

# Deep Learning Approach For Sign Language Recognition

*By Bambang Krismono Triwijoyo*

# Deep Learning Approach For Sign Language Recognition

Bambang Krismono Triwijoyo<sup>1</sup>, L Yuda Rahmani Karnaen<sup>2</sup>, Ahmat Adil<sup>3</sup>

1,2,3  
17  
Universitas Bumigora, Jl. Ismail Marzuki No.22, Mataram 83127, Indonesia

## ARTICLE INFO

### Article history:

Received October 15, 2022  
Revised January 15, 2023  
Accepted

### Keywords:

Deep Learning;  
Sign Language;  
CNN;

## ABSTRACT

Sign language is a method of communication that uses hand gestures between people with hearing loss. Each hand sign represents one meaning, but several terms don't have sign language, so they have to be spelled alphabetically. Problems occur when communicating between normal people with hearing loss, because not everyone understands sign language, so a model is needed to recognize sign language as well as a learning tool for beginners who want to learn sign language, especially alphabetic sign language. This study aims to create a hand sign language recognition model for alphabetic letters using a deep learning approach. The main contribution of this research is to produce a real-time hand sign language image acquisition, and hand sign language recognition model for Alphabet. The model used is a seven-layer Convolutional Neural Network (CNN). This model is trained using the ASL alphabet database which consists of 27 categories, where each category consists of 3000 images or a total of 87,000 hand gesture images measuring 111 x 200 pixels. First, the background correction process is carried out and the input image size is changed to 32 x 32 pixels using the bicubic interpolation method. Next, separate the dataset for training and validation respectively 75% and 25%. Finally the process of testing the model using data input of hand sign language images from a web camera. The test results show that the proposed model has good performance with an accuracy value of 99%. The experimental results show that image preprocessing using background correction can improve model performance.

2

This work is licensed under a Creative Commons Attribution-Share Alike 4.0



### Corresponding Author:

Bambang Krismono Triwijoyo, Universitas Bumigora, Jl. Ismail Marzuki No.22, Mataram 83127, Indonesia  
Email: [bkrismono@universitasbumigora.ac.id](mailto:bkrismono@universitasbumigora.ac.id)

## 1. INTRODUCTION

The research background is communication is very important in the process of social interaction. Communication leads to better understanding among the community, including the deaf [1]. Hand gesture recognition serves as the key to overcoming many difficulties and providing convenience for human life, especially for the deaf [3]. Sign language is a structured form of hand movement that involves visual movement and the use of various body parts namely fingers, hands, arms, head, body, and facial expressions to convey information in the communication process. For the deaf and speech-impaired community, sign language serves as a useful tool for everyday interactions [4]. However, sign language is not common among normal people, and only a few people understand sign language. This creates a real problem in communication between the deaf community and other communities, which has not been fully resolved to this day [3]. Not all words have sign language, so special words that do not have sign language must be spelled using a letter sign one by one [5]. Based on the background, this study aims to develop a sign language recognition model for letters of the alphabet using a deep learning approach. The deep learning approach was chosen because deep learning methods are popular in the field of computer science and are proven to produce a good performance for image classification [6][7]. The novelty of this study is the

5

application of resizing and background correction of the input image for training and testing to improve model performance, where the results of testing the model we propose are better than previous similar studies.

The related work from past research is as follows. There have been many studies to recognize sign language, using various methods and varied datasets. Researchers [5] proposed a language recognition system using a 3D motion sensor by applying a k-nearest neighbor and support vector machine (SVM) to classify 26 letters in sign language. The average results of the highest classification levels of 72.78% and 79.83% were achieved by k-nearest neighbors and support vector machines, respectively. Based on previous studies that proposed hand sign language recognition models for alphabets, the results were not optimal, this due to the complexity of lighting factors and other objects that appear in hand gesture images [5]. While there have been quite many studies on sign language recognition using a deep learning approach, here are several related studies including Study [8] has proposed a recognition system using a convolutional neural network (CNN) that can recognize 20 Italian gestures with high accuracy. Meanwhile, the following researchers introduced a sign language recognition (SLR) model using a deep learning approach. Study [9] implements transfer learning to improve accuracy. While study [10] has proposed the Restricted Boltzmann Machine (RBM) for automatic hand sign language recognition from visual data. The experimental results using four datasets show that the proposed multi-modal model achieves a fairly good accuracy. In the work [11], have proposed a deep learning-based framework for analyzing video features (images and optical flow) and skeletons (body, hands, and faces) using two sign language datasets. The results reveal the advantages of combining frame and video features optimally for SLR tasks.

