

Near-Perfect Alpha-Numeric ASL Recognition Using Salient Object Detection and 2D CNNs

Forrest Moulin

B.Sc. in Information Sciences and Technology

Independent Researcher

contact@mainstreamstudios.ai

Abstract

This research presents SignWave, an innovative alpha-numeric American Sign Language (ASL) recognition system developed by the author, which uniquely combines computer vision, machine learning, and salient object detection techniques. Leveraging background noise removal preprocessing and a deep convolutional neural network (CNN), SignWave recorded an outstanding 99.97% validation accuracy on a compact dataset of 2.25 GB, representing 36 ASL gesture classes A-Z and 1-10. Trained on a cloud virtual machine without a graphics processing unit (GPU), The model's efficiency is highlighted by its use of only 20 epochs, and the results demonstrate that high-accuracy sign language recognition (SLR) can be achieved with accessible computing resources. Designed for real-time sign language interpretation software, SignWave aims to enhance communication between sign language users and non-signers. By showcasing noticeable accuracy, efficiency, and a broader range of classes than some contemporary methods, this research has the potential to advance communication accessibility.

Introduction

American Sign Language is a visual-spatial language with unique linguistic rules, used commonly in deaf communities of North America [2, 3], making it an ideal language for computer vision interpretation use cases. In the United States, approximately 11 million individuals were reported to be deaf or have serious difficulty hearing in 2021 [15]. The language has seen significant growth in popularity, with ASL courses experiencing a 6,583% enrollment increase from 1990 to 2016, making it the third most-studied language on U.S. college campuses [3].

Sign language characters are communicated via fingerspelling, which involves manual representation of letters or numbers using specific handshapes. Fingerspelling serves as a critical component of ASL communication, allowing signers to spell out proper nouns, acronyms, and words that lack a corresponding sign [9]. In the context of SLR systems like SignWave, fingerspelling classification is fundamental, as it provides the initial framework for identifying simpler signs, paving the way for more advanced word gesture recognition and interpretation capabilities.

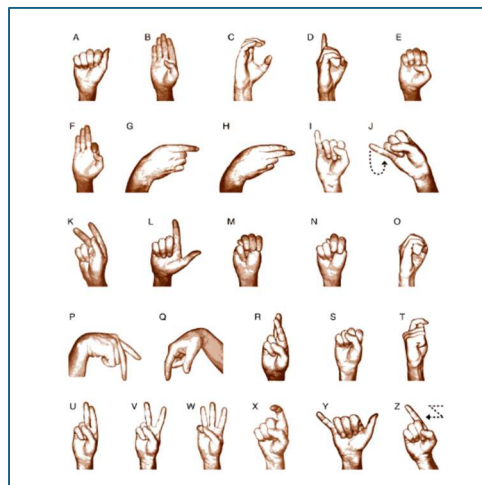


Figure 1: ASL Alphabet Fingerspelling [18]

SLR systems have potential to play a crucial role in bridging communication gaps between deaf individuals and the hearing community worldwide. According to the National Geographic Society [16], over 300 distinct sign languages are used worldwide, serving a community of more than 72 million deaf or hard-of-hearing individuals. These systems can enhance accessibility, improve educational outcomes, and foster social inclusion.

One challenge with existing SLR systems is that datasets may have diverse backgrounds that introduce noise and prevent effective training results. Furthermore, traditional methods often require larger datasets as well as extensive preprocessing and data augmentation, making them resource intensive. For example, preprocessing of the CNN model published by Kumar et al. [8] required that 49,000 images be removed from the initial dataset.

SignWave aims to address this gap by utilizing a single preprocessing step—background noise removal—resulting in a more lightweight model that operates efficiently with smaller, non-augmented datasets. Unlike traditional SLR methods that rely on extensive datasets, which can exceed 100,000 images and require 100 epochs, the proposed approach achieved significant results with only 36,000 images and 20 epochs [4].

Potential SLR applications include sign language proficiency assessments as well as vision-based translation and interpretation through Sign-to-Text (STT) or Sign-to-Speech (STS) systems for transcription or vocalization. The vision-based approach is notably more accessible compared to sensor-based glove systems, which can be impractical due to the inconvenience of wearing and removing gloves lack of portability, and the expenses associated with custom hardware for various hand sizes.

