

Article

Not peer-reviewed version

---

# Sign Language Recognition with Multimodal Sensors and Deep Learning Methods

---

[Chenghong Lu](#) , Misaki Kozakai , [Lei Jing](#) \*

Posted Date: 21 September 2023

doi: 10.20944/preprints202309.1462.v1

Keywords: Sign Language Recognition; sensor fusion



Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

*Article*

# Sign Language Recognition with Multimodal Sensors and Deep Learning Methods

Chenghong Lu <sup>1</sup>, Misaki Kozakai <sup>2</sup> and Lei Jing <sup>\*</sup>

Graduate School of Computer Science and Engineering, University of Aizu, Tsuruga, Ikki-machi, Aizuwakamatsu City 965-8580, Japan

\* Correspondence: leijing@u-aizu.ac.jp

**Abstract:** Sign language recognition is essential in hearing-impaired people's communication. Sign language recognition is an important concern in computer vision and has been developed with rapid progress in image recognition technology. However, sign language recognition using a general monocular camera has problems with occlusion and recognition accuracy in sign language recognition. In this research, we aim to improve accuracy by using a 2-axis bending sensor as an aid in addition to image recognition. We aim to achieve higher recognition accuracy by acquiring hand keypoint information of sign language actions captured by a monocular RGB camera and adding sensor assist. To improve sign language recognition, we need to propose new AI models. In addition, the amount of dataset is small because it uses the original data set of our laboratory. To learn using sensor data and image data, we used MediaPipe, CNN, and BiLSTM to perform sign language recognition. MediaPipe is a method for estimating the skeleton of the hand and face provided by Google. In addition, CNN is a method that can learn spatial information, and BiLSTM can learn time series data. Combining the CNN and BiLSTM methods yields higher recognition accuracy. We will use these techniques to learn hand skeletal information and sensor data. Additionally, the 2-axis Bending sensor glove data support training AI model. Using these methods, we aim to improve the recognition accuracy of sign language recognition by combining sensor data and hand skeleton data. Our method performed better than using skeletal information, achieving 96.5% accuracy in Top-1.

**Keywords:** sign language recognition; sensor fusion

## 1. Introduction

Recognition of hand motion capture is an interesting topic. Hand motion can represent many gestures. In particular, sign language plays an important role in the daily lives of hearing-impaired people. About 2.5 billion people are expected to have some degree of hearing loss by 2050, according to the WHO, and more than 1 billion young people are at risk of permanent hearing loss [1]. In addition, due to the impact of infectious diseases in recent years, online communication has become important. Facilitating communication between sign language users and non-users via video calls remains a pertinent research focus. However, the intricate nature of sign language gestures presents challenges to achieving optimal recognition solely through wearable data gloves or camera-based systems.

Both wearable data gloves and camera-based systems have been extensively explored for sign language recognition. Bending sensors glove only focus on finger bending degree. Consequently, several sign language words exhibiting similar curvature patterns become indistinguishable. This limitation curtails the utility of such devices. Given the significance of hand and arm gestures in sign language, it is imperative for vision-based approaches to prioritize the extraction of key points data from the hands, thereby reducing interference from extraneous background elements. Occlusion presents a significant challenge to vision-based methodologies. During the acquisition of hand key points, monocular cameras may fail to capture certain spatial information due to inter-finger occlusions. Such occlusions often act as impediments, constraining the potential for enhancement in recognition accuracy. In gesture recognition, it is easy for fingers to block each other, objects to block hands, or even parts to be nearly blocked due to overexposure or too dark, resulting in unrecognizability. As

shown in Figure 1, the occlusion problem is less effective in obtaining key points. Integration with bending sensors offers a solution, enabling precise measurement of finger angles, even in regions overlapped by external entities.

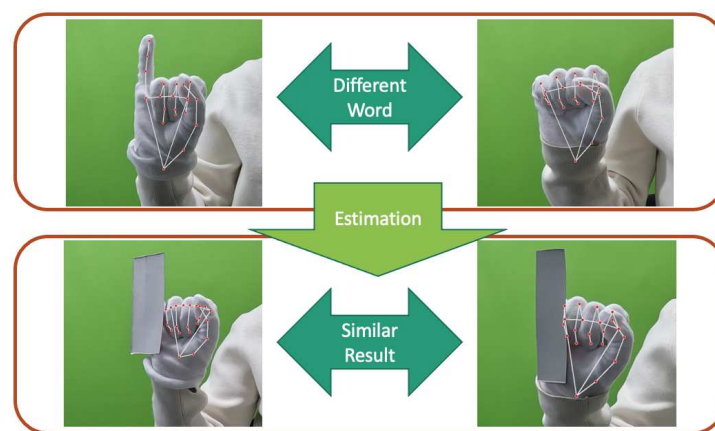
In this research, we integrate a wearable sensor-based system with a camera-based approach to enhance the precision of hand sign language capture. One inherent challenge in extracting skeletal information for sign language is addressing occlusions among fingers and accessing spatial data unattainable by standalone camera systems.

