

A Sign Language Recognition System with Pepper, Lightweight-Transformer, and LLM

JongYoon Lim¹, Inkyu Sa, Bruce MacDonald¹, and Ho Seok Ahn¹

¹ CARES, Department of Electrical, Computer and Software Engineering, University of Auckland
{jy.lim, b.macdonald, hs.ahn}@auckland.ac.nz, enddl22@gmail.com

Abstract

This research explores using lightweight deep neural network architectures to enable the humanoid robot Pepper to understand American Sign Language (ASL) and facilitate non-verbal human-robot interaction. First, we introduce a lightweight and efficient model for ASL understanding optimized for embedded systems, ensuring rapid sign recognition while conserving computational resources. Building upon this, we employ large language models (LLMs) for intelligent robot interactions. Through intricate prompt engineering, we tailor interactions to allow the Pepper Robot to generate natural Co-Speech Gesture responses, laying the foundation for more organic and intuitive humanoid-robot dialogues. Finally, we present an integrated software pipeline, embodying advancements in a socially aware AI interaction model. Leveraging the Pepper Robot’s capabilities, we demonstrate the practicality and effectiveness of our approach in real-world scenarios. The results highlight a profound potential for enhancing human-robot interaction through non-verbal interactions, bridging communication gaps, and making technology more accessible and understandable.

1 Introduction

Each day in the United States, approximately 33 infants are born with irreversible hearing loss [CDC, 2010], with around 90% of these infants born to parents with average hearing ability and potentially lacking proficiency in American Sign Language (ASL) [Mitchell and Karchmer, 2004]. The absence of sign language exposure places these infants in peril of Language Deprivation Syndrome, a condition defined by the absence of accessible, naturally acquired language within their critical language development period [Hall et al., 2017]. This syndrome

has profound implications, affecting various life aspects, including relationships, education, and employment.

To ensure accessible learning of sign language and address the potential challenges of lack of language exposure, various platforms exist to facilitate the learning of sign language. [Süzgün et al., 2015] [Martins et al., 2015]. Notably, multimodal platforms like robots have emerged as highly effective in language instruction, attributed to their interactive and adaptable learning settings [Uluer et al., 2015]. These platforms can meet individual learning necessities and preferences, presenting a multifaceted approach to acquiring language that surpasses conventional instructional methodologies. Integrating Social Human Robots emerges as a pivotal solution [Zakipour et al., 2016]. These robots are envisaged to mitigate the challenges inherent in learning sign language. By utilizing these advanced technologies, it is feasible to construct more inclusive and adjustable learning experiences, allowing a broader spectrum of individuals to communicate proficiently via sign language and consequently reducing the negative impacts of a lack of language acquisition.

However, recognizing sign language and generating human-like gestures in robotic systems is inherently computationally intensive and incredibly challenging for platforms with limited computational resources [Joksimoski et al., 2022]. The demand for real-time data processing, inherent to sign language recognition and natural gesture generation, necessitates high computational throughput and low latency [Sabyrov et al., 2019]. Additionally, deploying sophisticated machine learning algorithms, such as deep neural networks for feature extraction, recognition, and capturing temporal sequences, imposes an additional computational burden.

To address the previously highlighted challenges, we have created a comprehensive system for understanding sign language and making gestures, specifically designed for the Pepper robot. Our main contributions are outlined below:

