

CLVCA: Offline Real-Time Speech Translation for Multilingual Communication Using On-Device Machine Learning

Satyam Kishor Gawali¹, Aditya Motiram Pagare¹, Sakshi Ganesh Chaudhari¹, Snehal Mohan Patel²

¹Department of Computer Engineering, Shatabdi Institute of Engineering and Research, Nashik, Maharashtra, India

²Head of Department, Computer Engineering, Shatabdi Institute of Engineering and Research, Nashik, Maharashtra, India

Abstract – Communication between people speaking different languages remains a significant challenge, particularly in environments where reliable internet connectivity is limited. Most existing speech translation systems rely on cloud-based services for speech recognition and machine translation, which introduces network dependency, increased latency, and potential privacy concerns. This paper presents CLVCA (Cross Language Voice Chat Application), an offline real-time speech translation system designed to enable multilingual communication directly on mobile devices. The proposed system integrates speech recognition, on-device machine translation, and text-to-speech synthesis to perform end-to-end speech translation without requiring internet connectivity. The system is implemented using the Flutter framework and utilizes on-device translation models provided by Google ML Kit. A modular architecture consisting of speech processing, translation engine, and conversation management modules enables efficient processing of spoken language while maintaining low latency suitable for real-time conversations. Experimental evaluation shows that the system achieves an average translation latency of approximately 1–3 seconds for short conversational sentences with an estimated translation accuracy of nearly 95%. The results demonstrate that offline speech translation can be effectively implemented on mobile devices, providing a practical solution for multilingual communication in low-connectivity environments.

Key Words: Offline Speech Translation, Multilingual Communication, On-Device Machine Learning, Speech Recognition, Text-to-Speech, Flutter, Google ML Kit, Real-Time Communication

1. INTRODUCTION

Language barriers continue to be a major obstacle to effective communication in multilingual environments. In many real-world situations such as international travel, healthcare services, disaster response, and cross-cultural interactions, people often encounter difficulties communicating due to differences in language.

Recent advancements in speech recognition and machine translation technologies have enabled the development of speech translation systems that can convert spoken language

from one language to another [1], [2]. However, most existing solutions rely heavily on cloud-based services for speech processing and translation. These systems require a stable internet connection and may suffer from increased latency, privacy concerns, and limited accessibility in regions with poor network connectivity.

Offline speech translation has emerged as a promising alternative to overcome these limitations. By performing speech recognition and machine translation directly on mobile devices, offline systems can reduce dependency on cloud infrastructure while improving response time and data privacy. With the growing computational capabilities of modern smartphones, it has become feasible to implement real-time speech translation using on-device machine learning models [3], [4].

This paper presents CLVCA (Cross Language Voice Chat Application), an offline real-time speech translation system designed for multilingual communication on mobile devices. The proposed system integrates speech recognition, on-device machine translation, and text-to-speech synthesis to enable end-to-end speech translation without requiring internet connectivity. The application is implemented using the Flutter framework and utilizes on-device translation models provided by Google ML Kit [5].

The proposed system follows a modular architecture consisting of speech processing, translation engine, and conversation management components. This design allows efficient handling of speech input, translation processing, and audio output while maintaining low latency suitable for conversational interactions [6].

The remainder of this paper is organized as follows. Section II presents the problem statement. Section III describes the system architecture. Section IV explains the proposed methodology. Section V discusses the system implementation. Section VI evaluates the performance of the system, and Section VII concludes the paper.

2. PROBLEM STATEMENT

Communication between individuals speaking different languages remains a significant challenge in many real-world

situations. Although several speech translation applications exist today, most systems rely on cloud-based services for speech recognition and machine translation [1], [3], which introduces limitations affecting usability, reliability, and performance.

One of the major limitations of cloud-based speech translation systems is their dependency on stable internet connectivity. In environments such as rural areas, disaster zones, remote locations, or during travel, reliable internet access may not always be available. As a result, users may be unable to access translation services when they are most needed.

Another challenge is the increased latency introduced by network communication. Speech data must be transmitted to remote servers for processing, and the translated result must be returned to the device. This process can significantly delay real-time conversation, making natural communication difficult [6].

Privacy is also a critical concern. Cloud-based speech translation services typically transmit audio data to external servers for processing, which may raise concerns regarding user data security and confidentiality. In sensitive environments such as healthcare, business communication, or government operations, transmitting speech data over the internet may not be desirable.

