# Bridging Language Barriers: The Role of AI in Real-Time Multilingual Translation for Video Conferencing

Pratiksha Patare[1], Prachi Said[2], Afrin Shaikh[3], Prof. Shweta Shah[4]

[1,2,3] *Student, Department of Computer Engineering, SCTR's Pune Institute of Computer Technology, Pune, Maharashtra, India*

[4]*Professor, Department of Computer Engineering, SCTR's Pune Institute of Computer Technology, Pune, Maharashtra, India*

*Abstract—With the rise of global communication the need for effective multilingual translation in video conferencing has become crucial to break language barriers. Traditional video conferencing tools lack real time multilingual translation capa- bilities, limiting themselves in international collaborations. We present Bridging language barriers: The role of AI in real-time multilingual translation for video conferencing a video conferencing system with live audio translation. This system facilitates communication between two people who speak different languages bridging the language barriers. Using AI based language models,our system captures live audio, process translations, and delivers them in real time reducing the latency up to 4-5 seconds andsometimes 3 second. The platform makes use of Neural Machine Translation NMT Algorithm which is a type of deep learning architecture which can learn to translate language.*

*Index Terms—Video Conferencing, AI, WebRTC, Real-time Language Translation, MERN Stack, NLP, Facial Recognition, Noise Cancellation*

## I. INTRODUCTION

In today's globalized world, video conferencing has become an essential tool for communication, facilitating interactions across geographical and linguistic barriers. However, as the de- mand for international collaboration increases, existing video conferencing platforms face significant limitations, particularly in the areas of multilingual communication, background noise interference, and latency. Effective communication in multi- lingual settings is critical not only for businesses but also for educational institutions, healthcare providers, and diplomatic missions, where precision and clarity in communication are paramount. Unfortunately, current platforms offer limited support for real-time translation, often leading to misunder- standings and inefficient collaboration when participants speak different languages.

Moreover, the challenges of communication are exacer- bated by background noise and high latency. Background

Noise whether from bustling office environments, home set- tings, or public places can severely disrupt the flow of conversation, making it difficult for participants to follow and understand each other.

This research aims to overcome these limitations by de- veloping a comprehensive video conferencing system that integrates real-time multilingual translation, advanced noise cancellation, and significantly reduced latency. By leveraging cutting edge AI and machine learning techniques, the system delivers accurate translations across multiple languages, intel- ligently filters out background noise, and ensures seamless, low latency communication. This platform empowers users to engage in conversations without concerns about language differences, distractions, or delays, fostering more productive collaboration across borders.

The platform is built using the MERN stack ensuring scalability, robustness, and high performance. Real-time com- munication is facilitated through WebRTC technology, al- lowing peer to peer connections that ensure smooth video and audio transmissions with minimal delay. For multilingual communication, the system uses Natural Language Processing(NLP) and Speech Recognition APIs for real-time translation, transcription, and meeting summarization.

Additionally, the system utilizes several state-of-the-art toolsto further enhance its capabilities. Next.js is employed for building server-side rendered applications, improving perfor- mance and SEO. TypeScript ensures strong typing, enhancing code reliability and maintainability. LiveKit serves as the infrastructure for real-time video and audio, and Pusher powers real-time updates and bidirectional communication via WebSockets. The platform also

utilizes TailwindCSS for responsive, utility first design and PlanetScale for a highly scalable, distributed database solution. To manage the platform's data with precision and type safety, Prisma ORM is used, interfacing with MySQL for efficient data handling.

The platform's real-time language translation is powered by the Google-translate-browser API for accurate multilingual conversion, while the Web Speech API handles transcription and text-to-speech synthesis. OneAI's Summarizer provides intelligent meeting summaries, aiding productivity by delivering concise insights from lengthy conversations.