To address this, our proposed system leverages hand skeletons as delineated by MediaPipe for sign language prediction. We adopt a hybrid methodology, intertwining Convolutional Neural Networks (CNN) and Bidirectional Long Short-Term Memory (BiLSTM) models, to bolster our sign language recognition capabilities. CNN is good at extracting relationships between features, and BiLSTMs are adept at temporal data feature comprehension, rendering them ideal for action-oriented tasks such as sign language interpretation. Through this CNN + BiLSTM amalgamation, we have achieved superior recognition accuracy compared to single-sensor solutions.

This research contribution is shown below itemization.

Our devised system integrates visual and bending sensor inputs. Visual data is utilized to extract essential key points and joint angles while eliminating redundancy. This approach mitigates the influence of background and lighting variations, enhancing the system's generalizability and data efficiency. The flex sensor captures finger flexion patterns, enabling adaptability across diverse environments.

We amalgamated key point coordinates, finger joint angles, and curvature features, strategically combining multifaceted information at the feature level. This integration forms the foundation for our CNN-BiLSTM model, facilitating information synergy and effectively enhancing recognition rates.



**Figure 1.** Occlusion Problem in Hand Sign Language.

This paper consists of 1~6 sections. Section 1 explains sign language recognition, the goals of this research issues and solutions, and contribution. Section 2 introduces related works. In related works, we will introduce papers on sign language recognition and hand skeleton prediction, and clarify the purpose of this research. Section 3 and 4 describes the methodology of this research. A specific method for sign language recognition using image recognition and a 2-axis bending sensor is described. Sections 5 and 6 describe the experimental methods and results. Sign language vocabulary is used to compare the results of skeletal information and sensor recognition with the results of skeletal information only. Section 7 presents the discussion and conclusions of this work. Also, we present problems existing in current research and future research.

## 2. Related Works

In recent years, the evolution of wearable hand measurement devices has been evident, predominantly driven by miniaturization processes and advancements in algorithms. Notably, data gloves, including IMU [2] and bending sensors [3,4], have demonstrated significant advancements in wearability, accuracy, and stability metrics. Such advancements have consequently led to marked enhancements in the results of sign language recognition leveraging these measurement apparatuses.

With the evolution of deep learning algorithms, the extraction and analysis of features from visual data, including bone key point prediction, have substantially improved. While sign language recognition has experienced significant advancements, occlusions in images remain a notable challenge in computer vision. Himanshu and Sonia's review discusses the effects of occlusion on the visual system [5]. There are ways to avoid occlusion problems by using a depth camera, multiple cameras, or labeling invisible objects. There are also methods to detect occlusion, such as using shadows of objects and learning information before and after occlusion using time series data.

Therefore, the complementary information of the bending sensor system and the vision system is used to improve accuracy and stability. The application model is shown in Figure 2.

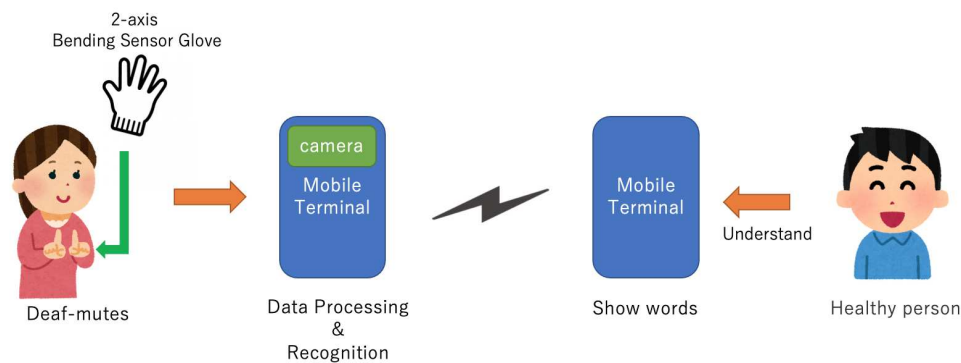


Figure 2. Application Model.

Himanshu and Sonia present a review on occlusion [5]. There are ways to avoid occlusion problems by using a depth camera, multiple cameras, or labeling invisible objects. There are also methods to detect occlusion, such as using shadows of objects and learning information before and after occlusion using time series data.

Avola et al. [6] uses SHREC [8] for the dataset to perform sign language recognition. SHREC is a dataset that uses a depth camera to acquire gesture skeletons. DLSTM, a deep LSTM, is used for sign language recognition. SHREC is used and the angles formed by the fingers of the human hand are used as features. From the predicted skeleton, the finger angles are calculated and used as features. The training using SHREC and DLSTM enables highly accurate sign language recognition.

Liuhao [7] et al. explained the prediction of the skeleton of the hand from image recognition. It estimates the complete 3D hand shape and poses from a monocular RGB image, rather than a depth camera. It uses the original graph convolutional neural network for training. In some cases in this research, recognition accuracy is reduced due to blind spot problems.