Alternatively, widespread availability of smartphones and computers with front-facing cameras and sufficient random-access memory (RAM) specifications enhances the feasibility of vision-based sign language solutions. According to Pew Research Center [17], approximately 9 out of 10 Americans own a smartphone. Due to the high level of accessibility of smartphone and computer cameras, vision-based systems are a promising avenue for improving communication and inclusivity for deaf or hard-of-hearing individuals.

Related Work

Recent advances in machine learning and computer vision have significantly enhanced the accuracy and sophistication of SLR systems, particularly for alpha-numeric recognition. Contemporary SLR systems often leverage state-of-the-art open-source tools such as MediaPipe for hand landmark detection or gesture recognition. MediaPipe,

an open-source framework developed by Google, facilitates the construction and deployment of machine learning pipelines and is widely used in SLR research [12].

In the early stages of SLR research, models struggled with lower accuracy rates due to limited computational power and less sophisticated algorithms. For example, early systems in the 1990s and early 2000s achieved recognition rates around 80-90% using techniques like Hidden Markov Models (HMMs) and basic feature extraction methods [20, 8]. These systems were often limited by the small datasets and the simplistic nature of the models employed.

Sundar and Bagyammal [21] achieved 99% accuracy on 26 ASL classes by combining MediaPipe with Long Short-Term Memory (LSTM) networks. However, their dataset was relatively shallow, comprising an equivalent total of 3,380 images. Kumar et al. [8] reported an impressive 99.95% accuracy on the same number of ASL classes using MediaPipe with a CNN. Similarly, Barbhuiya et al. [1] recorded a 99.82% accuracy on 36 ASL classes using a combination of the pre-trained AlexNet model and Support Vector Machines (SVM).

Method	Classes	Parameters	Accuracy (%)
MediaPipe & CNN[8]	26	Not specified	99.95
Simple CNN [4]	29	2,029,150	99.89
AlexNet & SVM[1]	36	Not specified	99.82
MediaPipe & LSTM[21]	26	188,090	99.00

Table 1: Related SLR Model Comparison

Zhou et al. [23] explored a sensor-based approach, achieving 98.63% accuracy on 660 sign gestures. Despite its high accuracy, the method's limitations in terms of portability, hardware costs, and the inconvenience of using specialized gloves pose significant challenges. Additionally, Li et al. [11] proposed a word-level deep sign language recognition system from video, demonstrating the potential of large-scale datasets and comparing different methods for enhanced performance in SLR tasks.

Methodology

The dataset used for training and validation was sourced from Synthetic ASL Alphabet and Synthetic ASL Numbers [5, 6]. Initially, the sourced dataset required 9.97 GB of storage and consisted of 37,000 images with a resolution of 512 x 512 pixels, including 'Blank' classes with random backgrounds. After the 'Blank' classes were removed, the dataset was organized into 36 classes representing ASL letters A-Z and numbers 1-10. Each class contained 900 images for training and 100 images for validation.

Preprocessing involved background removal using salient object detection via the rembg Python library to crop the hand from the image foreground and remove the background. This approach significantly reduced dataset size and complexity, resulting in a streamlined dataset of 2.25 GB — over a 75% reduction from the original size. This preprocessing step aimed to enhance model performance by eliminating background noise and focusing the CNN on relevant gesture features.

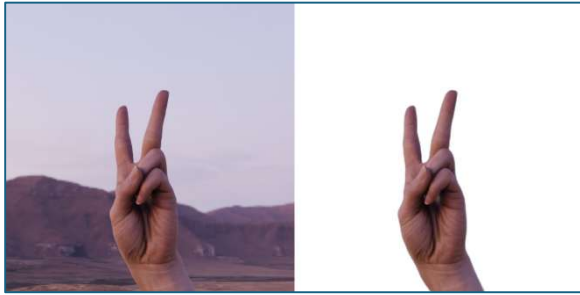


Figure 2: ASL Letter V Background Removal with Salient Object Detection

SignWave was developed using the Python 3.11.2 programming language and implemented with TensorFlow 2.16.1, a versatile deep learning framework providing robust support for designing and training neural networks. This study used a Google Cloud E2 virtual machine configured with 16 GB RAM, 4 virtual CPUs, and 20 GB mounted storage, which contrasts with more resource intensive environments commonly used in similar studies