- Sign Language Recognition: We developed a

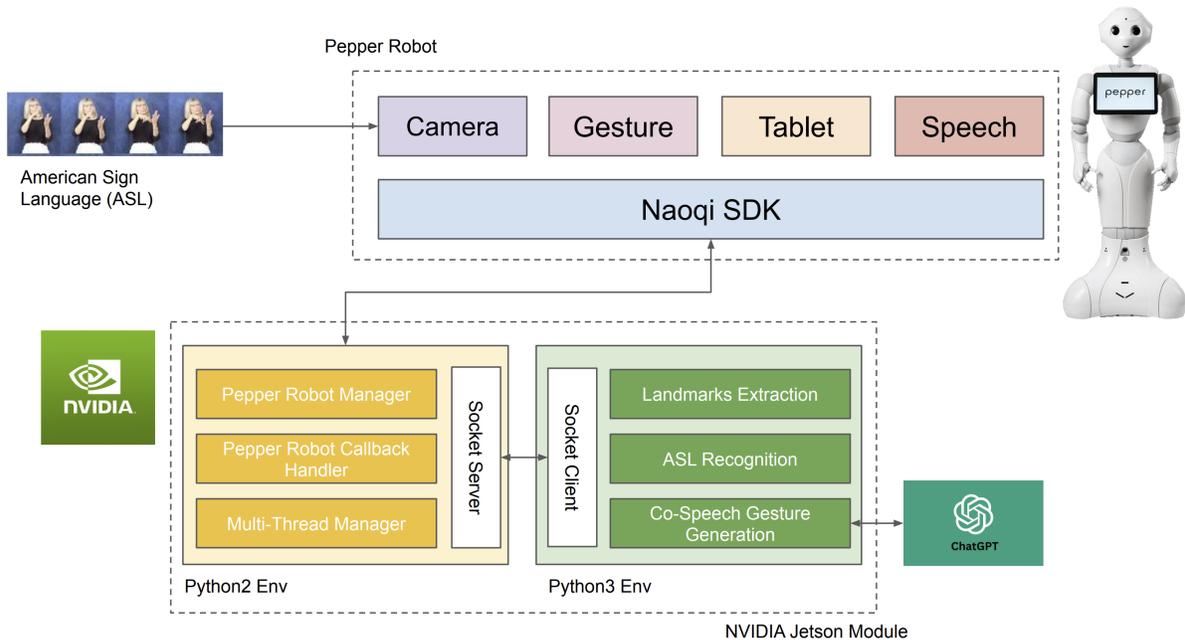


Figure 1: System Overview: frames capturing signs from Pepper are conveyed to the Jetson module, where landmarks are extracted and relayed to the ASL Recognition model. Subsequently, Co-Speech Gesture outputs are derived from ChatGPT and transmitted back to Pepper, enabling the execution of corresponding gestures and dialogue.

lightweight Deep Neural Networks (DNNs) model for understanding American Sign Language, optimized for systems with limited computing power.

- **Smart Interactions:** We employed low-level motions and carefully designed prompts to enable Pepper to interact intelligently, producing appropriate and context-aware gestures using a Large Language Model (LLM) such as ChatGPT.
- **Complete Integration:** We have built a fully integrated approach that combines these elements to enable social interactions between Pepper and humans, paving the way for more advanced human-robot interactions in the future.

2 Related Works

2.1 Sign Language Understanding using Deep Neural Networks

Sign languages, natural languages conveyed through gestures and facial expressions, present unique challenges and opportunities in computer vision and AI. The evolution of DNNs has propelled advancements in the accurate recognition and translation of sign languages [Zuo et al., 2023] [Hu et al., 2021] [Boháček and Hruží, 2022]. Initial efforts in gesture recognition heavily relied on traditional computer vision techniques until the introduction of Convolutional Neural Networks (CNNs) [Rao

et al., 2018], which demonstrated enhanced proficiency in recognizing gestures by focusing on the spatial understanding of signs. Integrating Recurrent Neural Networks (RNNs) [Guo et al., 2018] and transformer networks [Boháček and Hruží, 2022] has proven effective in analyzing the sequential flow of sign gestures to capture the inherent temporal dynamics of sign language. Beyond gesture recognition, the capability of DNNs extends to end-to-end sign language translation, directly converting sign language to text or speech. Recognizing the multimodal nature of sign language [Kagirov et al., 2019], involving not just hand movements but also facial expressions and body posture, multimodal deep learning approaches have been advocated, amalgamating data from diverse sensors to refine recognition accuracy. The development of expansive datasets has been crucial in propelling this research, offering a diverse range of sign languages and signers for robust training and evaluation of DNNs [Li et al., 2020] [Ronchetti et al., 2016] [Albanie et al., 2021]. However, despite these advancements, challenges persist, including data scarcity, signer variability, and the complexities of recording non-manual signs.

