.

[4] Hadian, H., Povey, D., & Vesely, K. (2018). Multilingual chain acoustic modeling for low-resource languages. In 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (pp. 4859-4863).

[5] Pal, M., & Plahl, C. (2015). Multilingual speech recognition for low resource languages using an extensible ASR toolkit. In 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (pp. 4975-4979).

[6] Cui, L., Zhu, X., & Chen, H. (2017). Multi-level knowledge driven recurrent neural networks for multilingual speech recognition. In 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (pp. 4935-4939).

[7] Ramabhadran, B., Hakkani-Tür, D., Bennett, C., Beaufays, F., Fernandez, R., & Plahl, C. (2012). The 2012 Babel system: The technical report. IARPA Babel Program.

[8] Vesely, K., Ghoshal, A., Burget, L., Povey, D., Rendel, A., Saon, G., ... & Zeyer, A. (2016). The Kaldi speech recognition toolkit. In 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (pp. 5180-5184).

[9] Xue, X., Espy-Wilson, C. Y., & Droppo, J. (2013). An investigation of multilingual deep neural networks for large vocabulary speech recognition. In 2013 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (pp. 7715-7719).

[10] Karimunnisa, S., Pachipala, Y. An AHP based Task Scheduling and Optimal Resource Allocation in Cloud Computing (2023) International Journal of Advanced Computer Science and Applications, 14 (3), pp. 149-159.

[11] Karimunnisa, S., Pachipala, Y. Task Classification and Scheduling Using Enhanced Coot Optimization in Cloud Computing (2023) International Journal of Intelligent Engineering and Systems, 16 (5), pp. 501-511.

[12] Ganesan, V., Sobhana, M., Anuradha, G., Yellamma, P., Devi, O.R., Prakash, K.B., Naren, J. Quantum inspired meta-heuristic approach for optimization of genetic algorithm (2021) Computers and Electrical Engineering, 94, art.

[13] Annamareddy, N., Donepudi, L.G., Parvathaneni, L., Brahma Rao, K.B.V., Putta, J., Yellamma, P. Comparison of Various Face Recognition Algorithms in ML/DS (2023) 2nd International Conference on Sustainable Computing and Data Communication Systems, ICSCDS 2023 - Proceedings, pp. 126-131.

[14] Kousik, K.V., AsishTony, M., Krishna, K.S., Narisety, S., Yellamma, P. An E-Commerce Product Feedback Review using Sentimental Analysis (2023) 6th International Conference on Inventive Computation Technologies, ICICT 2023 - Proceedings, pp. 608-613.

[15] Sarwar, S. M., & Babu, K. S. (2015). Multilingual speech recognition using deep learning. In 2015 International Conference on Signal Processing and Communications (SPCOM) (pp. 1-6).

571