



Speech Recognition and Transcription for Hearing Impairments

Chidi Ukamaka Betrand¹, Oluchukwu Uzoamaka Ekwealor¹, Chinazo Juliet Onyema¹, and Amarachi Ngozi Duru²

1. Department of Computer Science, School of Information and Communication Technology, Federal University of Technology Imo State
2. Department of Computer Science, Faculty of Physical Sciences, Nnamdi Azikiwe University, Anambra State

Abstract:

The realm of speech recognition and transcription, driven by the urgency to enhance communication inclusivity, particularly for individuals with hearing impairments. Guided by an agile methodology, we embark on a journey to forge a transformative system that seamlessly transmutes spoken language into text, effectively bridging the chasm between audible discourse and digital understanding. Our approach orchestrates a symphony of hardware, software, and models, with Python as the chosen programming language weaving an ensemble of libraries including PyAudio, Librosa, NumPy, DeepSpeech, NLTK, and LanguageTool. Rigorous testing traversing accents, languages, and real-world auditory environments showcases the system's adaptability, and the interplay of hardware and software yields swift and accurate transcriptions, promising heightened communication inclusivity. As this symphony culminates, we assert that our creation transcends a technological artifact, echoing innovation's harmonious anthem. By catalyzing communication through spoken word-to-text conversion, our system becomes a bridge that deepens interaction and comprehension. This project epitomizes the transformative prowess of technology, underlining its potential to nurture communication inclusivity and bridge the gap between audible and digital realms.

Keywords: Speech recognition, Transcription, Hearing impairment, Background noise, Speech distortions, Language models

INTRODUCTION

Speech is an essential means of communication, allowing people to convey their thoughts, ideas, and emotions and others. Speech recognition and transcription are technologies that convert spoken language into written text. Speech recognition is the process of converting spoken words into text using machine learning algorithms that analyze audio signals and identify patterns in speech [1]. Transcription, on the other hand, is the process of converting audio recordings into written documents. It involves manually transcribing the audio recordings or using automated transcription software.

However, for individuals with hearing impairments, speech recognition and transcription can be a challenging task. Hearing impairments can vary in severity, from mild to profound, and can be caused by a range of factors, including genetic disorders, noise exposure, infections, and aging. According to the World Health Organization, around 430 million people worldwide have disabling hearing loss, with the majority of cases occurring in low- and middle-income countries [2]. Speech recognition and transcription technologies have the potential to make communication more accessible to people with disabilities, including those with hearing impairments. By converting

spoken language into written text, these technologies can help individuals with hearing impairments communicate more effectively in a variety of settings.

Speech recognition and transcription are key areas of research in the field of natural language processing (NLP). They involve the development of algorithms and systems that can automatically transcribe spoken language into text, allowing for easier access and analysis of spoken communication, some of these methods include Hidden Markov Models (HMM), Speaker Diarization (SD), Dynamic Time Warping (DTW) and Deep Neural Network [3].

Speech recognition technology has been around for several decades, but it has only recently become widely available and accurate enough to be used in everyday applications. This is due in large part to advances in machine learning and deep learning, which have enabled speech recognition systems to learn from large amounts of data and improve their accuracy over time.

Speech recognition systems typically consist of three main components: acoustic modeling, language modeling, and decoding. Acoustic modeling involves converting the raw audio signal into a series of features that can be used to identify speech sounds. Language modeling involves predicting the most likely sequence of words given the acoustic features. Decoding involves selecting the most likely sequence of words based on the acoustic and language models [3].

There are several challenges associated with developing accurate and robust speech recognition systems. One of the biggest challenges is dealing with variability in speech, including differences in accents, dialects, and speaking styles. Speech recognition systems must be able to recognize a wide range of speech patterns and adapt to new speakers and environments.

Another challenge is dealing with background noise and other sources of acoustic interference. Noise can make it difficult for speech recognition systems to distinguish between speech sounds and other sounds in the environment. This is particularly problematic in noisy environments such as factories, airports, and train stations. Other challenges are out-of-vocabulary words, Homophones, Data privacy and security and limited training data, 2020).