Multimodal sensor data fusion methods are crucial in systems that combine curved sensors and vision. CNN [8] and BiLSTM [9] methods, which can obtain information from spatial and time series data. Fusion of CNN and BiLSTM [10] has been used in the field of Natural language processing. Also, The skeleton of the hand using a method called MediaPipe [11] from videos. In addition, by using the sensor, we can expect to measure the angle of the finger more accurately even in the part that overlaps other objects. Therefore, combining sensor data with sign language recognition will make it possible to accurately predict hand movements.

3. Method

In this paper, we use hand skeletons predicted by MediaPipe to predict sign language. To improve the prediction of sign language recognition, a combination of CNN and BiLSTM methods are used for prediction. Predict words using CNN and BiLSTM from the skeleton predicted by MediaPipe. CNN is a commonly used method in image recognition. In addition, BiLSTM can learn time-series data, so it is suitable for learning motions such as sign language. CNN learns spatial features of symbols, while BiLSTM learns temporal features. By combining CNN + BiLSTM, we obtained higher recognition accuracy than only with one kind of sensor.

3.1. MediaPipe

We use MediaPipe to predict skeletons from images. MediaPipe can predict face, posture and hand skeleton with high accuracy. This method is intended for use with GPUs for real-time inference. However, there are also lighter and heavier versions of the model to deal with CPU inference on mobile devices which is less accurate than running on desktops [12]. Fig:3.1 is the output of MediaPipe hand skeleton data. (a) are the predicted 21 keypoint positions. In (b), the points in (a) correspond to the numbers. (c) is an example of using MediaPipe. In this research, 21 keypoints indicated by red dots are used as skeleton data and used as a dataset.

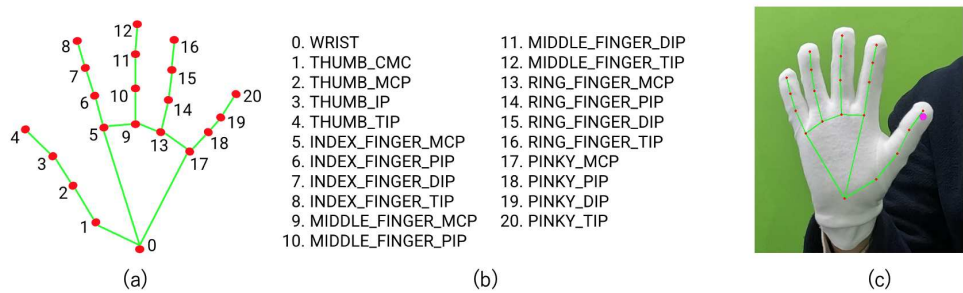


Figure 3. Skeleton and Bending Sensor Data Fusion.

3.2. CNN+BiLSTM

Since video data is used for sign language recognition, a method that processes both spatial information and time series data is effective. Spatial information is learned using CNN, and time series information is learned using BiLSTM. First, a sign language dataset is input to MediaPipe. MediaPipe outputs the keypoint data of the sign language, which is used as skeleton data. The skeleton data is then input to the CNN to extract spatial information, and then temporal information is extracted by BiLSTM. The spatial and temporal information is learned and used as a model. By combining CNN and BiLSTM, we have achieved higher recognition accuracy by learning spatial and temporal features than only with one kind of them.

3.3. 2-Axis Bending Sensor

The sensor used is a 2-axis bending sensor3.3 developed by Bend Labs. Compared to conventional sensors, this sensor measures angular displacement with higher accuracy in terms of power loss. The sensor output is the angular displacement as computed from the vectors defined by the ends of the sensor (v1 and v2). [13]





Figure 4. 2-axis bending sensor from Bend Labs.

4. System Design

4.1. System Outline

The sign language recognition system in this study supplements the results of MediaPipe with sensor data to improve the occlusion problem in sign language recognition. The model of this method is shown in Figure 4.1. So we need a sensor system, a dataset, and an AI model. First, we produced the 2-axis bending sensor glove for collecting finger angles. Collecting video data with the camera, and collecting sensor data with the sensor glove is at the same time. Next, collected data is fused. The data to be fused are keypoint data, joint angle, and sensor data. The fused data is first trained with a CNN to extract features. Then it is input to BiLSTM through LeakyReLU. LeakyReLU is an extension of the activation function ReLU. If the neuron’s output is zero, the gradients of the weights and parameters are zeroed by ReLU. We choose the LeakyReLU activation function to avoid dying neurons. Finally, through BiLSTM, the predicted words are output