A continuous deep learning-based sign language recognition model has also been introduced [12] which has proposed a 3D convolution residual network architecture and two-way LSTM, as a grammatical rule-based classification problem. The model has been evaluated on the benchmark of Chinese continuous sign language recognition with better performance. Other deep learning approach models have also been carried out for sign language recognition. A study [13] has proposed a ResNet50 Based Deep Neural Network architecture to classify finger-spelled words. The dataset used is the standard American Sign Language Hand gesture which produces an accuracy of 99.03%. While the study [14] Densely Connected Convolutional Neural Networks (DenseNet) to classify sign language in real-time using a web camera with an accuracy of 90.3%. The following studies [15][16][17][18][19] have implemented the CNN model for sign language recognition and tested using the American Sign Language (ASL) dataset, with an accuracy rate of 99.92%, 99.85%, 99.3%, 93%, and 99.91%.

Based on previous related studies, most of the sign language recognition methods use a deep learning approach. This study focuses on the introduction of hand sign language from the letters of the alphabet, which is used as a means of communication with the deaf. This study also uses the CNN model but with a different model architecture from previous studies. The CNN model was chosen because previous studies showed relatively better accuracy for image recognition [20][21][22]. The contributions of this research are: first, produce a real-time hand sign language image acquisition model by capturing each frame using a webcam video. Second, produce a hand sign language recognition model for the Alphabet, using a seven-layer CNN that has been trained using the ASL dataset and by applying resizing and background correction to the input image.

## 2. METHOD

This research is quantitative experimental research to measure the performance of hand sign language recognition models based on training datasets. Fig. 1 shows the proposed method for hand sign language recognition. In general, the proposed method consists of four stages, each of which is data acquisition, preprocessing, training and testing.

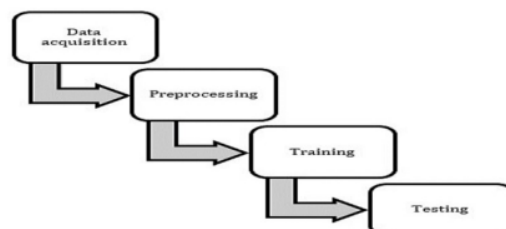


Fig. 1. Proposed Method

Based on the methodology applied in this study, as shown in Fig. 1, the first stage is data acquisition, where the data used in this study is the image. Image acquisition is the action of retrieving an image from an external source for further processing [23]. In this stage, the dataset used as model input is hand sign images, which are divided into 26 classes consisting of 26 sets of alphabets from A to Z. In this study, the dataset used as model input is hand sign images, which are divided into 26 classes consisting of 26 sets of alphabets from A to Z. The second stage is preprocessing. At this stage, the image size transformation is carried out to reduce the complexity of the model architecture. In this study, the transformation of the image size of the training data images from the initial size of 200x200 pixels was resized to 32x32 Pixels. In this study, we apply the bicubic interpolation method for resizing images as proposed by [24]. This resizing process is to reduce the computational time required for model training. To improve the segmentation accuracy of the hand sign language image below for complex lighting conditions, we apply a correction background method with luminance partition correction and adaptive threshold [25][26]. Furthermore, at the training stage, the CNN model architecture and its hyperparameters will be determined first. In this study, we use hyperparameter tuning to control the behavior of the machine learning model to produce optimal results. [27][28], then model training will be carried out using the dataset from the hand sign language image preprocessing. The last stage is model testing. At this stage, the model will be tested with hand sign language images in real-time using a webcam. The research proposed method flowchart is presented in Fig. 2.

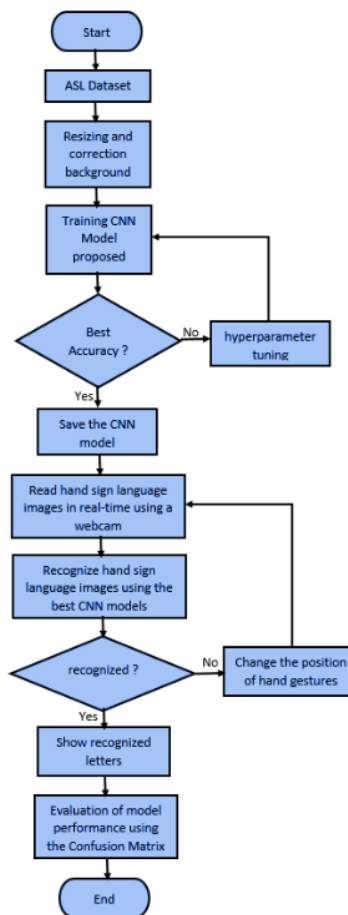


Fig 2. The Flowchart of Research Methodology

Furthermore, the results will be measured using a confusion matrix to determine the performance of the model [29]. Confusion matrices create result representations such as true positives (TP), true negatives (TN),