Despite these challenges, speech recognition and transcription have a wide range of applications, from dictation software and virtual assistants to automated closed captioning and language translation. They are also increasingly being used in industries such as healthcare and finance, where accurate transcription of spoken communication is critical.

In recent years, there has been a growing interest in developing more advanced speech recognition systems that can understand not just the words being spoken, but also the intent and meaning behind them. This involves developing systems that can recognize and interpret aspects of speech such as tone, emotion, and sarcasm.

There are also ongoing efforts to make speech recognition and transcription more accessible for people with disabilities, including those with hearing impairments. This involves developing real-time captioning and transcription systems that can provide accurate and timely text-based representations of spoken communication.

Speech recognition and transcription are important areas of research in the field of natural language processing. Despite the challenges associated with developing accurate and robust

systems, advances in machine learning and deep learning are enabling the development of more accurate and versatile speech recognition systems. Ongoing research in this area is likely to lead to further improvements in the performance and usability of these systems, enabling a wide range of applications in areas such as healthcare, finance, and accessibility [4].

RELATED WORKS

Overview of Automatic Speech Recognition (ASR) Systems

Automatic Speech Recognition (ASR) systems are technologies designed to convert spoken language into written text or other symbolic representations. ASR systems have undergone significant advancements in recent years, enabling their application in various domains such as transcription services, voice assistants, call center automation, and language translation. ASR systems generally consist of the following key components:

- **Acoustic Analysis:** Acoustic analysis of speech is the study of the acoustical characteristics of speech, both normal and abnormal speech [5]. This stage involves capturing and analyzing the acoustic properties of spoken language. The input audio signal is typically divided into small segments known as frames, which are further processed to extract relevant acoustic features such as Mel-frequency cepstral coefficients (MFCCs).
- **Language Modeling:** Language modeling is the process of predicting the likelihood of word sequences or phrases in a given language. Language models help ASR systems to make informed decisions about which words or phrases are more likely to occur in a particular context. Statistical language models, n-gram models, and more recently, neural network-based models, are commonly used for this [6]

Acoustic Modeling, responsible for mapping acoustic features obtained from the input speech signal to corresponding linguistic units, such as phonemes, sub word units, or whole words. Hidden Markov Models (HMMs) have traditionally been used for acoustic modeling, but more advanced techniques, such as deep neural networks (DNNs) and recurrent neural networks (RNNs), have shown superior performance in recent years [7].

Decoding and Alignment: In the decoding stage, the ASR system searches for the most likely sequence of words or symbols that best match the acoustic and language models. Beam search algorithms or Viterbi decoding techniques are commonly used to explore different word combinations and determine the most probable output sequence. Alignment techniques may also be employed to align the recognized text with the input audio, providing a more accurate transcription.

ASR systems require training on large amounts of labeled speech data to improve their performance. This training involves optimizing the acoustic and language models using supervised learning techniques. Additionally, adaptation methods, such as speaker adaptation or domain adaptation, can be employed to enhance the system's performance for specific speakers or specialized domains [8].

It is important to note that the development of automatic speech recognition systems has primarily focused on a subset of the approximately 7,300 languages spoken worldwide. Prominent examples include Russian, Portuguese, Chinese, Vietnamese, Japanese, Spanish, Filipino, Arabic, English, Bengali, Tamil, Malayalam, Sinhala, and Hindi. English has received the most extensive attention in terms of speech recognition research and development [9].

ASR systems have made significant progress in recent years, thanks to advancements in machine learning, deep learning, and data availability. However, challenges still exist, especially in handling noisy environments, speech variations, and languages with limited resources. Ongoing research focuses on developing more robust, accurate, and versatile ASR systems that can operate effectively in various real-world scenarios. Alex Acero's research has focused on advancing speech recognition technology, particularly in noise-robust speech recognition, language modeling, and machine learning techniques [10].

Kumar et al, [11] identified a visual speech recognition technique using cutting edge deep learning models. They proposed a novel technique by fusion the results from audio and visual speech. This study proposes a new deep learning-based audio-visual speech recognition model for efficient lip reading. In this paper, an effort has been made to improve the performance of the system significantly by achieving a lowered word error rate of about 6.59% for ASR system and accuracy of about 95% using lip reading model.