The development process also incorporated several open-source libraries to enhance functionality and streamline tasks. Matplotlib 3.9.0 was utilized for line graph visualizations, providing clear and informative visual representation of data. NumPy 1.26.2 facilitated numerical operations, ensuring

efficient computation and manipulation of large datasets. Pandas 2.1.4 was employed for organizing training history into DataFrames, making data handling more manageable. Rembg 2.0.57 was used for background removal, a crucial preprocessing step in this research's image data preparation. Scikit-learn 1.5.1 enabled the computation of confusion matrices, essential for evaluating model performance, and Seaborn 0.13.2 was leveraged for visualizing these confusion matrices, enhancing the interpretability of classification results.

SignWave includes three distinct variants and architectures:

Model	Classes	CNN Layers	Trainable Parameters
ABC	26 (A-Z)	10	1,014,394
123	10 (1-10)	10	1,012,330
ABC-123	36 (A-Z & 1-10)	15	2,086,116

Table 2: SignWave Model Comparison

The ABC and 123 models were designed using a simpler 10-layer CNN architecture most comparable to the LeNet-5 structure. These models incorporate two convolutional layers followed by pooling layers, and two fully connected layers. ABC-123 utilized a 15-layer CNN architecture similar to a simplified version of VGG16, with multiple convolutional layer pairs followed by pooling layers for handling the larger number of classes [19, 22]. All three models were trained for 20 epochs with a batch size of 25 images, which were rescaled to 50 x 50 pixels to reduce the total number of parameters and accelerate the training process.

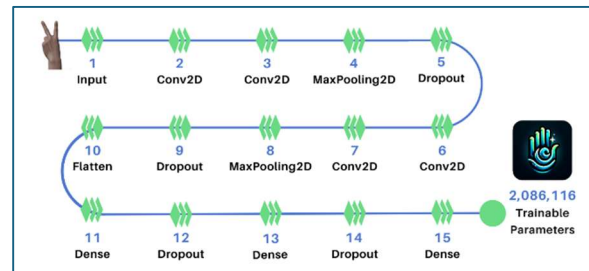


Figure 3: SignWave ABC-123 Deep CNN Architecture

A custom Keras Callback class was implemented to monitor model performance. This callback generated and updated accuracy and loss curve plots for both training and validation datasets using Matplotlib. The plots, along with relevant tables, were compiled into a single PDF report updated at the end of each epoch. This approach provided real-time insights into the model's progress and facilitated the analysis of training dynamics. After each epoch, the custom callback saved the model with a unique file name that included the model type, version, epoch number, and timestamp. This naming convention ensured efficient management of model iterations, making it easier to retrain or utilize the Keras files in an organized manner.

To provide a comprehensive view of model performance, several visualizations were generated, including accuracy curves, cross-entropy loss curves, confusion matrices, and data tables. Model performance was analyzed by comparing training and validation metrics to ensure consistent accuracy. This evaluation process involved cross-referencing validation accuracy with training results to validate model reliability and guide necessary adjustments. This methodology ensured continuous assessment of performance metrics, enabling iterative improvements and robust model development.

Performance evaluation included plotting a confusion matrix to analyze prediction performance by class. The results were visualized using Matplotlib and Seaborn, providing a comprehensive view of the model's effectiveness in recognizing and classifying ASL gestures.

Results

The SignWave models achieved exceptional performance metrics, indicating their effectiveness and efficiency in ASL recognition tasks. Notably, the models recorded near-perfect validation accuracy as well as impressive cross-entropy loss values. The validation results were as follows:

Model	Validation Accuracy (%)	Validation Loss
ABC	99.96	0.0023
123	100.0	0.0002
ABC-123	99.97	0.0009

Table 3: SignWave Model Validation Metrics

The training and validation accuracy curves for SignWave ABC-123 depicted in Figure 4 illustrate the model's effective learning process. The plateauing difference between the training and validation curves suggests slight underfitting, which can lead to improved generalization of unseen data [7].