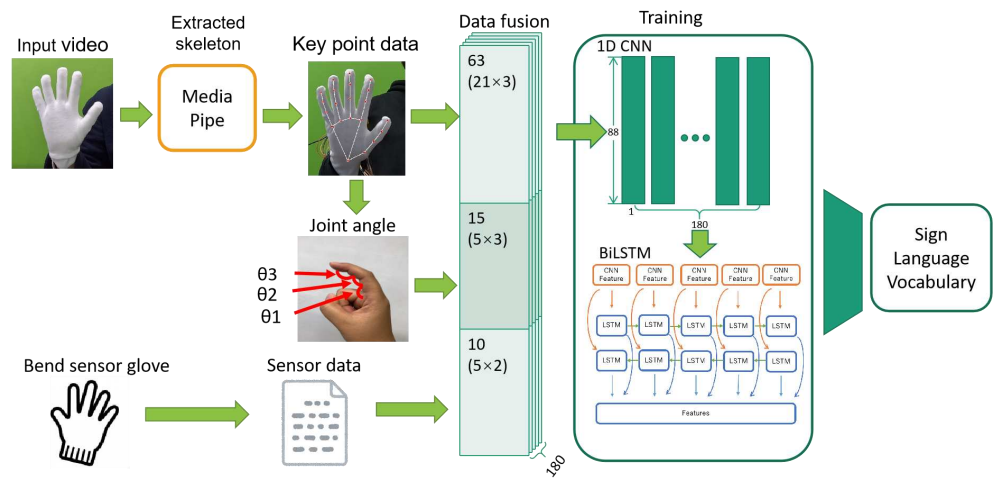


Figure 5. Research Method: Data collection and Training.

5. Implementation

5.1. Outline

Predicts hand movements using image and sensor data. First, we create a dataset. The dataset is a video of the sign language and the finger angles from the sensor data. Next, the hand skeleton is predicted from the sign language video. The hand skeleton is estimated using a MediaPipe. Next, the sensor data and the skeletal data are fused and trained with BiLSTM + CNN Finally, a model for gesture estimation is formed.

5.2. Sign Language Dataset

First, create a dataset of sign language videos to create skeleton data. The dataset is original data from laboratory members. Sign language words are used in 32 Japanese sign language vocabulary(SLV). The Japanese language is represented by 46 letters. They are represented by vowels (a, i, u, e, o), and consonants (k, s, t, n, h, m, y, r, w). The letter list used in this research is shown below 5.2. Japanese has letters with vowels only, vowels and consonants, and special characters represented by "nn". The table shows consonants in columns and vowels in rows. The first column from the right is for vowels only("/" mean no consonants), and "nn" appears at the end of the column for the consonant n.

5.3. Image Data Collection

The dataset has videos of 4 people for each word shot at 60fps with a green screen background. The sensor glove is put on the righ hand. SLV is basically fixed, such as clenching a fist or raising the only index finger, and the hand is not moved. However, some SLV is expressed by moving the hand. "ri", "no", "nn", and "mo" in the wordlist table5.2. "mo" is a finger movement only, but "ri", "no" and "nn" are expressed by moving the wrist.

w	r	y	m	h	n	t	s	k	/	
										a
										i
										u
										e
										o
										n

Figure 6. Japanese Sign Language Letter List.

5.3.1. Key Point Estimation

We predict skeletal data from videos of sign language wearing sensor gloves. MediaPipe estimate 21 key points and make them skeleton data. Keypoint coordinates are 3D(x, y, z) and 60 frames are acquired per second.

5.3.2. Calculating Joint Angle

Calculate finger angles from skeleton data obtained with MediaPipe. This is useful for data argumentation of the dataset. There is one finger angle for each joint, and angles are calculated by the inner product. For example, to calculate the angle of the pinky finger, the keypoint k is predicted by the media pipe and calculated using

5.4. Collecting Sensor Data

This section describes the original Bending sensor glove and finger angle data collection. We made an original Bending sensor glove to collect finger angles. The glove is worn on the right hand. The data collected while wearing the glove is saved as a text file on the main computer along with time

stamps and angles of five fingers 2 axis angles. In addition, a video of the sign language is also filmed at the same time as the bending sensor data is collected. The angle of the finger acquired at the same time as the bending sensor data and the image acquired at the same time support image recognition.

5.5. Bending Sensor Glove Structure

This part describes the design of the original glove, the sensors, the sensor controllers, and the sensor data structures. Figure 5.4 is the actual 2-axis bending sensor glove. Secure the fingertips and loosely secure the rest so that the sensor does not come loose. Therefore, fixing parts was created with a 3D printer. The fingertip part is designed so that the sensor can be inserted and fixed. Also, if every part fixes the sensor, the movement of the finger will be restricted, making it impossible to express sign language. Therefore, the part other than the tip is not fixed. Also, when actually using it, wear white gloves to hide the sensor. This will prevent from MediaPipe not recognizing the sensor glove as a hand. Then Raspberry pi pico is used as a controller to control the sensor. Note that sensor gloves have different values depending on the person using the same hand pose.