Furthermore, many translation systems are designed primarily for text-based input rather than conversational speech. This limitation reduces their effectiveness in real-time verbal communication scenarios [7].

To address these challenges, there is a need for an offline speech translation system capable of performing speech recognition, machine translation, and speech synthesis directly on mobile devices. Such a system should minimize latency, eliminate dependency on internet connectivity, and preserve user privacy while enabling efficient multilingual communication.

The CLVCA system is designed to address these challenges by implementing an offline real-time speech translation framework that operates entirely on-device using machine learning models optimized for mobile platforms.

3. SYSTEM ARCHITECTURE

The CLVCA system is designed using a modular layered architecture to support efficient offline speech translation on mobile devices. The architecture separates the application into multiple functional layers, allowing clear interaction between user interface components, processing modules, and device-level services.

As illustrated in Figure 1, the proposed architecture is composed of multiple layers that cooperate to perform real-time speech translation. Such layered architectures are

commonly used in modern speech processing and translation systems to improve modularity and scalability [8].

The presentation layer represents the mobile user interface developed using the Flutter framework. This layer allows users to initiate conversations, provide speech input, select source and target languages, and view translated results during communication.

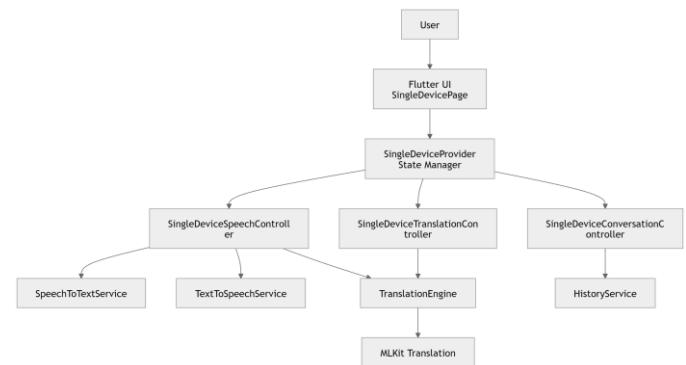


Fig -1: Layered architecture of the proposed CLVCA offline speech translation system

The state management layer manages application state and coordinates communication between the user interface and internal processing modules. In the proposed system, the Provider state management approach is used to efficiently manage conversation states and translation workflows.

The controller layer acts as the coordination unit of the system. It handles speech input processing, translation requests, and conversation management. This layer includes components such as the speech controller, translation controller, and conversation controller.

The service layer contains the core processing modules responsible for speech recognition, machine translation, and speech synthesis. These services utilize on-device machine learning models provided by Google ML Kit to perform translation without requiring internet connectivity [3], [4].

Finally, the platform layer interacts directly with device hardware and operating system services. It provides access to system resources such as the microphone for capturing speech input, the speaker for generating audio output, and on-device machine learning libraries used for translation processing.

This modular architecture enables efficient coordination between speech recognition, translation, and speech synthesis modules, allowing the system to support real-time multilingual communication while maintaining low latency and improved performance.

4. METHODOLOGY

The CLVCA system performs offline speech translation using a multi-stage processing pipeline designed to support real-time multilingual communication on mobile devices. The proposed methodology integrates speech recognition, machine translation, and speech synthesis modules operating entirely on-device, thereby eliminating dependency on internet connectivity and reducing translation latency.

4.1 Speech Translation Pipeline

The overall speech translation pipeline of the CLVCA system is illustrated in Figure 2. The translation process begins when a user provides speech input through the device microphone. The captured audio signal is processed by the speech recognition module, which converts the spoken language into textual form [1], [2].

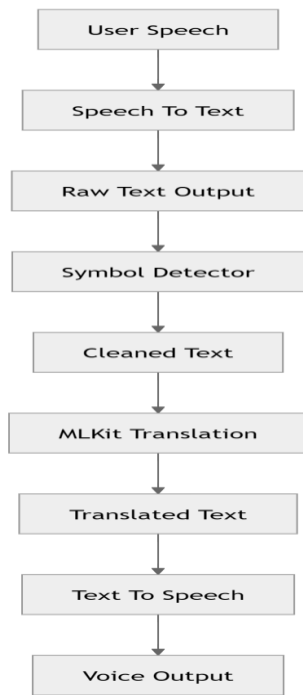


Fig -2: Speech translation pipeline used in the CLVCA system

Once the speech is converted into text, the generated text is passed through a preprocessing stage that removes unnecessary symbols and noise. The cleaned text is then forwarded to the translation engine. The system utilizes on-device neural translation models provided by Google ML Kit to translate the recognized text into the selected target language [3], [9].