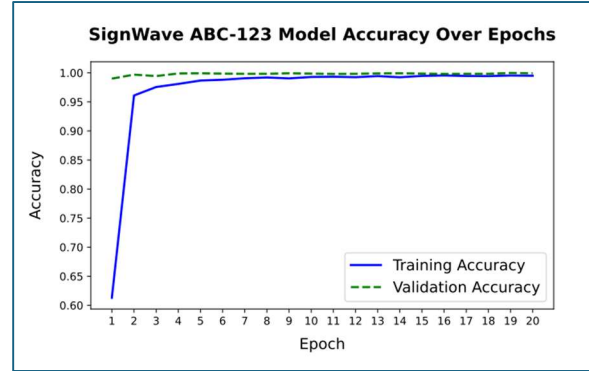


Figure 4: SignWave ABC-123 Training and Validation Accuracy Curves

The confusion matrix for SignWave ABC-123 presented in Figure 5 details the model's classification performance across the 36 ASL gesture classes. The true prediction values along the diagonal indicate the model perfectly classified 35/36 classes of the ABC-123 model. The lone misclassification occurred with the letter R, which was mistaken for the letter C only once.

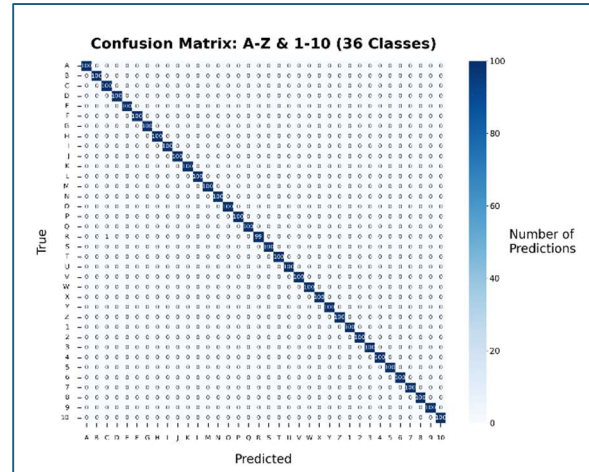


Figure 5: SignWave ABC-123 Confusion Matrix

Class	Precision	Recall	F1-Score
R	0.99	1.0	0.99
C	1.0	0.99	0.99
All Others	1.0	1.0	1.0

Table 4: ABC-123 Model
Class Performance Metrics

Precision, recall, and F1-scores were calculated for each class, shown in Table 4. Precision is the ratio of correctly predicted positive observations to the total predicted positives, while recall is the ratio of correctly predicted positive observations to all observations in the actual class. F1-score is the weighted average of precision and recall.

Discussion

The SignWave models were benchmarked against known ASL recognition methods to assess their relative performance. Traditional methods typically rely on extensive datasets, prolonged training periods, and advanced computational resources, such as GPUs.

For instance, SignWave utilized a streamlined dataset of 2.25 GB for the ABC-123 model, which contained 60% less images than the dataset used by a comparable model that only classified 29 ASL classes [4]. This reduction in dataset size without compromising accuracy underscores the effectiveness of the salient object detection technique to remove background noise.

Unlike traditional models that often require more costly CPUs and GPUs for training, SignWave models were trained on a cloud virtual machine without a GPU. For example, research by Elsayed et al. [4] utilized an Intel Core i9 CPU and NVIDIA GeForce RTX 2080 Ti GPU, suggesting a significantly more complex and costly computing environment. By utilizing a more accessible configuration, our approach demonstrates that high-performance results can be achieved without the necessity for such expensive hardware. This approach highlights the potential for data scientists and machine learning specialists to conduct advanced research and development with more affordable computing resources.

Traditional ASL recognition models often involve 50 or more epochs for training [4, 8, 21]. Conversely, SignWave's achieved superior validation

accuracy in less than 20 epochs, demonstrating a substantial improvement in training efficiency when background noise is removed. In fact, the ABC model reached 99.96% accuracy in just 5 epochs, and the 123 model reached 100% accuracy in just 13 epochs.

However, a limitation of using salient object detection as a preprocessing step, is that this same preprocessing is required before the image can be passed to the model for prediction in the context of a live ASL interpretation or video-based ASL translation. While this technique is near perfect with static images, performance with video feed frames may vary. Thus, at least 16 GB RAM is recommended for deployment in a live recognition system. Additionally, background removal does reduce the dataset storage size, but the process can be time consuming, potentially taking longer than the training process itself. Lastly, while the data appears to be representational of real-world images of hands, it was synthetically generated.