false positives (FP), and false negatives (FN). TP means a positive result that is predicted by machine learning correctly. TN means negative outcome predicted by machine learning Correct. While FP means positive results predicted by machine learning are wrong, and FN means negative outcomes predicted by machine learning are wrong. Confusion Matrix Performance evaluation with a confusion matrix results in accuracy, precision, and recall [30][31]. Accuracy is the number of data points that machine learning correctly predicts among all data. Can calculate as eq.1:

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN} \quad (1)$$

Precision is the percentage of relevant elements that can indicate the number of times the model can predict correctly and can be calculated as eq. 2.

$$Precision = \frac{TP}{TP + FP} \quad (2)$$

Meanwhile, recall is the percentage of relevant elements that are correctly classified by the model above all relevant elements. Recall calculation can be done using eq. 3.

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

### 3. RESULTS AND DISCUSSION

#### 3.1. Data Acquisition

The data is secondary data downloaded from kaggle.com in the ASL Alphabet repository which contains image datasets. in JPG format containing 29 folders with each containing 3,000 hand sign images. The total image data obtained amounted to 87,000 images [32]. Fig. 3 shows examples of hand sign language images sourced from the ASL Alphabet repository. The 87,000 images were distributed for the training process. The data is divided into two 75% for training and 25%. The train-test split is a technique for evaluating the performance of a machine-learning algorithm [33]. In the test set, so for each class, there are 2,400 images for training and 600 images for testing, or a total of 69,600 images for training and 17,400 images for testing.

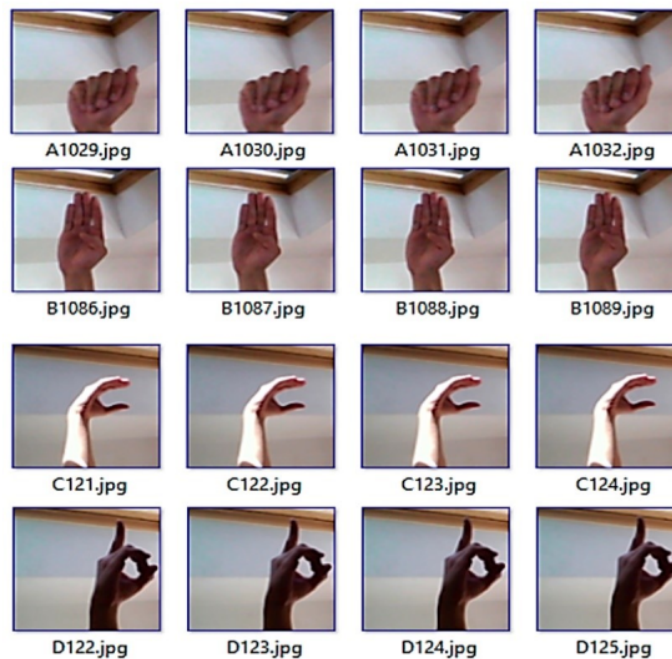


Fig. 3. Examples of hand sign language images from Kaggle



### 3.2. Preprocessing

At the preprocessing stage, the resizing process is carried out using the bicubic interpolation method [24], the results of this process produce an image size of 32x32 pixels, from the original image size of 200x200 pixels. This step is taken to reduce the time complexity during model training. the next preprocessing step is image background correction, to produce better accuracy using luminance partition correction and adaptive threshold [25]. Fig 4 shows examples of preprocessing results.