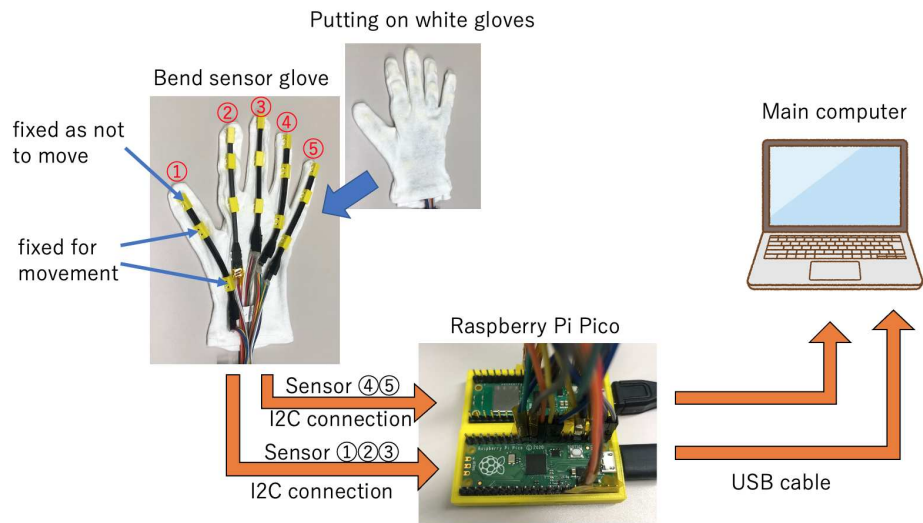


Figure 7. Sensor Glove Design.

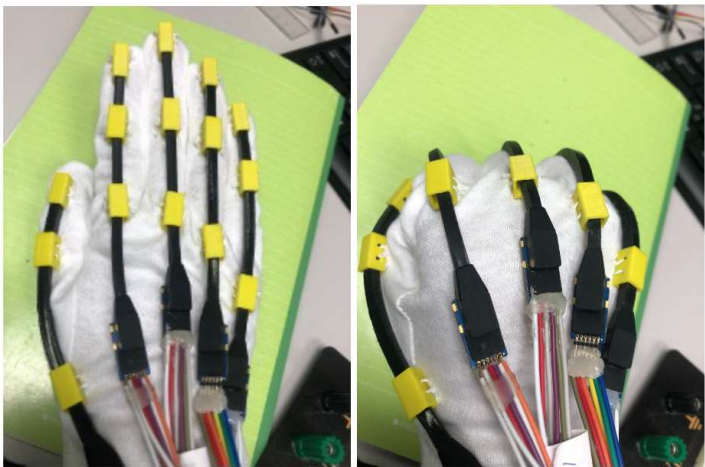


Figure 8. Sensor Glove.



5.6. Data Fusion

First, Skeleton Data is acquired by MediaPipe, finger joint angles are calculated from Skeleton data, and sensor data is fused. Skeleton data is 63 (21 key points \* 3 dimensions), finger joint angle is 15 (5 fingers \* 3 joint angles), and sensor data is 10 (5 fingers \* 2-axis).

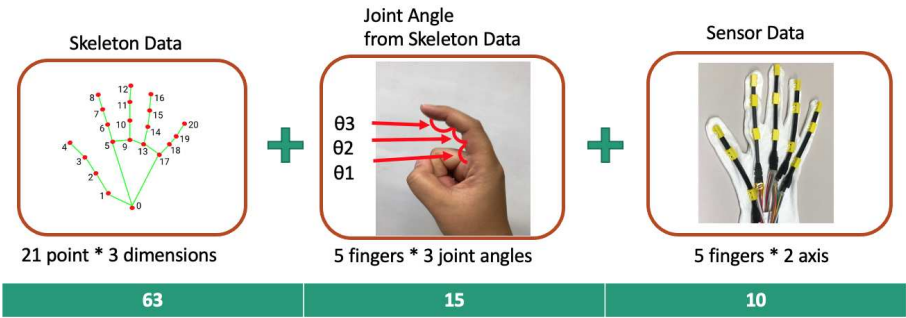


Figure 9. Data Fusion.

6. Experiment and Evaluation

6.1. Experiment Purpose

In this experiment, we will design the experiment with the aim of improving the blind spot problem and recognition rate. In addition, we will evaluate the Bending sensor gloves' performance and show the gloves' usefulness. Experimental evaluation and discussion will be made by comparing results with and without Bending sensor gloves.

6.2. Experiment Design

First, we evaluate the MediaPipe only. We assigned ID's to words and output the results. The table6.1 shows the words used in this research. The dataset has 10 videos of the vocabulary for every 4 people. From the training results we found that the recognition rate for the blind spot problem and fingerprints that require motion is low. As an example, "tu", and "i" are similar and considered difficult to recognize based on temporal features alone. We also found that the recognition rate of letters that require wrist and arm movements is low. Although CNN learns spatial features, there is a limit using only a two-dimensional monocular camera.

From the above, to solve the occlusion problem, we add an occlusion point in sign language motion. First, use a video with occlusion added and perform recognition with MediaPipe. Input occlusion images to the trained model and compare the results with image-only data and fused data.

6.3. Experiment Setting

We will add about the experimental environment. First, we prepared a Bending sensor glove and a camera to collect data. The camera uses GoPro Hero10. High resolution, fast and small. Also, use a green screen for the background and unify the background colors. Wear sensor gloves and collect sensor data and video data. There are 32 sign language words, and the words shown in the table6.1 are used.

To demonstrate the effectiveness of the sensor, occlusion is generated in the SLV video and recognized by MediaPipe. First, we generate an occlusion in the image. Occlusion is expressed by randomly selecting the coordinates of the keypoint obtained with MediaPipe and displaying a black square on it. Occlusion is (80, 80) for image size (1080, 1920).