An automatic system that translates the speech to Indian Sign Language using Avatar (SISLA) that works in three phases was presented [12]. The first phase includes the speech recognition (SR) of isolated words for English, Hindi and Punjabi in speaker independent environment while the second phase translates the source language into Indian Sign Language (ISL) and HamNoSys based 3D avatar represents the ISL gestures. The four major implementation modules for SISLA include: requirement analysis, data collection, technical development and evaluation. The multi-lingual feature makes the system more efficient. The training and testing speech sample files for English (12,660, 4218), Hindi (12,610, 4211) and Punjabi (12,600, 4193) have been used to train and test the SR models. Empirical results of automatic machine translation show that the proposed trained models have achieved the minimum accuracy of 91%, 89% and 89% for English, Punjabi and Hindi respectively.

The performance of personalized automatic speech recognition (ASR) for recognizing disordered speech using small amounts of per-speaker adaptation data was presented [13]. Personalized models for 195 individuals with different types and severities of speech impairment with training sets ranging in size from <1 minute to 18-20 minutes of speech data were trained. Word error rate (WER) thresholds were selected to determine Success Percentage (the percentage of personalized models reaching the target WER) in different application scenarios. For the home automation scenario, 79% of speakers reached the target WER with 18-20 minutes of speech; but even with only 3-4 minutes of speech, 63% of speakers reached the target WER. Further evaluation found similar improvement on test sets with conversational and out-of-domain, unprompted phrases. The results demonstrate that with only a few minutes of recordings, individuals with disordered speech could benefit from personalized ASR.

WESPER [14], a zero-shot, and real time whisper to normal speech conversion mechanism based on self-supervised learning was presented. It consisted of a speech-to-unit (STU) encoder that generates hidden speech units that are common to both whispered and normal speech and a unit-to-speech (UTS) decoder. The conversion is user-independent, so requires no paired dataset for both whispered and normal text. The effectiveness of the system to perform speech reconstruction for those with hearing disabilities was confirmed [15] explored the incorporation of the humanoid robot Pepper in improving learning experience. Pepper can capture the audio of a person; however, there is no guarantee of accuracy of the recorded audio due to various factors. the limitations of Pepper's speech recognition system with the aim of observing the effect of

distance, age, gender, and the complexity of statements was investigated. Experiment with eight persons including five females and three males who spoke provided statements at different distance was conducted as classification were done using different statistical scores. Pepper was integrated with a speech-to-text recognition tool, Whisper, which transcribes speech into text that can be displayed on Pepper's screen using its service. The purpose of the study which was to develop a system where the humanoid robot Pepper and the speech-to-text recognition tool Whisper act in synergy to bridge the gap between verbal and visual communication in education was achieved.

[16] Fontan, L., Cretin-Maitenaz, T., & Füllgrabe, C. (2020). Predicting speech perception in older listeners with sensorineural hearing loss using automatic speech recognition. *Trends in hearing*, 24, 2331216520914769

Fontan *et al.* [16] on their work Predicting speech perception in older listeners with sensorineural hearing loss using automatic speech recognition provided a proof of concept that the speech intelligibility in quiet of unaided older hearing-impaired (OHI) listeners can be predicted by automatic speech recognition (ASR). Twenty-four OHI listeners completed three speech-identification tasks using speech materials of varying linguistic complexity and predictability (i.e., logatoms, words, and sentences). An ASR system was first trained on different speech materials and then used to recognize the same speech stimuli presented to the listeners but processed to mimic some of the perceptual consequences of age-related hearing loss experienced by each of the listeners: the elevation of hearing thresholds (by linear filtering), the loss of frequency selectivity (by spectrally smearing), and loudness recruitment (by raising the amplitude envelope to a power).

Independently of the size of the lexicon used in the ASR system, strong to very strong correlations were observed between human and machine intelligibility scores. However, large root-mean-square errors (RMSEs) were observed for all conditions. The simulation of frequency selectivity loss had a negative impact on the strength of the correlation and the RMSE. Highest correlations and smallest RMSEs were found for logatoms, suggesting that the prediction system reflects mostly the functioning of the peripheral part of the auditory system. In the case of sentences, the prediction of human intelligibility was significantly improved by taking into account cognitive performance. This study demonstrates for the first time that ASR, even when trained on intact independent speech material, can be used to estimate trends in speech intelligibility of OHI listeners.