## II. LITERATURE SURVEY

| Area of Study | Paper Name (Year) | Problem Identified | Solution |
|---|---|---|---|
| WebRTC | WebRTC: APIs and RTCWEB Protocols of the HTML5 Real-Time Web (2012) [1] | High latency and poor quality communication. | Integrate WebRTC for low-latency video/audio. |
| Noise Suppression | Deep Complex Networks for Real-Time Noise Suppression (2019) [2] | Background noise disrupts conversations. | Use deep learning models for effective noise cancellation. |
| Facial Recognition | DeepFace: Closing the Gap to Human-Level Performance in Face Verification (2014) [3] | Insecure authentication and low engagement. | Implement facial recognition for secure login and emotion detection. |
| Real-Time Translation | Google's Multilingual Neural Machine Translation System (2017) [4] | Language barriers hinder communication. | Employ NMT for instant multilingual translation. |
| Automated Summarization | Towards a Real-Time Meeting Assistant (2005) [5] | Inefficient manual note-taking. | Use NLP and speech-to-text for automatic meeting summaries. |
| Low-Latency Models | Low-Latency Machine Learning Models for Real-Time Communication (2022) [6] | Delays in audio/video affect user experience. | Implement low-latency ML models for smooth communication. |
| Speech Translation | Latency-Optimized Speech Translation with Noise Reduction (2019) [7] | Poor translation accuracy in noisy environments. | Combine noise reduction with real-time translation. |
| Adaptive Filtering | Adaptive Filtering for Real-Time Noise Reduction (2017) [8] | Inconsistent audio quality across users. | Use adaptive filtering for dynamic noise suppression. |
| Translation Challenges | Challenges in Real-Time Machine Translation (2020) [9] | Trade-offs between speed and accuracy. | Optimize models for efficient, accurate translations. |
| Latency Minimization | Latency Minimization Techniques in Multilingual Communication (2018) [10] | High latency in multilingual settings. | Use optimized algorithms to minimize latency. |

TABLE I
LITERATURE SURVEY

## III. EXISTING SYSTEM

Current video conferencing platforms have made signifi- cant advancements in virtual communication, providing high- quality audio and video. However, their ability to address critical issues such as real-time multilingual translation, effec-tive noise cancellation, and reduced latency is limited. Several studies highlight these shortcomings and the gaps in current systems.

- Some platforms have begun integrating translation fea- tures, such as Skype Translator and Google Meet's live captions, but these are rudimentary and often lack the precision and fluency needed for professional communi- cation. Research by Zhang et al. (2021) and Lee Kim (2022) discusses how current machine translation models struggle with real-time processing, contextual accuracy, and handling of domain-specific jargon or regional di-alects.

- Many platforms offer basic noise suppression features. These platforms employ AI-driven noise cancellation al- gorithms, but research by Kumar Gupta (2017) indicates that these systems struggle with non-repetitive and com- plex background noises, such as construction sounds or multiple overlapping voices in shared spaces. While noise reduction has improved, its effectiveness decreases when dealing with more unpredictable noise sources.

- Latency remains one of the most challenging problems in video conferencing, particularly when real-time multilin- gual translation is involved. The

time required for speechrecognition, translation, and synthesis increases delays, causing unnatural pauses in conversations. According to Wang et al. (2021), many video conferencing platforms are unable to sufficiently reduce latency during complex tasks like translation, particularly in environments with suboptimal network conditions.

## IV. METHODOLOGY

### A. Proposed System

To address the challenges in the existing systems in video conferencing, this paper proposes a comprehensive, integratedsystem that combines AI-driven technologies and optimized algorithms. The core of the proposed system is an advanced multilingual translation engine, powered by a hybrid model that combines Neural Machine Translation (NMT) and Statistical Machine Translation (SMT). This hybrid approach leverages the strengths of both models: NMT excels at understanding context and semantics through deep learning, while SMT provides robust performance in handling domain-specific terminology and less common language pairs.