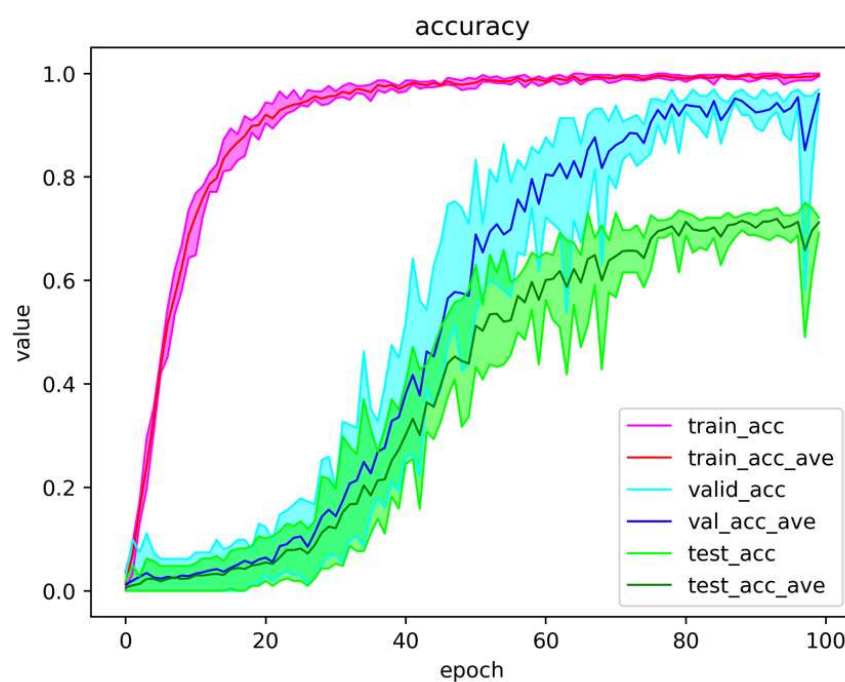
#### 6.4. Experiment Process

For example, the words "tu" and "i" are very similar, but the difference is whether the hand is fully clasped or not 6.1. The sensor data also shows that "ti" has less movement. From the sensor data, we can expect to improve the system of spatial prediction, which has been limited by image recognition.

Occlusion is generated at random positions for each video 6.2. Occlusion was generated by inserting black squares at random positions in the image. However, MediaPipe may not be able to get the Keypoint if an occlusion occurs. MediaPipe acquires skeleton data for each frame, but if the keypoint cannot be acquired in the first frame, the output result of MediaPipe is  $(x, y, z) = (-1, -1, -1)$ . If the frame is in the middle, the output result is the value of the previous frame.

#### 6.5. Experiment Results

The model was trained with k-Fold cross validation. For training with a small data set, the training accuracy during training could be higher. If this is the case, the accuracy in training may be high, but the accuracy in testing may be lowered, resulting in over-fitting. To prevent this situation, there is a technique called k-Fold cross validation. In k-Fold cross-validation, data is divided into k pieces, some of which are used for validation data and others for training data. Since all the divided data are used once for validation data, training is performed k times. The average of the k training accuracies is calculated as the result. The cross entropy method is calculated for the loss function. If the probability distributions of p and q are approximate, the cross-entropy loss is smaller. In other words, the closer the learning accuracy approaches 1, the closer the result approaches 0. Results for the MediaPipe only are shown below. There were 2282 number of samples extracted from the MediaPipe. There are 261 test data, and the remaining data is training and evaluation data. Also, training data and evaluation data are split at a ratio of 4:1, there are 640 training data and 275 evaluation data. Training data is used for training, evaluation is used for evaluation during training, and test data is used for model evaluation.



**Figure 10.** Accuracy curve of Only Skeleton Data.

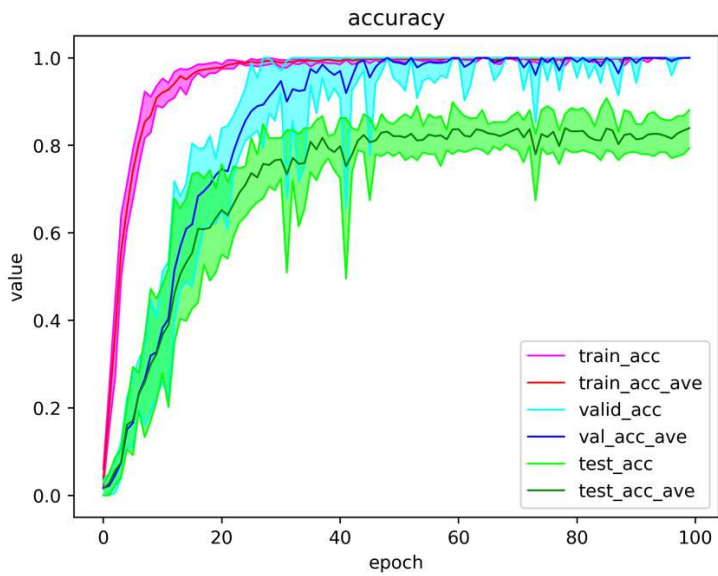


Figure 11. Accuracy curve of Fusion Data.