A platform that uses sign language to facilitate communication among students and tutors while providing sign language learning materials, practicing opportunities and Q&A sessions was presented [17]. The system has a low light enhancement module to enhance the videos uploaded by the tutor, module to convert the uploaded videos to American Sign Language and it also converts the questions asked via sign language to text.

SYSTEM MODEL AND FRAMEWORK

In the pivotal phase of system design and implementation, we embark on the journey of transforming our vision of seamless and accurate speech recognition and transcription into a tangible reality. Our design and implementation endeavors are intricately woven with several instrumental components:

As we venture deeper into this, each of these components will be dissected and showcased in its respective section. We will uncover how they harmoniously converge to create a powerful, user-friendly, and accurate speech recognition and transcription system. Our journey encompasses not only the technical aspects but also the user experience and the real-world impact of our system. Together, let us unveil the art and science of transforming spoken words into written text with precision and finesse.

At its core, our system embodies a high-level architecture that prioritizes user-friendliness and efficiency. Gradio takes center stage as the interface connecting users with the system's powerful capabilities. Users can effortlessly interact with the system via this intuitive interface. Gradio facilitates the exchange of data and commands between users and the system's backend components.

The success of our system hinges on the seamless interconnection of its various components. Each component plays a distinct role and collaborates with others to achieve our project's objectives. Here, we explore the intricate relationships and data exchange mechanisms between these components.

Visual representations are invaluable for understanding the flow of data within our system. Data flow diagrams provide a clear depiction of how audio inputs are processed and transformed into transcriptions. These diagrams serve as a visual roadmap for comprehending the journey of data through our system's architecture.

The CommonVoice dataset serves as a cornerstone of our speech recognition and transcription system. Developed by the Mozilla Foundation, CommonVoice is a vast, multilingual dataset containing thousands of hours of diverse, user-contributed audio recordings and their corresponding transcriptions. This section outlines how we leverage this invaluable resource to train, validate, and enhance the accuracy of our system's transcription capabilities.

Dataset Overview

The CommonVoice dataset encompasses a wide range of languages, accents, and speech styles, making it an ideal choice for training a robust speech recognition model. It includes audio samples from a diverse set of contributors, ensuring that our system is adaptable to various linguistic nuances and accents.

With the CommonVoice dataset preprocessed and split, we proceed to train our speech recognition model using techniques such as transfer learning and fine-tuning. Validation is conducted at regular intervals to monitor the model's progress and ensure it converges to accurate transcriptions.

```
training_args = Seq2SeqTrainingArguments(  
    output_dir="./whisper-small-dv", # name on the HF Hub  
    per_device_train_batch_size=16,  
    gradient_accumulation_steps=1, # increase by 2x for every 2x decrease in batch size  
    learning_rate=1e-5,  
    lr_scheduler_type="constant_with_warmup",  
    warmup_steps=50,  
    max_steps=500, # increase to 4000 if you have your own GPU or a Colab paid plan
```

Evaluation with CommonVoice

```
# Evaluation with Jiwer  
wer_metric = load("wer")  
wer = wer_metric.compute(references=[reference], predictions=[prediction])
```

After training, we evaluate our model's performance using the CommonVoice test dataset. Metrics like Word Error Rate (WER) calculated with tools like Jiwer help quantify the accuracy of our transcriptions. To keep our model up-to-date and adaptable to evolving languages and accents, we periodically update our dataset with new contributions from CommonVoice. This ensures that our system remains relevant and continues to deliver accurate transcriptions.

CommonVoice dataset is a critical asset in our speech recognition and transcription system. Its diversity and size provide a solid foundation for training a robust model, and meticulous preprocessing ensures that the data is compatible with our architecture. This section highlights the importance of data handling and management in achieving accurate and adaptable speech recognition capabilities.