2.2 Social Human-Robot Interaction (HRI)

Research in human-robot interactions (HRI) has sparked interest in robotics and social sciences [Lemaignan et al., 2017], evolving from task-oriented interactions to socio-

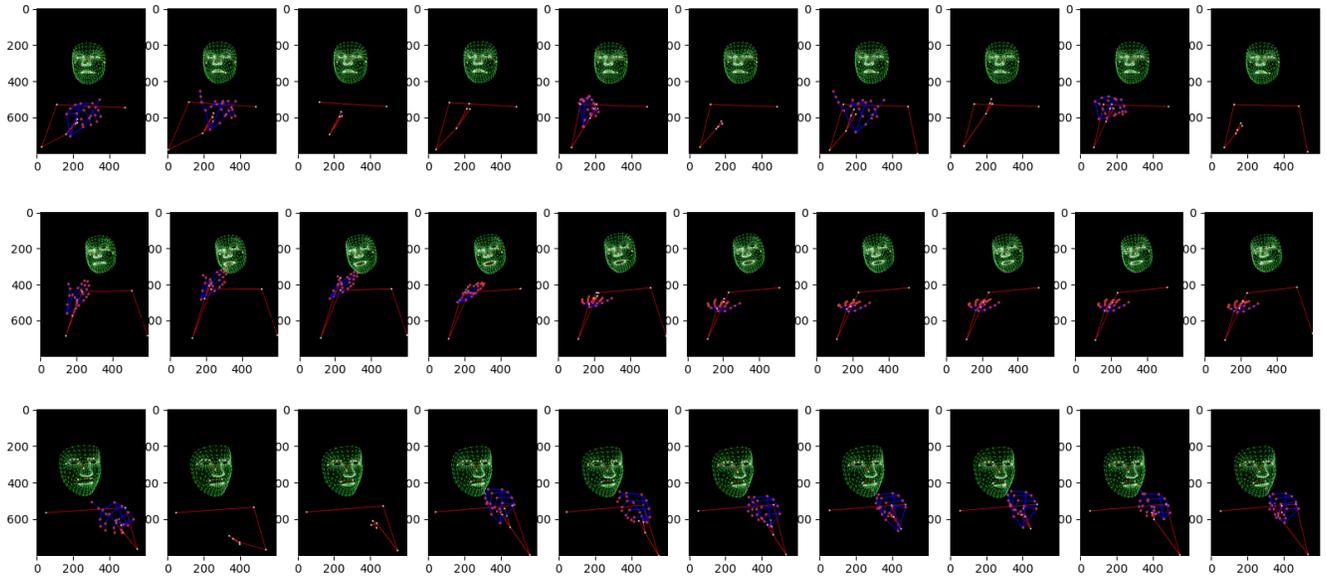


Figure 2: Extraction of Landmarks Using Mediapipe: The top row represents the sign for the word 'same', the middle row depicts the sign for 'bad', and the bottom row illustrates the sign for 'nuts'.

emotional exchanges resembling human-to-human interactions [Johanson et al., 2019]. Advancements in emotion recognition using deep learning enable robots to understand and respond to human emotions, facilitating seamless interactions. The development of robotic empathy has been crucial in fostering genuine human-robot connections, particularly in elderly care and education, where emotional support is vital [Gasteiger et al., 2022]. However, most current HRI focus on verbal communication, overlooking the significance of non-verbal cues like body language and facial expressions in enhancing interaction quality. The ability of robots to undertake perspective-taking improves collaborative work by considering human viewpoints and feelings. Applications of social HRI have yielded impressive results in fields like tutoring and counseling, underscoring the effectiveness of robots possessing socio-emotional skills. Nonetheless, achieving truly social and emotionally resonant HRI poses challenges, with areas such as the uncanny valley effect and the balance between robot autonomy and user control remaining key research domains.

2.3 Large Language Model in Robotics

The fusion of Large Language Models (LLMs) and robotics has sparked extensive research focusing mainly on prompt engineering, aiming to facilitate seamless human-robot interaction [Billing et al., 2023]. Studies in prompt engineering within LLMs have paved the way for enhanced model responses, establishing foundational communication protocols between humans and robots. Researchers have demonstrated that integrating LLMs

in robotic systems allows for the interpretation and execution of complex commands, emphasizing the critical role of optimal prompts [Yu et al., 2023]. Additionally, advancements in multimodal integration enable richer, context-aware interactions by combining visual and linguistic data. However, this integration has brought forth ethical concerns, such as bias and responsible deployment of technologies, necessitating meticulous consideration in their development and application. The practical applications of these integrations are extensive, with notable advancements in healthcare and education. Future research is directed towards refining prompt engineering techniques and developing more coherent interaction paradigms to effectively bridge the gap between natural language understanding and robotic responsiveness.