The system captures speech in real-time, converting it to text using advanced speech recognition models that utilize state-of-the-art techniques, such as end to end deep learning architectures. This process ensures high accuracy in tran- scribing spoken language, even in environments with varying levels of background noise. Following transcription, the text is translated into the target language with minimal delay, thanks to a carefully designed pipeline that minimizes processing timeat each stage. The translation module is trained using vast multilingual datasets, including diverse sources like academic papers, conversational transcripts, and technical documents, to improve contextual understanding and maintain accuracy across various topics and dialects.

To tackle the persistent issue of background noise, the system incorporates AI-driven noise cancellation algorithms that operate in tandem with the translation module. The noise cancellation feature employs advanced deep learningmodels such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) to effectively filter out unwanted noise from audio streams in real-time. These models are designed to recognize and separate human speech from various noise sources such as traffic, keyboard typing, and environmental sounds while preserving speech clarity. The system leverages a combination of frequency domain analysis and time frequency masking techniques, allowing it to adapt dynamically to changes in the acoustic environment.

Moreover, the integration of these noise cancellation al- gorithms is facilitated by a feedback loop that continuously analyzes audio input, enabling the system to adjust its filtering strategies in real-time. This ensures that users experience crystal-clear audio even in challenging settings, such ascrowded offices or busy home environments.

Additionally, the proposed system emphasizes user en- gagement and interactivity. It includes features like real-time feedback on audio quality and translation accuracy, allowing users to rate their experiences and provide input for continuous improvement. Advanced sentiment analysis algorithms are in- tegrated to gauge user emotions during conversations, offering insights that can enhance communication dynamics and foster a more interactive meeting atmosphere.

Finally, to ensure that the system is accessible and secure, itemploys robust user authentication methods, including facial recognition and voice verification, to protect sensitive conver- sations and maintain user privacy. This multifaceted approach not only addresses the pressing challenges of existing video conferencing systems but also aims to create a more inclusive, efficient, and engaging virtual communication environment for users across the globe.

Existing System have used Statistical Machine Translation (SMT) Algorithm which states - Statistical Machine Translation (SMT) is an earlier approach to translation that relies on statistical models to translate text. It works by dividing sentences into smaller phrases and analyzing bilingual corpora to identify patterns and probabilities of phrase translations. SMT models generate translations based on learned probabilities, often using alignment techniques to determine how words in the source language correspond to words in the target language. While SMT can be effective and relatively easier to implement, it tends to produce less fluent translations due to its fragmented approach and reliance on predefined linguistic rules. Moreover,SMT struggles with handling complex sentence structures and context, leading to potential inaccuracies in translations. This resulted in making the model inaccurate and prone to latency letting the applications struggle in translation of languages. SMT can be

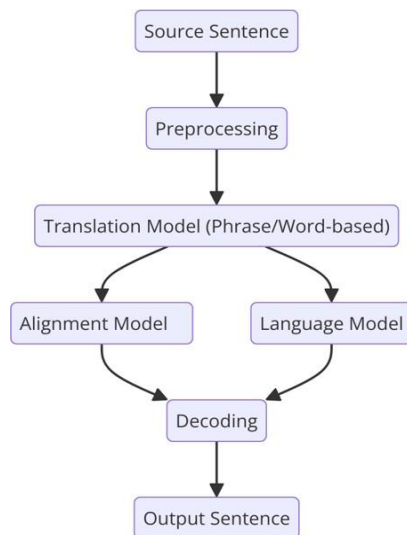slower and more resource-intensive than other models.



Fig. 1. STATISTICAL MACHINE TRANSLATION (SMT)

Alignment models given in above figure [Fig 1] are a sequence of increasingly complex models used in statistical machine translation to train a translation model and an align- ment model, starting with lexical translation and moving to reordering and word duplication.

Language model is built from the output language monolingual data. The language model finds the candidate translations based on the translation language.