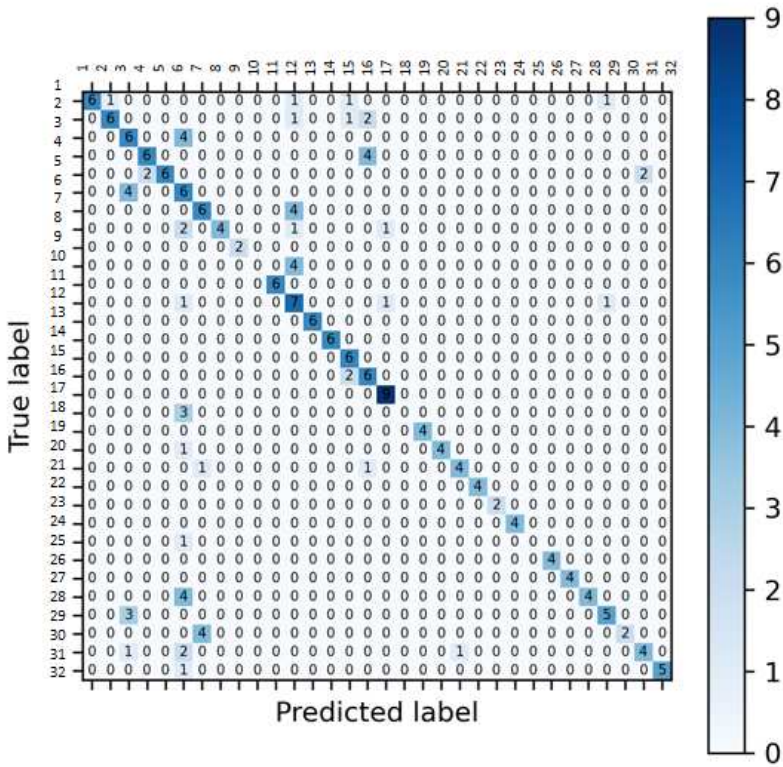


Figure 12. Confusion Matrix: Only Skeleton Data.

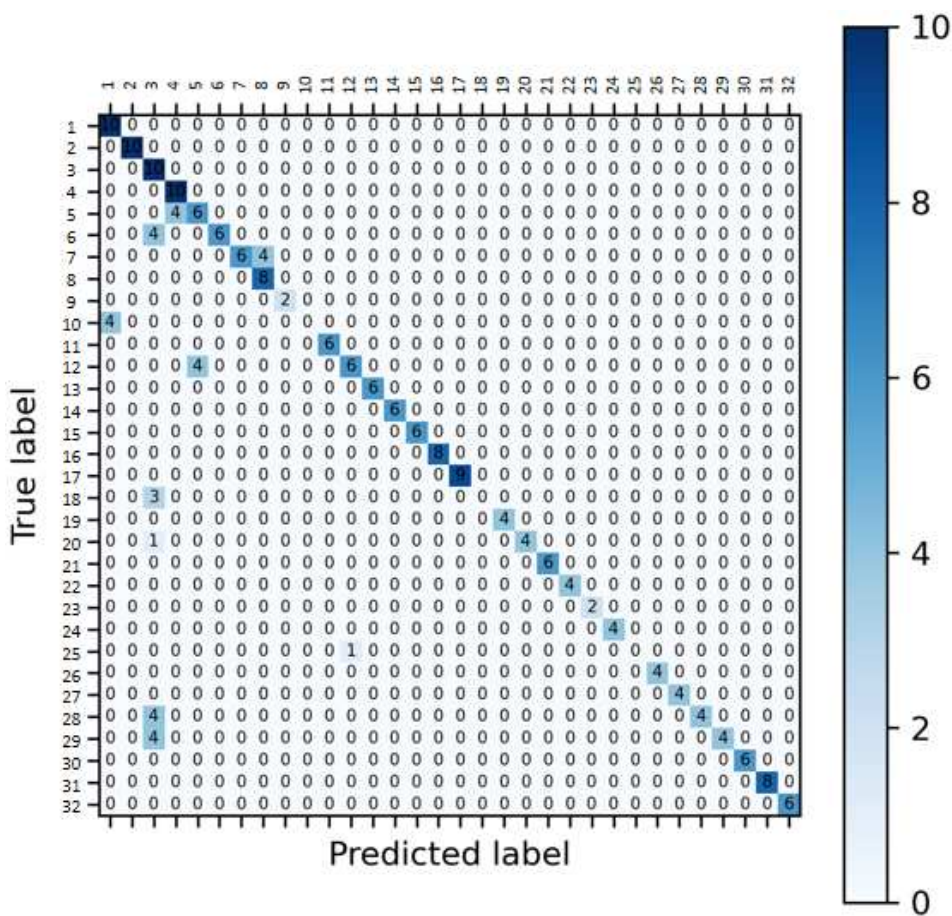


Figure 13. Confusion Matrix: Fusion Data.