2.4 Humanoid Robots in Education

Humanoid robots, with their human-like appearance and dynamic interaction abilities, are increasingly being integrated into educational environments, from primary schools to universities, enhancing teaching methods and student engagement. These robots, as explored in studies [Leyzberg et al., 2014] and [Kennedy et al., 2016], serve as effective tutors, providing personalized, consistent, and adaptive learning experiences. They have proven particularly beneficial in language acquisition, offering immersive learning environments for students, especially in learning second languages. Additionally, their utility extends to special education, improving social interaction and focus for children with autism. In STEM education, humanoid robots act as educational

tools for coding and robotics and as agents promoting problem-solving and critical thinking. Integrating humanoid robots necessitates understanding human-robot interaction, with studies [Scheutz, 2011] investigating the social dynamics, trust, rapport, and emotional bonding possibilities between students and robots. However, despite the multitude of benefits, challenges persist in areas like maintenance, teacher training, and balancing human and robot-led instruction, which are crucial to address for maximizing learning outcomes.

2.5 Lightweight Deep Neural Networks in Robotics

Integrating lightweight DNNs with embedded systems like NVIDIA Jetson modules is an important advancement in robotics, drawing significant scholarly interest for its potential to enhance robotic capabilities. Research in this field has extensively focused on designing and optimizing lightweight DNNs to operate efficiently on resource-constrained systems [Ghimire et al., 2022], with studies showcasing the deployment intricacies and advantages of utilizing NVIDIA Jetson modules for improved computational efficiency and power consumption in robotic applications. Significant work has been undertaken to integrate these optimized DNNs with robotic systems, enriching autonomous capabilities and enabling advanced real-time decision-making and environmental perception. Implementing these networks has allowed for real-time object detection and navigation, and advances in multi-sensor fusion have improved the robustness and accuracy of robotic perception modules. However, this domain faces challenges, especially in model optimization and resource allocation, with innovative solutions being proposed to overcome the limitations of embedded systems. Numerous application-specific developments have underscored the versatility and impact of lightweight DNNs in healthcare, agriculture, and industrial automation.

3 Methodology

The proposed system architecture (Figure 1) revolves around enabling Pepper Robot to interpret and interact using ASL. Users initiate communication through sign language, positioning themselves for clear visibility. The robot’s camera sensor captures the user’s gestures and postures, which are processed using Google’s Mediapipe holistic tool to extract human body landmarks. These landmarks are relayed to a DNN model on an NVIDIA Jetson module, which identifies and classifies the signed word or phrase, considering the nuances of hand movements and facial expressions. The identified ASL is inputted into an LLM, like ChatGPT, which generates a corresponding verbal response and suggests appropriate gestures for Pepper Robot. These suggestions are con-

verted into executable instructions using the Naoqi SDK, allowing Pepper Robot to respond to the user with verbal communication and corresponding gestures, offering an interactive and immersive experience.

The Isolated Sign Language Recognition Corpus (version 1.0) is an extensive compilation of approximately 100,000 videos featuring isolated signs. It encompasses hand and facial landmarks, created through Mediapipe version 0.9, and is articulated by 21 Deaf signers who predominantly use American Sign Language, employing a lexicon of 250 signs. The dataset contains columns denoting the frame number in the raw video, the type of landmark (which can be one of ‘face’, ‘left hand’, ‘pose’, ‘right hand’), the landmark index number, and the normalized spatial coordinates of the landmark represented by $[x/y/z]$.

3.1 Sign Language Recognition

Dataset

As shown in Figure 3, only the coordinates of lips, hands, and arm pose are utilized in our approach. The landmarks are normalized using the mean and standard deviation of all landmarks, enhancing the model’s overall performance. To further optimize performance, data augmentation plays a crucial role. Random resampling of the original length and random masking are employed for temporal augmentation. Additionally, spatial augmentation is implemented by applying horizontal flips and random affine transformations, which encompass scaling, shifting, rotating, and shearing.