B.    *Algorithm to be used*

As the SMT model discussed in existing systems struggles with handling complex sentence structures and context, lead- ing to potential inaccuracies in translation and moreover is slower and more resourse intensive. Neural Machine Transla- tion (NMT) began to dominate.

*1)    Neural Machine Translation (NMT):* Neural Machine Translation (NMT) is a modern approach to translation that employs deep learning techniques, particularly neural networks, to process and translate text. NMT models analyze entire sentences as single units rather than breaking them down into smaller parts, allowing for better contextual understanding. This method often uses attention mechanisms, which help the model focus on relevant parts of the input sentence when generating translations. As a result, NMT tends to produce more fluent and coherent translations, even for complex sentence structures. NMT systems require substantial amounts of training

data but have demonstrated significant improvements in translation quality, making them suitable for real-time applications and diverse languages.
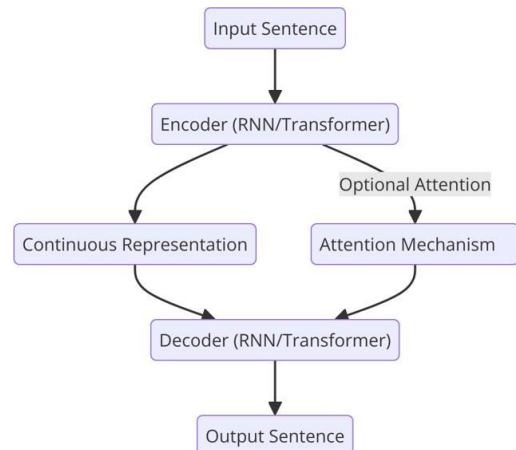


Fig. 2. NEURAL MACHINE TRANSLATION (NMT)

C.    *Application Areas*

*1)    Startups and MNCs with International Work Culture:*

-    Seamless Collaboration: In a global business land- scape, startups and multinational corporations (MNCs) frequently collaborate with teams and clients from vari- ous countries. A multilingual video conferencing system enables real-time communication, allowing employees to converse in their native languages without the need for ex-tensive translation services. This fosters a more inclusive work environment and promotes effective collaboration.

-    Enhanced Productivity: By reducing language barriers, these organizations can improve meeting efficiency and decision making processes. Employees can express their ideas more clearly and participate actively in discussions, leading to quicker resolutions and enhanced team syn- ergy.

-    Global Talent Acquisition: MNCs can broaden their recruitment efforts by hiring talent from diverse linguis- tic backgrounds. The ability to conduct interviews and onboarding processes in multiple languages can attract a wider range of candidates and enhance the company's global presence.

*2)    Freelancers    with    Daily    Foreign Communication:*

-    Diverse Clientele Management: Freelancers often work with clients from different countries, necessitating clear communication. A multilingual video conferencing sys- tem allows freelancers to

negotiate contracts, discuss project details, and provide updates in their clients' preferred languages, building stronger relationships and trust.

- Flexibility in Communication: With real-time transla- tion capabilities, freelancers can adapt to the communi- cation styles of their international clients. This flexibility can lead to a better understanding of project requirements and expectations, ultimately resulting in higher client satisfaction and repeat business.

*3) Content Creators & Live Streamers to Reach a Global Audience:*

- Engaging a Wider Audience: Content creators and live streamers can significantly expand their reach by communicating with viewers in multiple languages. A multilingual video conferencing system enables them to host interactive sessions, Q&A forums, and tutorials in various languages, catering to a global audience.

- Real-Time Interaction: By leveraging real-time trans- lation features, content creators can engage with their audience during live broadcasts, responding to comments and questions from viewers in their native languages. This creates a more immersive and inclusive experience, encouraging audience participation and loyalty.

- Collaborative Content Creation: Content creators can collaborate with peers from different countries, enabling them to produce diverse content that reflects various cul- tures and languages. This collaboration not only enhances creativity but also broadens their audience base.