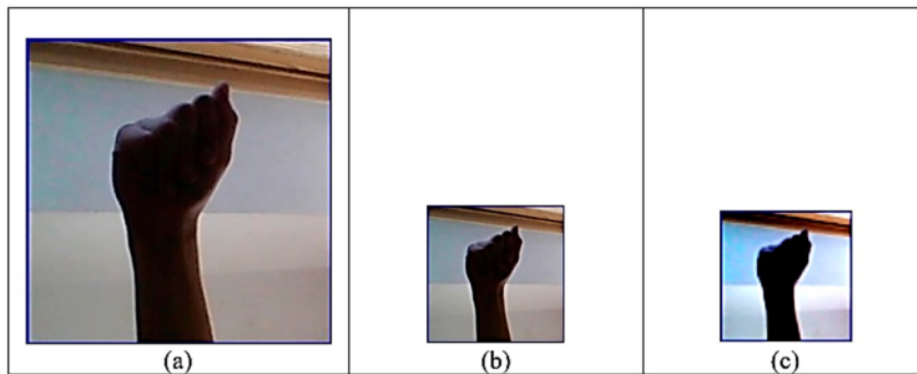


Fig. 4. (a). Original image size 200x200 pixels, (b). The result of resizing to a size of 32x32 pixels, (c). Background correction results.

In Fig.4 it can be seen that the dimensions of the hand sign language image size for the training dataset are reduced in size as well as the increase in color brightness and contrast resulting from the background correction process.

### 3.3. Training

At this stage, the CNN model design used in the training process is carried out to produce an appropriate model to classify hand sign language images. The CNN model applied uses hyperparameter values consisting of the learning rate, epoch, loss function, and optimizer. Table 1 below is the CNN architecture specification used in this study.

Table 1. CNN Model Architecture

Layer type	Description	Size
Input Layer	Input Image	32x32x3
Conv 1	Convolutional ReLU MaxPooling	8 kernel, 3x3 Window size 2x2
Conv 2	Convolutional ReLU MaxPooling	16 kernel, 3x3 Window size 2x2
Conv 3	Convolutional ReLU MaxPooling	32 kernel, 3x3 Window size 2x2
Flatten	Flatten	512
Fully Connected Layer	Dense ReLU	512
Output Layer	Dense Softmax	29

The CNN architectural design is as shown in table 1, there are input layers, 3 convolution layers, flatten, fully connected layers, and output layers. The Input Layer requires input with a size of 32x32x3 where 32x32 pixels 3 layers RGB, then in Conv 1 using 8 kernels measuring 3x3 after that using ReLU activation [34], and using Maxpool with a window size of 2x2, as well as Conv 2 and Conv 3, only on Conv 2 uses 16 kernels and Conv 3 uses 32 kernels. Next is the Flatten layer and Fully Connected Layer with 512 nodes, and

the last is the output layer containing Dense Softmax with 29 nodes [35], this is adjusted to the number of classes in the hand signal dataset of the alphabet. 4

The training process will be carried out through a series of iterations whose number of repetitions is determined by the maximum epoch value [36]. One epoch is a process when all training data has been used and passes through all network nodes once. While the hyperparameter used in this study is the Adam Optimizer [37], with The Batch Size is 32 and Epochs are 20. The model training process is carried out in a hardware and software environment with specifications for Dell Latitude E7440 Laptop, 12 GB DDR4 RAM, Processor Intel® Core™ i5-7200U CPU @ 2.50GHz, Nvidia GeForce 940MX GPU, 256 GB SSD. The operating system used is Windows 10 Professional and the training and testing model algorithms are implemented using python code by utilizing the Tensorflow library. The results of the training model as shown in Fig. 5.

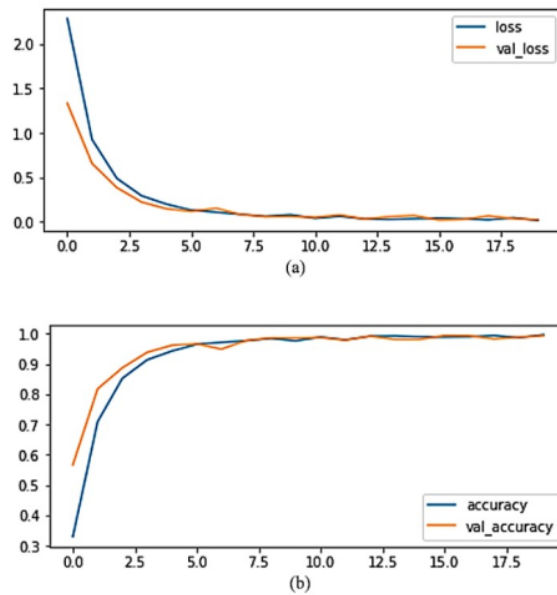


Fig. 5. Graphics of the Training Model Process: (a) Loss of Training and Validation, (b) Accuracy of Training and Validation