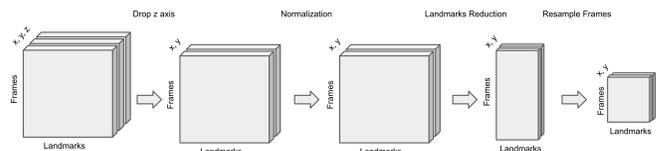


Figure 3: Data Preprocessing. The input data frames undergo a series of transformations: dropping the z-axis, normalization, retaining only the required landmarks, and finally, resampling the frames.

Model

This study employs a specialized model to extract features from landmarks. The initial phase of feature extraction involves the use of multiple dense layers, where each dense layer is succeeded by Layer Normalization and ReLU activation functions. The resulting extracted feature is then forwarded to four layers of a Transformer encoder, integral for processing sequential data and particularly potent for natural language processing tasks due to its self-attention mechanism.

The Transformer encoder processes the input sequence in this model architecture (Figure 4). It compresses the

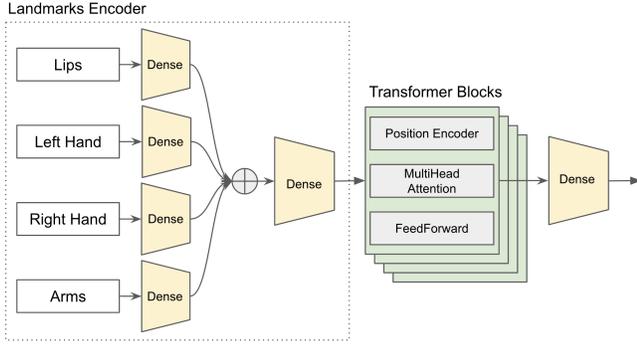


Figure 4: Model Architecture. The preprocessed data is initially passed through a feature extractor and then combined. Subsequently, it is channelled through Transformer blocks before being fed into the classifier.

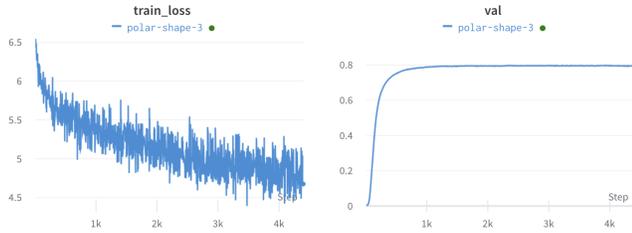


Figure 5: Loss and Validation

information into "context" or "memory," which a decoder usually would use to produce an output sequence in a typical Transformer model. However, the decoder is skipped in this research to ensure parameters and inference time efficiency. Instead, the output from the Transformer encoder layers is directly forwarded to a dense layer to obtain the logits for the classes, avoiding using an activation function in the final dense layer. This approach maintains model efficacy while optimizing computational resources and processing time. The model has a total of 2,562,970 parameters, which is relatively small, yet it still performs reasonably in recognizing ASL.

	Playtime (second)
mean	4.09
std	3.24
min	0.51
25%	2.02
50%	2.90
75%	4.90
max	24.4

Table 1: Play Time Statistics for the recorded 430 gestures executed by Pepper.

3.2 Co-speech gesture Dialogue generation using LLM

Given the capabilities of Large Language Models (LLMs) in understanding and interpreting context within sentences, extracting emotions, sentiments, and other nuanced aspects of language, they serve as powerful tools for enriching interactions with humanoid robots like Pepper. LLMs can be instrumental in generating appropriate and meaningful gestures for Pepper, synchronized with its spoken subtitles, enhancing the overall communicative experience.

Gesture Tag	Thinking
Description	The robot gently taps its head with its right hand, moving carefully and smoothly, like it's deep in thought.
Playtime(s)	2.17
Moving Body Parts	Eyes, Neck, Right Arm, Right Hand

Table 2: A sample of Pepper's gesture descriptor. It includes a gesture tag, a human-authored description, its play time in seconds, and the specific robot body parts required to execute the gesture.

A two-step request (Table 3) to the model can be employed to leverage LLMs for integrating meaningful gestures. Initially, dialogue can be converted to speech using models like ChatGPT. Subsequently, the output from the first step can be prompted to incorporate gesture tags around specific words or sentences, creating a richer, more immersive interaction by aligning gestures with the spoken content. Providing a prompt to the LLM is crucial for generating natural and socially aware outputs. In the prompt instruction, we incorporate gesture descriptors (Table 2) to convey more detailed information about Pepper's predefined gestures. The playtime statistics are displayed in Table 1.