The training curve is shown below. The blue line shows the accuracy of Training, the orange line shows the accuracy of Validation, and the green line shows the accuracy of validation with Test data. For the Skeleton-only validation, cross-validation was performed 5 times, with an average training accuracy of 85.9% when training and 73.5% when using test data 6.3. For the Fusion data validation, cross-validation was performed 5 times, with an average training accuracy of 99.2% when training and 96.5% when using test data 6.5.

Figure6.4 shows the loss curve during learning. The blue line shows the loss during learning, the orange line shows the validation loss, and the green line shows the validation loss of the test data. It can be seen that as the number of training epochs increases, the learning loss and validation loss decrease and stabilize at about 100 epochs. the closer the loss curve is to 0, the better the learning accuracy

6.6. Discussion

By merging the datasets, we were able to improve the recognition accuracy. However, there were some words that could not be recognized even with the fusion data. Words that could not be recognized had a poor recognition rate in MediaPipe in the first place, and there was little training data. In addition, when comparing the misrecognized sign language, the sensor data values were similar, and the learning result was even lower. When MediaPipe’s recognition rate is poor, we expected sensor support but learned that many values were complemented with (-1, -1, -1) when complementing loss values in MediaPipe Accuracy may be lost. Also, since the data set this time was a motion without hand movement, it is necessary to verify motion with motion in order to bring it closer to real sign language.



## 7. Conclusion

In this research, we aimed to improve sign language recognition with occlusion accuracy by combining CNN+BiLSTM and also combining bending sensor data with skeleton data. The combination of CNN+BiLSTM method allowed us to perform finger character recognition better than using it alone. However, there were limitations in acquiring spatial information, such as blind spot problems. Therefore, we used a 2- axis bending sensor to assist with spatial information. The performance evaluation of the original 2-axis bending glove further strengthened the spatial information of sign language. By using sensor data, we were able to improve sign language recognition accuracy in the presence of occlusion compared to skeleton data alone. Our future tasks are how to deal with the sensor when the recognition rate in MediaPipe is poor, and how to deal with motion.

## References

1. WorldHealthOrganization. World report on hearing, 2021. Available online: <https://www.who.int/news-room/fact-sheets/detail/deafness-and-hearing-loss>.
2. Lu, C.; Dai, Z.; Jing, L. Measurement of Hand Joint Angle Using Inertial-Based Motion Capture System. *IEEE Transactions on Instrumentation and Measurement* **2023**, *72*, 1–11.
3. Faisal, M.A.A.; Abir, F.F.; Ahmed, M.U.; Ahad, M.A.R. Exploiting domain transformation and deep learning for hand gesture recognition using a low-cost dataglove. *Scientific Reports* **2022**, *12*.
4. Lu, C.; Amino, S.; Jing, L. Data Glove with Bending Sensor and Inertial Sensor Based on Weighted DTW Fusion for Sign Language Recognition. *Electronics* **2023**.
5. Purkait, P.; Zach, C.; Reid, I.D. Seeing Behind Things: Extending Semantic Segmentation to Occluded Regions. *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* **2019**, pp. 1998–2005.
6. Avola, D.; Bernardi, M.; Cinque, L.; Foresti, G.L.; Massaroni, C. Exploiting Recurrent Neural Networks and Leap Motion Controller for the Recognition of Sign Language and Semaphoric Hand Gestures. *IEEE Transactions on Multimedia* **2018**, *21*, 234–245.
7. Ge, L.; Ren, Z.; Li, Y.; Xue, Z.; Wang, Y.; Cai, J.; Yuan, J. 3D Hand Shape and Pose Estimation From a Single RGB Image. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* **2019**, pp. 10825–10834.
8. O'Shea, K.; Nash, R. An Introduction to Convolutional Neural Networks. *ArXiv* **2015**, *abs/1511.08458*.
9. Zhang, S.; Zheng, D.; Hu, X.; Yang, M. Bidirectional long short-term memory networks for relation classification. *Proceedings of the 29th Pacific Asia Conference on Language, Information and Computation*, 2015, pp. 73–78.
10. Chiu, J.P.C.; Nichols, E. Named Entity Recognition with Bidirectional LSTM-CNNs. *Transactions of the Association for Computational Linguistics* **2015**, *4*, 357–370.
11. Lugaresi, C.; Tang, J.; Nash, H.; McClanahan, C.; Uboweja, E.; Hays, M.; Zhang, F.; Chang, C.L.; Yong, M.G.; Lee, J.; Chang, W.T.; Hua, W.; Georg, M.; Grundmann, M. MediaPipe: A Framework for Building Perception Pipelines. *ArXiv* **2019**, *abs/1906.08172*.
12. Zhang, F.; Bazarevsky, V.; Vakunov, A.; Tkachenka, A.; Sung, G.; Chang, C.L.; Grundmann, M. MediaPipe Hands: On-device Real-time Hand Tracking. *ArXiv* **2020**, *abs/2006.10214*.
13. Soft Angular Displacement Sensor Theory Manual. 2018.

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.