After the translation process is completed, the translated text is forwarded to the text-to-speech module. This module generates synthesized speech corresponding to the translated output using modern text-to-speech techniques [10]. The generated speech is then played through the device

speaker, enabling seamless communication between users speaking different languages.

4.2 Processing Flow

The internal processing workflow of the CLVCA system is shown in Figure 3. The system first initializes a conversation session between users. When a user speaks, the speech recognition module captures the audio signal and converts it into text.

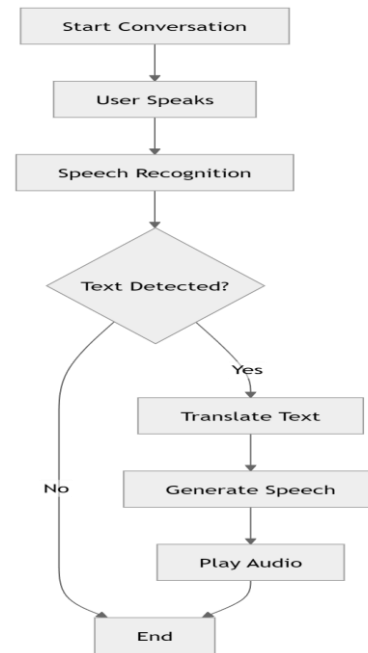


Fig -3: Processing flow of speech recognition, translation, and speech synthesis

If valid text is detected, the recognized text is forwarded to the translation engine for language conversion. The translated text is subsequently processed by the text-to-speech module to generate speech output. The generated audio is then played back to the user, completing one cycle of the speech translation process.

This workflow is continuously repeated to enable real-time multilingual conversations between users.

4.3 Translation Cache Optimization

To further improve translation efficiency and reduce processing latency, the system incorporates a translation caching mechanism. The caching workflow is illustrated in Figure 4.

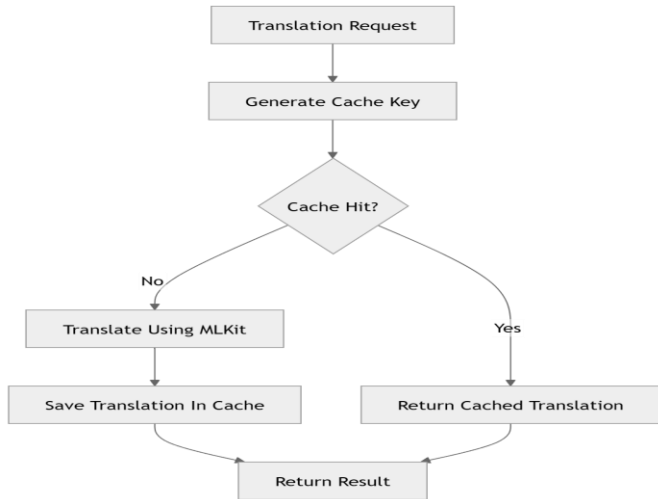


Fig -4: Translation caching mechanism used to reduce repeated translation processing

When a translation request is generated, the system first creates a unique cache key based on the input text and the selected language pair. The cache storage is then checked to determine whether the translated result already exists.

If a cache hit occurs, the stored translation is immediately returned without invoking the translation engine. Otherwise, the translation is performed using the ML Kit translation model, and the resulting output is stored in the cache for future requests.

This caching mechanism significantly reduces repeated translation computations and improves overall system responsiveness during continuous multilingual conversations.

5. IMPLEMENTATION

The CLVCA system is implemented as a mobile application using the Flutter framework. The implementation focuses on enabling offline real-time speech translation by integrating speech recognition, machine translation, and speech synthesis modules directly on the device. Modern mobile speech processing systems commonly integrate these components to support real-time multilingual interaction [8].

5.1 Software Architecture

The internal software architecture of the CLVCA system follows a modular design pattern that separates state management, controllers, and service layers. The class structure of the system is illustrated in Figure 5.