3.3 Deployment Model to NVIDIA Jetson module

To deploy a trained model to the NVIDIA Jetson module, a transformation of the PyTorch model into TensorRT is essential. TensorRT, developed by NVIDIA, stands out as a high-performance deep learning inference library, fine-tuned to enhance the speed and efficiency of deep learning models during the inference phase. It is specifically designed to optimize and accelerate the deployment of models in environments like embedded systems, which is characteristic of the NVIDIA Jetson module. This conversion assures optimal utilization of the board's resources and guarantees swift and efficient model inferences, making it highly suitable for embedded boards where enhanced performance and resource

First Step) Input to LLM
[INSTRUCTIONS] #### A signer accurately depicted a cloud with a 90% accuracy rate.
First Step) Output from LLM
Great! You drew a cloud sign, but the weather today is really nice. Just look up at the sky.
Second Step) Input to LLM
[INSTRUCTIONS—Gesture Descriptors] #### Great! You drew a cloud sign, but the weather today is really nice. Just look up at the sky.
Second Step) Output from LLM
[Yes] Great! [/Yes] You drew a cloud sign, but [Excited] the weather today is really nice [/Excited]. Just [ShowSky] look up at the sky [/ShowSky].

Table 3: Table illustrating the two-step processing approach to generate Co-Speech Gesture using ChatGPT. The initial step utilizes the recognized word and its accuracy to generate a prompt with a specific INSTRUCTION. In the second stage, the returned output is processed using specific INSTRUCTION and Gesture Descriptors. The concluding output is text interspersed with gesture tags.

optimization are crucial.

3.4 Communication between Pepper and Jetson module

The Pepper robot operates using Python 2, while the Jetson module, assigned the task of interpreting sign language, utilizes Python 3. A socket network program establishes effective communication between Pepper and the Jetson module. This approach is founded on network protocols, typically TCP/IP, enabling data exchange between the applications running on Pepper and the Jetson module, which operate as different machines in the network. Socket programming is integral in this setup as it allows for creating scalable and robust network applications, providing a bidirectional communication link between the endpoints. This method is efficient, swift, and integrates seamlessly with other processes, ensuring smooth and responsive interaction between different system components.

4 Results

In the pursuit of bridging communication gaps using AI and robotics, this research has generated noteworthy findings. As we leveraged a confluence of technologies ranging from DNNs to Large Language Models, our empirical observations highlighted the strengths and chal-

lenges inherent in our approach. Our observations underscore a significant step forward in using sign language in human-robot interaction. While certain areas, such as depth prediction for landmarks extraction, require further refinement, the overarching results signify a promising foundation for future enhancements.

4.1 ASL Recognition on Jetson module

Our custom-developed ASL recognition model, optimized for the NVIDIA Jetson module, demonstrated a commendable accuracy rate. Upon testing, the model achieved an accuracy of 79.8% as shown in Figure 5. This is particularly promising, considering the computational constraints of the Jetson module and the complexity inherent in recognizing the nuances of sign language.

4.2 Mediapipe Holistic’s Performance

The Google Mediapipe holistic tool was employed for human body landmarks extraction. Our experiments indicated that the tool strongly predicted landmarks’ x and y positions. However, its capabilities exhibited a limitation when it came to depth prediction. This aspect warrants further investigation and may necessitate supplementary techniques or sensors for robust three-dimensional understanding.

4.3 ChatGPT’s Multimodal Features

One of the more intriguing observations came from deploying the ChatGPT LLM. ChatGPT showcased the ability to generate multimodal features. It was proficient in crafting dialogues while simultaneously generating a diverse array of gestures and emotions. This multifaceted interaction potential reinforces the applicability of LLMs in human-robot interaction scenarios.

4.4 Pipeline Integration

Our integrated pipeline, which amalgamates multiple stages from ASL recognition to robot interaction, functioned seamlessly. The coherence and efficiency of the pipeline validate our architectural choices and implementations. Furthermore, the system is poised for scalability, indicating that it’s ready to incorporate more meaningful experiments geared toward Social Human-Robot Interaction.