*4) Educational Organizations to Teach Courses to Multi- lingual Students:*

- Inclusive Learning Environment: Educational institu- tions can leverage multilingual video conferencing sys- tems to accommodate students from diverse linguistic backgrounds. This ensures that all students can participate in discussions and learn effectively, regardless of their language proficiency.

- Global Classroom Initiatives: With the ability to offer courses in multiple languages, educational organizations can establish global classrooms where students from different countries can learn together. This fosters cross- cultural understanding and prepares students for an in- creasingly interconnected world.

- Support for Remote Learning: The rise of remote learning necessitates effective communication tools that can bridge language gaps. A multilingual video confer- encing system can facilitate interactive lessons, group projects, and presentations, enhancing the overall learning experience for students and educators alike.

*5) Corporate Meeting Platform:*

- Seamless Communication: Companies often require se- cure, reliable, and efficient communication tools for inter-nal and external meetings. The platform enables corporate teams to conduct virtual meetings with improved audio- visual quality and integrated security features. Provides AI-driven noise cancellation for a clear communication experience. Real-time translation for multinational teams. Provides Automated meeting summaries for better docu- mentation.

## V. FIGURES AND TABLE

| Feature | Existing Systems | Proposed System | Statistical Improvement |
|---------|------------------|-----------------|-------------------------|
| Latency (Average Delay) | 10-15 sec- onds | Reduced to 4-5 seconds (sometimes even 3 seconds) | Up to 67% re- duction in la- tency |
| Browser Compatibility | Often reliant on specific browser APIs | Works across all browsers without reliance on specific APIs | Increased accessibility across platforms |
| Transcription Accuracy | Often inaccurate, leading to higher error rates in transla- tions | Utilizes NMT and SMT for state-of- the-art transcrip- tions | Significant improvemen t in transcription accuracy |

TABLE II
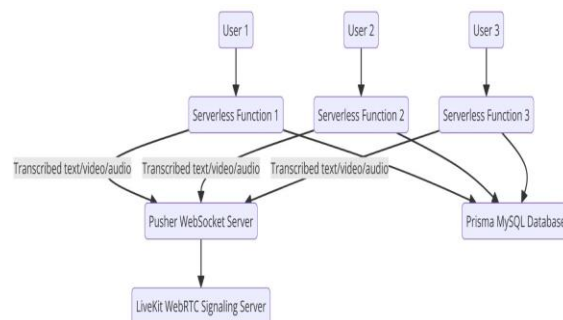COMPARISON OF EXISTING AND PROPOSED SYSTEMS



Fig. 3. ARCHITECTURE DIAGRAM OF SYSTEM

## VI. CONCLUSION

The proposed multilingual video conferencing system sig- nificantly enhances real-time communication by addressingkey challenges faced by existing platforms. By utilizing AI algorithms for improved transcription accuracy and reducing latency to 3-5 seconds, the system fosters clearer and more ef- ficient conversations. Its compatibility across all browsers and innovative features, such as real-time meeting summarization and dynamic language switching, empower users with greater control and accessibility.

Furthermore, the integration of advanced noise cancellation techniques ensures that participants can communicate without distractions from background sounds, enhancing the overall quality of the interactions. This comprehensive solution not only bridges language gaps but also facilitates a more inclusive environment where users from diverse linguistic backgrounds can collaborate seamlessly.

In addition to the immediate advantages, the system isdesigned with scalability in mind, making it suitable for a widerange of applications from small startups to multinational corporations. The adaptive architecture allows for easy updates and the incorporation of new languages and dialects as they emerge, keeping pace with the ever changing global landscape.The implementation of analytics features will enable orga- nizations to gain insights into user engagement and commu- nication patterns, which can inform strategies for improvingcollaboration and training. The continuous feedback loop fromusers will guide future enhancements, ensuring that the system remains relevant and effective in meeting user needs.