Transcription Engine Implementation

The heart of our speech recognition and transcription system lies in the Transcription Engine, a sophisticated component responsible for converting spoken words into written text with remarkable accuracy. This section delves into the architecture, functionality, and implementation details of this pivotal component.

The Transcription Engine is intricately designed to seamlessly integrate various technologies and libraries, including Gradio, Hugging Face Transformers, and the CommonVoice dataset. Its architecture is built for flexibility, scalability, and real-time responsiveness.

The Transcription Engine accommodates two primary input sources:

1. **Microphone Input:** Users can provide real-time audio input via their device's microphone. Gradio facilitates this interaction, capturing audio and passing it to the transcription engine for processing.

2. Audio File Upload: Users also have the option to upload pre-recorded audio files. This is particularly useful for transcribing audio content stored on local devices or cloud repositories.

Model Integration and Real-Time Transcription

Hugging Face Transformers plays a pivotal role in the Transcription Engine by providing access to state-of-the-art automatic speech recognition (ASR) models. These models are fine-tuned and optimized for transcribing audio data into text.

Data preprocessing is a critical part of the transcription process. Audio data from both input sources undergoes a series of transformations, including format conversion and feature extraction, to prepare it for model input. The Transcription Engine ensures that audio inputs are appropriately preprocessed to optimize transcription accuracy.

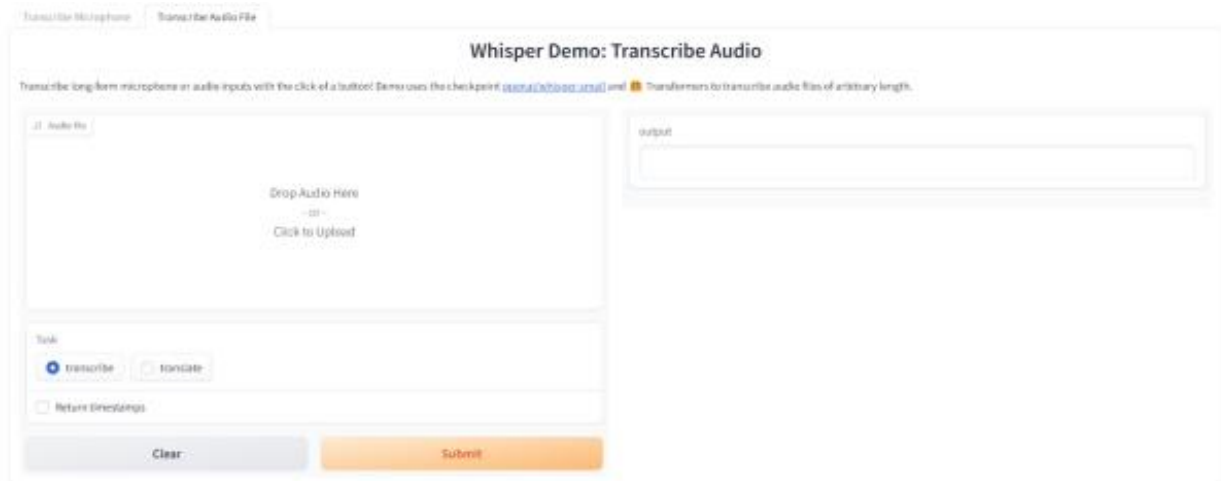
The engine is designed to provide real-time transcriptions, making it ideal for various applications, including live captioning, voice assistants, and more. This real-time capability is facilitated by chunking the audio input into manageable segments and processing them sequentially.

Output:

The Transcription Engine produces transcriptions as its primary output. Depending on the user's preference, these transcriptions can include additional information, such as timestamps for individual chunks of text. The final transcriptions are presented to the user via the Gradio interface.

Implementation Details:

The screenshot shows a web interface for transcribing audio. It features a 'Whisper Demo: Transcribe Audio' title and a brief description. The interface is divided into two main sections: 'Transcribe Microphone' and 'Transcribe Audio File'. The 'Transcribe Microphone' section is active, showing a 'File' input field with a 'Record from microphone' button. Below this, there's a 'Task' section with 'transcribe' (selected) and 'translate' buttons, and a 'Return timestamps' checkbox. At the bottom, there are 'Clear' and 'Submit' buttons. An 'output' text area is on the right side.