5 Discussion

In human-robot interaction, this research addresses the critical need for robots to comprehend and engage meaningfully with humans, especially those relying on ASL. A streamlined, resource-efficient model was developed for real-time ASL recognition, minimizing computational overhead in embedded systems. Incorporating LLMs allows for a deeper understanding of the intent, emotion, and context behind signs, enriching human-robot dialogues. The research’s integrated pipeline epitomizes the

collaboration of various AI technologies, establishing a foundation for socially aware AI interaction models enabling robots to relate to human users empathetically and intuitively. Future works aim at the system’s expansion and refinement, especially in educational sectors, and improved ASL recognition, driving the vision of empathy and understanding robots in human interaction.

Acknowledgments

This work was supported by the Science for Technological Innovation (UOAX2123, Developing a Reo Turi (Māori Deaf Language) Interpreter for Ngāti Turi, the Māori Deaf Community), and the Te Pūnaha Hihiko: Vision Mātauranga Capability Fund (Te Ara Auaha o Muriwhenua Ngāti Turi: The Journey of Muriwhenua Māori Deaf Innovation, UOAX2124) funded by the Ministry of Business, Innovation & Employment (MBIE, New Zealand). Ho Seok Ahn* is the corresponding author.

References

- [Albanie et al., 2021] Albanie, S., Varol, G., Momeni, L., Bull, H., Afouras, T., Chowdhury, H., Fox, N., Woll, B., Cooper, R., McParland, A., et al. (2021). Bbc-oxford british sign language dataset. *arXiv preprint arXiv:2111.03635*.
- [Billing et al., 2023] Billing, E., Rosén, J., and Lamb, M. (2023). Language models for human-robot interaction. In *ACM/IEEE International Conference on Human-Robot Interaction, March 13–16, 2023, Stockholm, Sweden*, pages 905–906. ACM Digital Library.
- [Boháček and Hruz, 2022] Boháček, M. and Hruz, M. (2022). Sign pose-based transformer for word-level sign language recognition. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 182–191.
- [CDC, 2010] CDC (2010). Identifying infants with hearing loss. <https://www.cdc.gov/mmwr/preview/mmwrhtml/mm5908a2.htm> [Accessed: (28 Sep 2023)].
- [Gasteiger et al., 2022] Gasteiger, N., Lim, J., Hellou, M., MacDonald, B. A., and Ahn, H. S. (2022). Moving away from robotic interactions: Evaluation of empathy, emotion and sentiment expressed and detected by computer systems. In *2022 31st IEEE International Conference on Robot and Human Interactive Communication (RO-MAN)*, pages 1365–1370. IEEE.
- [Ghimire et al., 2022] Ghimire, D., Kil, D., and Kim, S.-h. (2022). A survey on efficient convolutional neural networks and hardware acceleration. *Electronics*, 11(6):945.
- [Guo et al., 2018] Guo, D., Zhou, W., Li, H., and Wang, M. (2018). Hierarchical lstm for sign language translation. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
- [Hall et al., 2017] Hall, W. C., Levin, L. L., and Anderson, M. L. (2017). Language deprivation syndrome: a possible neurodevelopmental disorder with sociocultural origins. *Soc. Psychiatry Psychiatr. Epidemiol.*, 52(6):761–776.
- [Hu et al., 2021] Hu, H., Zhao, W., Zhou, W., Wang, Y., and Li, H. (2021). Signbert: pre-training of hand-model-aware representation for sign language recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11087–11096.
- [Johanson et al., 2019] Johanson, D. L., Ahn, H. S., MacDonald, B. A., Ahn, B. K., Lim, J., Hwang, E., Sutherland, C. J., and Broadbent, E. (2019). The effect of robot attentional behaviors on user perceptions and behaviors in a simulated health care interaction: randomized controlled trial. *Journal of medical Internet research*, 21(10):e13667.
- [Joksimoski et al., 2022] Joksimoski, B., Zdravevski, E., Lameski, P., Pires, I. M., Melero, F. J., Martinez, T. P., Garcia, N. M., Mihajlov, M., Chorbev, I., and Trajkovik, V. (2022). Technological solutions for sign language recognition: a scoping review of research trends, challenges, and opportunities. *IEEE Access*, 10:40979–40998.
- [Kagirov et al., 2019] Kagirov, I., Ryumin, D., and Axyonov, A. (2019). Method for multimodal recognition of one-handed sign language gestures through 3d convolution and lstm neural networks. In *International Conference on Speech and Computer*, pages 191–200. Springer.
- [Kennedy et al., 2016] Kennedy, J., Baxter, P., Senft, E., and Belpaeme, T. (2016). Social robot tutoring for child second language learning. In *2016 11th ACM/IEEE international conference on human-robot interaction (HRI)*, pages 231–238. IEEE.
- [Lemaignan et al., 2017] Lemaignan, S., Warnier, M., Sisbot, E. A., Clodic, A., and Alami, R. (2017). Artificial cognition for social human-robot interaction: An implementation. *Artificial Intelligence*, 247:45–69.
- [Leyzberg et al., 2014] Leyzberg, D., Spaulding, S., and Scassellati, B. (2014). Personalizing robot tutors to individuals’ learning differences. In *Proceedings of the 2014 ACM/IEEE international conference on Human-robot interaction*, pages 423–430.
- [Li et al., 2020] Li, D., Rodriguez, C., Yu, X., and Li, H. (2020). Word-level deep sign language recognition from video: A new large-scale dataset and methods