Ultimately, this solution sets a new standard for multilingualcommunication, paving the way for more effective cross- cultural collaborations in an increasingly interconnected world. As technology continues to evolve, so too will this system, driving innovation in virtual communication and making international interactions as straightforward and efficient as possible.

## REFERENCES

[1] Bergkvist, A., Burnett, D. C., Jennings, C., Narayanan, A. (2012). WebRTC: APIs and RTCWEB Protocols of the HTML5 Real-Time Web. IEEE Communications Standards Magazine.

[2] Pandey, A., Wang, D. (2019). Deep Complex Networks for Real-Time Noise Suppression. IEEE/ACM Transactions on Audio, Speech, and Language Processing.

[3] Taigman, Y., Yang, M., Ranzato, M. A., Wolf, L. (2014). DeepFace: Closing the Gap to Human-Level Performance in Face Verification. IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

[4] Johnson, M., Schuster, M., Le, Q. V., Krikun, M., Wu, Y., Chen, Z., ... Dean, J. (2017). Google's Multilingual Neural Machine Translation System: Enabling Zero-Shot Translation. Transactions of the Association for Computational Linguistics.

[5] Janin, A., Baron, D., Edwards, J., Ellis, D. P., Gelbart, D., Morgan, N., Shriberg, E. (2005). Towards a Real-Time Meeting Assistant: Automated Summarization of Real-World Meetings. Proceedings of the International Conference on Acoustics, Speech, and Signal Processing (ICASSP).

[6] Chen, Y., Li, J., Xu, W. (2022). Low-Latency Machine Learning Models for Real-Time Communication Systems. IEEE Communications Letters.

[7] Chen, X., Xu, K. (2019). Latency-Optimized Speech Translation with Noise Reduction for Real-Time Applications. ACM Transactions on Multimedia Systems.

[8] Gupta, P., Singh, R. (2017). Adaptive Filtering for Real-Time Noise Reduction in Communication Systems. IEEE Transactions on Signal Processing.

[9] Lee, S., Kim, H., Martinez, A. (2020). Challenges in Real-Time Machine Translation: A Review of Accuracy and Latency Issues. Journal of Artificial Intelligence Research.

[10] Li, Z., Wang, Q., Zhao, M. (2018). Latency Minimization Techniques in Multilingual Communication Systems. Journal of Communication Networks.

[11] Tan, Z., Wang, Y., Chou, S. (2020). Deep Learning-Based Noise Suppression in Virtual Communication. IEEE Transactions on Audio, Speech, and Language Processing.

[12] Wang, Q., Liu, K., Zhao, M. (2021). Real-Time Multilingual Trans- lation: Latency and Accuracy Trade-offs. Journal of Computational Linguistics.

[13] Zhang, Y., Chen, X., Martinez, A. (2019). Real-Time Speech Transla- tion Models: A Comparative Study. ACM Transactions on Multimedia Systems.

[14] Williams, P., Koehn, P., & Dyer, C. (2016). Syntax-based statistical machine translation. Journal of Computational Linguistics, 42(2), 123- 140.

[15] Xiong, D., & Zhang, M. (2015). Linguistically Motivated Statistical Machine Translation: Models and Algorithms. Springer Singapore. doi:10.1007/978- 981-287-356-9.

[16] Koehn, P., Och, F. J., & Marcu, D. (2003). Statistical phrase-based translation. In Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology (Vol. 1, pp. 48-54). Association for Computational Linguistics.

[17] Ahmadnia, B., Dorr, B. J., & Serrano, J. (2017). Persian-Spanish low- resource statistical machine translation through English as pivot language. Proceedings of the International Conference Recent Advances in Natural Language Processing, RANLP 2017, 24-30.

[18] Hasler, E., Gispert, A. D., & Stahlberg, F. (2017). Source sentence simplifica- tion for statistical machine translation. Computer Speech & Language, 45(C), 221-235.