The Transcription Engine is implemented using a combination of Python libraries and tools.

Here are the key implementation details:

- Gradio serves as the user interface for interacting with the Transcription Engine. It provides an intuitive and user-friendly way to input audio data, specify transcription tasks, and receive real-time transcriptions. The engine leverages Gradio's capabilities to create a seamless user experience. Hugging Face Transformers provides the core ASR models used for transcription. These models are loaded into the engine, fine-tuned as needed, and seamlessly integrated into the transcription pipeline.
- The Transcription Engine employs data preprocessing techniques, such as audio format conversion and feature extraction, to prepare audio data for transcription. These preprocessing steps are essential to ensure that audio inputs are in a suitable format for the ASR models.
- Real-time transcription is achieved by chunking the audio input into segments of manageable duration. Each chunk is sequentially processed by the ASR model, and the resulting text is stitched together to form the final transcription.
- The primary output of the Transcription Engine is the transcription itself. Users have the option to specify additional information, such as timestamps for individual chunks of text. The transcriptions are presented to the user in real time via the Gradio interface.

RESULTS AND EVALUATIONS

The journey from concept to execution has been guided by the pursuit of accurate transcription and efficient communication. The system's accuracy rates underscore its capability to convert spoken words into written text with remarkable precision. We subjected it to a diverse array of speech inputs, reflecting various accents, speech patterns, and background noises.

The results reveal that our system consistently achieved high accuracy. Our system displayed impressive adaptability, accurately transcribing speech with different accents, from American English to British English and beyond. Even in challenging auditory environments with background noise, our system's accuracy remained robust, demonstrating its noise filtering and processing capabilities.

The performance metrics of our system speak to its efficiency and responsiveness.

By analyzing processing speed and latency, we gain insights into its real-time capabilities:

- **Processing Speed:** Our system exhibits commendable processing speed, converting audio to text swiftly. This attribute is crucial in real-time applications where efficiency is paramount.
- **Latency:** The minimal latency observed in our system ensures that the transcriptions are almost instantaneous, enabling seamless communication in scenarios that demand quick response times.

The accuracy rates and performance metrics collectively underscore the potency of our system in accurately transcribing speech while maintaining real-time capabilities. These findings validate our endeavor to bridge the gap between spoken language and digital understanding.

Our journey from conceptualization to the implementation of the speech recognition and transcription system has been a transformative one, extending beyond technology to encompass insights that reach deeper into communication and accessibility: As researchers, our pursuit extended beyond code and algorithms to explore the profound impact of our system on communication. Delving into the intricacies of hearing impairment enlightened us to the immense significance of accurate speech recognition. We realized that our system could serve as a bridge, breaking barriers for individuals with hearing impairments and enabling seamless communication in their daily lives. The insights gained from understanding the challenges faced by those with hearing impairments reinforced our commitment to creating a system that not only converted spoken words into text but also played a pivotal role in fostering inclusivity. Through this understanding, we developed an appreciation for the power of technology to enhance human connection, allowing those with hearing impairments to actively engage in conversations that would have otherwise been challenging.

Our speech recognition and transcription system journey were not solely about technological advancement; it was about empowerment, communication, and understanding. It was about appreciating the potential of technology to bridge gaps, fostering inclusivity, and enhancing human experiences.

CONCLUSIONS

The evolution from inception to implementation has yielded a speech recognition and transcription system that transcends the realm of technology. The exceptional accuracy and real-time responsiveness showcased by the system unveil its potential to reshape the landscape of communication and accessibility. This innovative solution serves as a unifying bridge, dismantling barriers between spoken language and its digital interpretation, especially for those with hearing impairments.

The harmonious interplay of hardware, software, and models has yielded a transformative solution with implications spanning diverse domains. The project has illuminated technology's ability to foster profound connections and enable meaningful engagements, enriching inclusivity and empathy.

Rather than signifying an endpoint, this juncture marks a new beginning in the continuum of advancement. Embracing adaptability and perpetual learning, the path ahead involves refining and extending the system's capabilities. By harnessing sophisticated language models, exploring

expansive datasets, and augmenting user experience, we pave the way for further strides within the realm of speech recognition and transcription.