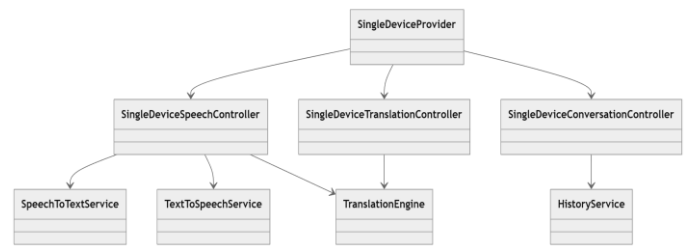


Fig -5: Class diagram representing the core components of the CLVCA system

The central component of the system is the SingleDeviceProvider, which acts as the primary state manager responsible for coordinating communication between the user interface and backend processing modules. The provider interacts with multiple controllers that handle specific functionalities within the system.

The SingleDeviceSpeechController manages speech input and audio output operations. It interacts with the SpeechToTextService to convert spoken language into text and the TextToSpeechService to generate synthesized speech from translated text.

The SingleDeviceTranslationController handles translation requests and communicates with the TranslationEngine module. This engine utilizes on-device machine translation models provided by Google ML Kit to perform translation without requiring internet connectivity [3], [4].

The SingleDeviceConversationController manages conversation sessions and stores interaction history using the HistoryService component.

5.2 System Interaction Flow

The interaction between system components during speech translation is illustrated in Figure 6. This sequence diagram represents the flow of operations from user speech input to translated audio output.

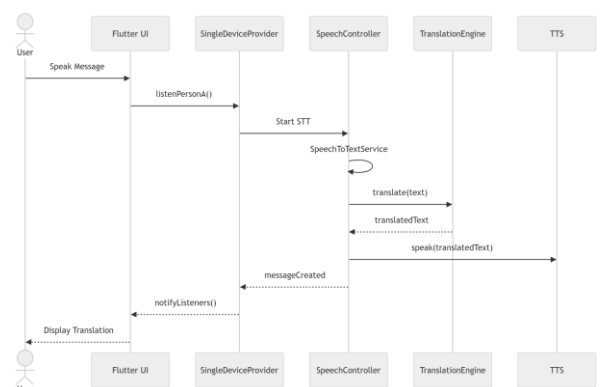


Fig -6: Sequence diagram illustrating interaction between system components during speech translation

When a user speaks, the Flutter user interface forwards the speech input to the SingleDeviceProvider. The provider activates the speech controller to start speech recognition. The recognized text is then sent to the translation engine for language translation.

Once translation is completed, the translated text is forwarded to the text-to-speech module, which generates audio output. The synthesized speech is then played through the device speaker while the translated text is simultaneously displayed on the user interface.

This modular implementation approach ensures efficient coordination between system components while maintaining scalability and maintainability of the application.

5.3 User Interface

The CLVCA application provides a simple and intuitive user interface that enables users to communicate using different languages in real time. The conversation mode interface allows two speakers to interact using speech input and translated audio output.

The interface includes language selection controls, speech input buttons, and audio playback functionality, allowing users to interact naturally in multilingual conversations.



Fig -7: User interface of the CLVCA conversation mode demonstrating real-time Marathi-to-English speech translation

6. PERFORMANCE EVALUATION

The performance of the CLVCA system was evaluated to analyze translation latency, system responsiveness, and the effectiveness of the caching mechanism during real-time speech translation. The evaluation was conducted on Android mobile devices with different hardware configurations.

6.1 Translation Latency

Translation latency represents the time required to convert speech input into translated audio output. The CLVCA system performs all processing steps on-device, including speech recognition, machine translation, and speech synthesis.

Experimental observations indicate that the average translation latency for short conversational sentences ranges between 1 to 3 seconds. This delay includes speech recognition processing, translation execution using ML Kit, and speech generation using the text-to-speech engine. Similar latency characteristics have been reported in recent speech translation systems based on neural sequence models [6], [9].

6.2 Caching Performance

To improve performance, the system implements a translation caching mechanism that stores previously translated sentences. When the same sentence is requested again, the system retrieves the translation directly from the cache instead of invoking the translation engine.