With a sufficient RAM configuration, the ABC-123 model can be deployed in conjunction with a real-time Python app that offers Sign-to-Text to transcribe the processed gesture, or Sign-to-Speech to read it aloud. In real-time SLR, there are limitations to the speed at which signs can be predicted and displayed due to the sequential nature of the process. Initially, the system must detect that a hand movement or sign change has occurred, which involves processing the video feed to identify and track hand gestures.

Once a sign is detected, the system then classifies the gesture, analyzing the movement and matching it to known sign patterns using machine learning models. After recognizing the gesture, the system generates and displays the prediction, converting the recognized sign into text or another form of output and updating the display [13]. Each of these stages introduces some delay, contributing to the overall latency of the system. This latency impacts how quickly the system can respond to and display the recognized signs, affecting the flow of communication in real-time SLR contexts.

The ABC-123 model can also predict 36 classes from appropriately cropped static images when backgrounds have already been removed. Because images were rescaled at training, dataset images of larger sizes can be rescaled down to 50 x 50 pixels before being passed to the model for class prediction.

Conclusion

This research has confirmed the efficacy of the SignWave system for recognizing alpha-numeric American Sign Language (ASL) gestures using salient object detection and 2D convolutional neural networks (CNNs). The successful application of these techniques in a static image processing setting underscores the potential of resource-efficient solutions for sign language recognition (SLR).

SignWave's high accuracy, achieved without the need for a GPU, underscores the potential of leveraging public cloud infrastructure for developing advanced SLR systems. This research demonstrates that high-performance models can be effectively developed and deployed in the cloud, highlighting the logistical and cost benefits associated with cloud computing for SLR. The successful application of resource-efficient techniques, such as innovative preprocessing methods and streamlined model architectures, showcases how public cloud environments can provide a practical and economical solution for creating high-accuracy SLR applications. These findings suggest that cloud-based approaches can significantly reduce the need to purchase expensive hardware while still achieving competitive performance in machine learning tasks.

Currently, most SLR systems are capable of handling basic Sign-to-Text for individual letters or numbers, which is analogous to vocalizing each letter or number sequentially. However, a more advanced, word-based solution would provide a more conversational experience, similar to Speech-to-Text technologies. This would enable a more natural and fluid communication style, similar to actual conversation rather than discrete sign recognition.

Future research can build upon these findings to explore more complex scenarios, such as live-feed gesture recognition and the development of word-based systems. The next steps include adapting the system for real-time sign language interpretation and Sign-to-Speech functionality, as well as expanding its capabilities to support a broader range of gestures and sign languages. This research paves the way for continued advancements in SLR, aiming to enhance communication accessibility and inclusivity for sign language users worldwide.

References

- [1] A. A. Barbhuiya, R. K. Karsh, and R. Jain, "CNN-based feature extraction and classification for sign language," *Multimedia Tools and Applications*, vol. 80, no. 2, pp. 3051–3069, 2021. Available: <https://doi.org/10.1007/s11042-020-09095-1>.
- [2] Canadian Association of the Deaf, "Language," 2022. [Online]. Available: <https://cad-asc.ca/issues-positions/language/>. [Accessed: Jul. 22, 2024].
- [3] Clemson University, "American Sign Language," 2024. [Online]. Available: <https://www.clemson.edu/cah/academics/languages/languages/asl.html>. [Accessed: Jul. 22, 2024].
- [4] N. Elsayed, A. Ibrahim, and M. Saleh, "Vision-based American Sign Language classification approach via deep learning," *arXiv*, 2022. [Online]. Available: <https://arxiv.org/pdf/2204.04235>. [Accessed: Jul. 22, 2024].
- [5] O. Fahey, "Synthetic ASL Alphabet [Data set]," Kaggle, Jun. 17, 2022. [Online]. Available: <https://www.kaggle.com/datasets/lexset/synthetic-asl-alphabet>. [Accessed: Jul. 22, 2024].
- [6] O. Fahey, "Synthetic ASL Numbers [Data set]," Kaggle, Jun. 17, 2022. [Online]. Available: <https://www.kaggle.com/datasets/lexset/synthetic-asl-numbers>. [Accessed: Jul. 22, 2024].
- [7] I. Goodfellow, Y. Bengio, and A. Courville, "Deep learning," MIT Press, 2016. Available: <https://www.deeplearningbook.org/>.
- [8] R. Kumar, A. Pandey, and A. Gupta, "Mediapipe and CNNs for real-time ASL gesture recognition," *arXiv*, 2023. [Online]. Available: <https://arxiv.org/pdf/2305.05296>. [Accessed: Jul. 22, 2024].
- [9] B. Lee and K. Secora, "Fingerspelling and its role in translanguaging," *Languages*, vol. 7, no. 4, p. 278, 2022. Available: <https://doi.org/10.3390/languages7040278>.
- [10] R. Liang and M. Ouhyoung, "A real-time continuous gesture recognition system for sign language," in *Proceedings of the Third International Conference on Automatic Face and Gesture Recognition*, 1998, pp. 558–565. Available: <https://doi.org/10.1109/AFGR.1998.670869>.