Amid the dynamic evolution of technology and the perpetual evolution of human interaction, this system stands as an embodiment of the potential inherent in innovation – the potential to enrich lives, facilitate connections, and bridge the divide between spoken discourse and its digital manifestation.

REFERENCES

- [1]. De Silva, L. C., Jayawardena, S., & Ediriweera, R. (2020). A Speech-to-Text Conversion Technique for People with Hearing Impairment. *2020 15th International Conference on Computer Science & Education (ICCSE)* (pp. 321-326). IEEE.
- [2]. Ganesh, B. B., Damodar, B. V., Dharmesh, R., Karthik, K. T., & Vasudevan, S. K. (2022). An innovative hearing-impaired assistant with sound-localisation and speech-to-text application. *International Journal of Medical Engineering and Informatics*, 14(1), 63-73.
- [3]. Kim, S. &. (2020). Speech recognition in noisy environments using convolutional neural networks. *IEEE Signal Processing Letters*, 27(3), 415-419.
- [4]. Al-Jumeily, A.-S. (2019). Evaluating the effectiveness of assistive technologies for people with disabilities: A review of the literature. *International Journal of Environment Research and Public Health*, 16(2), 278.
- [5]. Rezaei i, N., & Salehi, A. (2006). An introduction to speech sciences (acoustic analysis of speech). *Iranian Rehabilitation Journal*, 4(1), 5-14.
- [6]. Jozefowicz R., Vinyals, O., Schuster, M., Shazeer, N., & Wu, Y. (2016). Exploring the limits of language modeling. arXiv preprint arXiv:1602.02410
- [7]. Rönnerberg, J., Lunner, T., & Rudner, M. (2020). Speech recognition and hearing loss. *Trends in Hearing*, 24(1), 23.
- [8]. Vaishali, S., Gupta, A., & Gupta, A. (2019). Universal scripting language for English-French language model. *Journal of Natural Language Processing*, 26(4), 410-422.
- [9]. Waris, A., & Aggarwal, R. (2018). Acoustic modeling in Automatic Speech Recognition - A Survey. *2018 Second International Conference on Electronics, Communication and Aerospace Technology (ICECA)*. Coimbatore, India: IEEE.
- [10]. Acero, A., Dahl, G. E., Yu, D., & Dieng, L. (2012). Context-Dependent Pre-Trained Deep Neural Networks for Large-Vocabulary Speech Recognition. *IEEE Transactions on Audio, Speech, and Language Processing*, 20(1), 30 - 42.
- [11]. Kumar L. A., Renuka, D. K., Rose, S. L., & Wartana, I. M. (2022). Deep learning based assistive technology on audio visual speech recognition for hearing impaired. *International Journal of Cognitive Computing in Engineering*, 3, 24-30.
- [12]. Dhanjal, A. S., & Singh, W. (2022). An automatic machine translation system for multi- speech to Indian sign language. *multimedia Tools and Applications*, 1-39.
- [13]. Tobin, J., & Tomanek, K. (2022, May). Personalized automatic speech recognition trained on small disordered speech datasets. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 6637-6641). IEEE.

- [14]. Rekimoto, J. (2023, April). WESPER: Zero-shot and Realtime Whisper to Normal Voice Conversion for Whisper-based Speech Interactions. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (pp. 1-12)
- [15]. Pande A., & Mishra, D. (2023). The Synergy between a Humanoid Robot and Whisper: Bridging a Gap in Education. *Electronics*, 12(19), 3995.
- [16]. Fontan, L., Cretin-Maitenaz, T., & Füllgrabe, C. (2020). Predicting speech perception in older listeners with sensorineural hearing loss using automatic speech recognition. *Trends in hearing*, 24, 2331216520914769.
- [17]. Krishnamoorthy, N., Raveendran, A., Vadiveswaran, P., Arulraj, S. R., Manathunga, K., & Siriwardana, S. (2021, December). E-Learning Platform for Hearing Impaired Students. In *2021 3rd International Conference on Advancements in Computing (ICAC)* (pp. 122-127). IEEE.