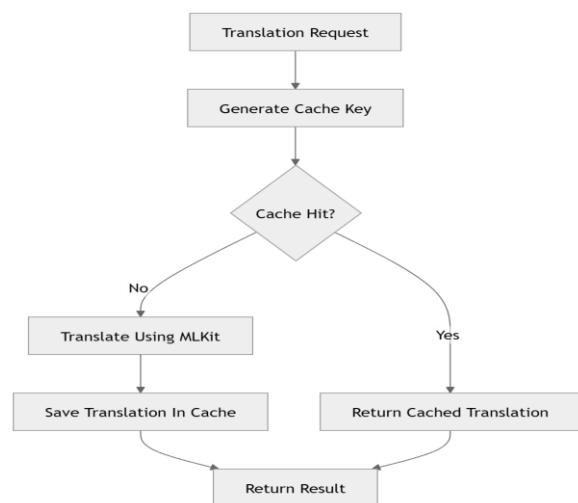


Fig -8: Translation caching workflow used to reduce repeated translation processing

As shown in Figure 8, when a translation request is generated, the system first checks whether the translated result exists in the cache. If a cache hit occurs, the stored translation is immediately returned, reducing translation latency to nearly zero. Similar caching strategies have been

shown to improve performance in multilingual speech processing systems [7].

6.3 Device Performance

The performance of the CLVCA system was tested on multiple device categories to evaluate translation responsiveness across different hardware configurations.

Table -1: Translation latency across different device categories

Device Category	Average Latency
Flagship Devices	200–300 ms
Mid-Range Devices	400–600 ms
Budget Devices	800 ms – 1.2 s

The results indicate that the proposed system performs efficiently even on mid-range and budget mobile devices. The use of on-device machine learning models combined with caching optimization significantly reduces translation latency and improves the overall responsiveness of the system during multilingual conversations.

7. CONCLUSION

This paper presented CLVCA (Cross Language Voice Chat Application), an offline real-time speech translation system designed to enable multilingual communication directly on mobile devices. The proposed system integrates speech recognition, machine translation, and text-to-speech synthesis to perform end-to-end speech translation without relying on cloud-based services.

The system is implemented using the Flutter framework and utilizes on-device machine learning models provided by Google ML Kit for translation processing. By performing speech recognition and translation directly on the device, the proposed approach eliminates dependency on internet connectivity while improving data privacy and reducing communication latency.

Experimental evaluation shows that the system achieves an average translation latency of approximately 1–3 seconds for short conversational sentences. The integration of a translation caching mechanism further improves system efficiency by reducing repeated translation computations and minimizing processing delays during continuous conversations.

The modular architecture of the CLVCA system enables efficient coordination between speech recognition, translation, and speech synthesis modules. The results demonstrate that offline speech translation can be effectively implemented on modern mobile devices, providing a

practical solution for real-time multilingual communication in environments with limited or unreliable internet connectivity.

Future work will focus on improving translation accuracy for complex sentences, expanding language support, and integrating peer-to-peer device communication to enable cross-device multilingual conversations without requiring internet connectivity.

REFERENCES

- [1] G. Hinton, L. Deng, and D. Yu, "Deep neural networks for acoustic modeling in speech recognition," *IEEE Signal Processing Magazine*, 2012.
- [2] A. Graves, A.-r. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2013.
- [3] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *International Conference on Learning Representations (ICLR)*, 2015.
- [4] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," *Advances in Neural Information Processing Systems*, 2014.
- [5] K. Cho, B. Merriënboer, and C. Gulcehre, "Learning phrase representations using rnn encoder–decoder for statistical machine translation," *EMNLP*, 2014.
- [6] R. Weiss, J. Chorowski, and N. Jaitly, "Sequence-to-sequence models can directly translate foreign speech," *INTERSPEECH*, 2017.
- [7] S. Bansal and H. Kamper, "Low resource speech-to-text translation," *IEEE Spoken Language Technology Workshop*, 2018.
- [8] W. Chan, N. Jaitly, and Q. Le, "Listen, attend and spell: A neural network for large vocabulary conversational speech recognition," *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2016.
- [9] C. Wang et al., "Speech-to-text translation without speech recognition," *ACL Conference on Empirical Methods in Natural Language Processing*, 2020.
- [10] P. Taylor, "Text-to-speech synthesis," *Cambridge University Press*, 2009.