- [11] D. Li, C. Rodriguez, X. Yu, and H. Li, "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*, 2020, pp. 1459-1469. Available: <https://doi.org/10.1109/WACV45572.2020.9072872>.
- [12] C. Lugaresi, J. Tang, H. Nash, C. McClanahan, E. Uboweja, M. Hays, F. Zhang, C.-L. Chang, M. Yong, J. Lee, W.-T. Chang, W. Hua, M. Georg, and M. Grundmann, "MediaPipe: A framework for perceiving and processing reality," in *Proceedings of the Third Workshop on Computer Vision for AR/VR at IEEE Computer Vision and Pattern Recognition (CVPR)*, 2019. [Online]. Available: <https://research.google/pubs/mediapipe-a-framework-for-perceiving-and-processing-reality/>. [Accessed: Jul. 22, 2024].
- [13] F. Moulin, "python-asl-detection," GitHub, Jun. 2024. [Online]. Available: <https://github.com/ffm5113/python-asl-detection>. [Accessed: Jul. 22, 2024].
- [14] National Association of the Deaf, "Learning American Sign Language," 2024. [Online]. Available: <https://www.nad.org/resources/american-sign-language/learning-american-sign-language/>. [Accessed: Jul. 22, 2024].
- [15] National Deaf Center on Postsecondary Outcomes, "How many deaf people live in the United States?" 2024. [Online]. Available: <https://nationaldeafcenter.org/faq/how-many-deaf-people-live-in-the-united-states/>. [Accessed: Jul. 22, 2024].
- [16] National Geographic Society, "Sign language," 2024. [Online]. Available: <https://education.nationalgeographic.org/resource/sign-language/>. [Accessed: Jul. 22, 2024].
- [17] Pew Research Center, "Mobile fact sheet," 2024. [Online]. Available: <https://www.pewresearch.org/internet/fact-sheet/mobile/>. [Accessed: Jul. 22, 2024].
- [18] Pocket Sign, "Sign language alphabet – ASL fingerspelling," n.d. [Online]. Available: <https://www.pocketsign.org/alphabet>. [Accessed: Jul. 22, 2024].
- [19] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *International Journal of Computer Vision*, vol. 120, no. 3, pp. 142-158, 2015. Available: <https://doi.org/10.1007/s11263-015-0816-y>.
- [20] T. Starner, J. Weaver, and A. Pentland, "Real-time American Sign Language recognition using desk and wearable computer-based video," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, no. 12, pp. 1371-1375, 1997. Available: <https://doi.org/10.1109/34.650113>.
- [21] B. Sundar and T. Bagyammal, "American Sign Language recognition for alphabets using MediaPipe and LSTM," *Procedia Computer Science*, vol. 215, pp. 642-651, 2022. Available: <https://doi.org/10.1016/j.procs.2022.12.066>.
- [22] Z. Tao, Z. Yang, B. Chen, W. Bao, and H. Cheng, "Protein sequence classification with LetNet-5 and VGG16," in *Intelligent Computing Theories and Applications*, D. S. Huang, K. H. Jo, J. Jing, P. Premaratne, V. Bevilacqua, and A. Hussain, Eds. Springer, 2022, vol. 13394, pp. 621-633. Available: https://doi.org/10.1007/978-3-031-13829-4_60.
- [23] Z. Zhou, K. Chen, X. Li, S. Zhang, Y. Wu, Y. Zhou, and J. Chen, "Sign-to-speech translation using machine-learning-assisted stretchable sensor arrays," *Nature Electronics*, vol. 3, pp. 571-578, 2020.