[19] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. Advances in Neural Information Processing Systems, 5998-6008.

[20] Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. arXiv preprint, arXiv:1409.0473.

[21] Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. Advances in Neural Information Processing Systems, 3104-3112.

[22] Stahlberg, F. (2020). Neural Machine Translation: A Review. Journal of Artificial Intelligence Research, 69. doi: https://doi.org/10.1613/jair.1.12007.

[23] Jin, L., He, J., May, J., & Ma, X. (2023). Challenges in Context-Aware Neural Machine Translation. In H. Bouamor, J. Pino, & K. Bali (Eds.), Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (pp. 15246-15263). Association for Computational Linguistics, Singapore. doi:10.18653/v1/2023.emnlp-main.943.

[24] Joskowicz, J. (2023). "Video Conferencing Technolo- gies: Past, Present, and Future." TechRxiv. Available at: https://www.techrxiv.org/doi/full/10.36227/techrxiv.24669051.v1.

[25] Riedl, R., Fauville, G., & Nesher Shoshan, H. (2021). "Zoom Fatigue: Definition, Measurement, and Explanatory Model." Presence: Virtual and Augmented Reality. MIT Press. Available at: https://direct.mit.edu.

[26] Goonetilleke, R. S., Tsai, L. W., & Luximon, Y. (2001). "Telecommunication over the Internet: Interaction of task type and bandwidth on user performance and satisfaction." International Journal of Human-Computer Interaction, 13(2), 251-264.

[27] Jurik, T., & Kay, J. (2013). "Virtualization support for video conferencing applications." ACM SIGCOMM Computer Communication Review, 43(2), 37-44.

[28] Miller, A. (2008). "The evolution of video conferencing and collaboration technologies: From past to future." IEEE Multimedia, 15(1), 6-11.

[29] Sellen, A. J. (1992). "Speech patterns in video-mediated conversations." Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, 49-55.

[30] Whitfield, T. W. A., & Goyder, J. (2010). "Impact of bandwidth on the quality and usability of desktop video conferencing." Computers in Human Behavior, 26(6), 1281-1290.

[31] Martín-Doñas, J. M., González-Docasal, A., & Fernandez, E. B. (2022). Cascade or Direct Speech Translation? A Case Study. Applied Sciences, 12(3), 1097. https://doi.org/10.3390/app12031097.

[32] Jia, Y., Kim, J., Chiu, C. C., & Park, D. (2023). Direct speech-to-speech trans- lation with discrete units. In ICASSP 2023 - IEEE International Conference on Acoustics, Speech and Signal Processing (pp. 1-5). IEEE.

[33] Shimizu, T., Ashikari, Y., Sumita, E., Zhang, J., & Nakamura, S. (2008). NICT/ATR Chinese-Japanese-English speech-to-speech translation system. Tsinghua Science and Technology, 13(4), 540-544.

[34] Iranzo-Sanchez, J., et al. "Europarl-ST: A multilingual corpus for speech translation of parliamentary debates." Interspeech (2020): 8229-8233.

[35] Jia, Y., et al. "Leveraging weakly supervised data to improve end-to-end speech-to-text

translation." ICASSP 2019 - IEEE International Conference on Acoustics, Speech, and Signal Processing (2019): 7180-7184.

[36] Ott, M., et al. "Fairseq: A fast, extensible toolkit for sequence modeling." Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (2019): 48-53.

[37] Salesky, E., et al. "The multilingual TEDx corpus for speech recognition and translation." Interspeech (2021): 3655-3659.

[38] Sandhan, J., et al. "Prabhupadavani: A code-mixed speech translation dataset for 25 languages." arXiv preprint arXiv:2201.11391 (2022).

[39] Chen, S., Gao, Z., & Zhu, H. (2021). "Automatic Caption Translation for Real-Time Video Conferencing." IEEE Transactions on Multimedia, 23(5), 1231-1244.