- comparison. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 1459–1469.
- [Martins et al., 2015] Martins, P., Rodrigues, H., Rocha, T., Francisco, M., and Morgado, L. (2015). Accessible options for deaf people in e-learning platforms: technology solutions for sign language translation. *Procedia Computer Science*, 67:263–272.
- [Mitchell and Karchmer, 2004] Mitchell, R. E. and Karchmer, M. (2004). Chasing the mythical ten percent: Parental hearing status of deaf and hard of hearing students in the united states. *Sign language studies*, 4(2):138–163.
- [Rao et al., 2018] Rao, G. A., Syamala, K., Kishore, P., and Sastry, A. (2018). Deep convolutional neural networks for sign language recognition. In *2018 conference on signal processing and communication engineering systems (SPACES)*, pages 194–197. IEEE.
- [Ronchetti et al., 2016] Ronchetti, F., Quiroga, F., Estrebou, C. A., Lanzarini, L. C., and Rosete, A. (2016). Lsa64: An argentinian sign language dataset. In *XXII Congreso Argentino de Ciencias de la Computación (CACIC 2016)*.
- [Sabyrov et al., 2019] Sabyrov, A., Mukushev, M., and Kimmelman, V. (2019). Towards real-time sign language interpreting robot: Evaluation of non-manual components on recognition accuracy. In *CVPR Workshops*.
- [Scheutz, 2011] Scheutz, M. (2011). 13 the inherent dangers of unidirectional emotional bonds between humans and social robots. *Robot ethics: The ethical and social implications of robotics*, page 205.
- [Süzgün et al., 2015] Süzgün, M., Özdemir, H., Camgöz, N., KINDIROĞLU, A., Başaran, D., Togay, C., and Akarun, L. (2015). Hospisign: an interactive sign language platform for hearing impaired. *Journal of Naval Sciences and Engineering*, 11(3):75–92.
- [Uluer et al., 2015] Uluer, P., Akalın, N., and Köse, H. (2015). A new robotic platform for sign language tutoring: Humanoid robots as assistive game companions for teaching sign language. *International Journal of Social Robotics*, 7:571–585.
- [Yu et al., 2023] Yu, W., Gileadi, N., Fu, C., Kirmani, S., Lee, K.-H., Arenas, M. G., Chiang, H.-T. L., Erez, T., Hasenclever, L., Humplik, J., Ichter, B., Xiao, T., Xu, P., Zeng, A., Zhang, T., Heess, N., Sadigh, D., Tan, J., Tassa, Y., and Xia, F. (2023). Language to rewards for robotic skill synthesis.
- [Zakipour et al., 2016] Zakipour, M., Meghdari, A., and Alemi, M. (2016). Rasa: A low-cost upper-torso social robot acting as a sign language teaching assistant. In *Social Robotics: 8th International Conference, ICSR 2016, Kansas City, MO, USA, November 1-3, 2016 Proceedings 8*, pages 630–639. Springer.
- [Zuo et al., 2023] Zuo, R., Wei, F., and Mak, B. (2023). Natural language-assisted sign language recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14890–14900.