20  
As seen in Fig. 5, The training process is carried out in 20 epochs, in the first epoch training accuracy is 33.04% and validation accuracy is 56.69%, while training loss is 28.26% and validation loss is 33.30%. in the tenth epoch, the training accuracy is 97.60% and validation accuracy is 98.57%, while the training loss is 8.11% and validation loss is 6.08%. Finally, in the last or twentieth epoch, the training accuracy is 99.60% and the validation accuracy is 99.68%, while the training loss is 1.64% and the validation loss is 2.55%. The main finding of this research is that the use of resizing and background correction methods, as well as setting hyperparameters can improve model accuracy.

### 3.4. Testing

After the model training process is complete, the next process is to test the model. The testing process is similar to the model training process, only the difference is that when the model testing process is not carried out backward pass or backpropagation iterations it does not change the weight or weight of the model as in the training process [38]. The Testing process is carried out using 23 data that is different from the training data set, to obtain valid testing results, this is following the recommendations of [39][40] regarding training and testing of the CNN model. Model testing is done by reading the hand signal image that is inputted by taking each frame from the webcam video, then the frame is identified based on the model that has been trained, then produces output in the form of identification results and accuracy values. Then the highest value from the identification results is directly written on the output board. Fig. 6 is a flowchart for the model testing process.

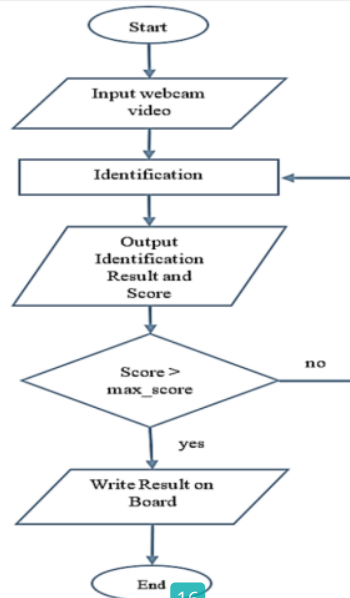


Fig. 6. Flowchart for the testing process

As shown in Fig. 6, the testing process is carried out by entering input in the form of hand signals via the Web camera on the Laptop Computer, where the hand makes a hand gesture in the Region of interest box on the Webcam display. Then the webcam display brings up the alphabetical prediction of the hand signal along with the score on the board display. Testing the model is carried out for each character letter 10 times, so a total of 290 times of testing. Fig. 7 and Fig. 8 Show the Confusion matrix value from the model testing results, for 29 types of hand sign language for each letter of the alphabet. Which consists of accuracy, precision, and recall.

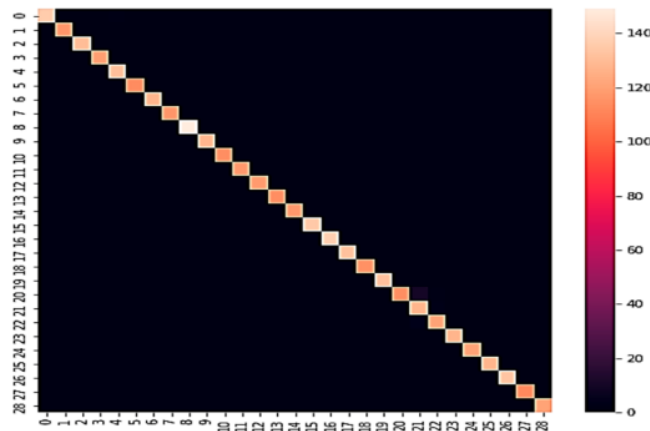


Fig. 7. Confusion Matrix of The Model Testing



A	1.00	0.99	0.99	138
B	1.00	1.00	1.00	116
C	1.00	1.00	1.00	130
D	1.00	1.00	1.00	119
E	0.98	1.00	0.99	131
F	1.00	1.00	1.00	113
G	1.00	1.00	1.00	127
H	1.00	0.99	1.00	117
I	1.00	1.00	1.00	149
J	0.99	1.00	1.00	128
K	1.00	1.00	1.00	114
L	1.00	1.00	1.00	117
M	1.00	1.00	1.00	118
N	1.00	1.00	1.00	114
O	1.00	1.00	1.00	117
P	1.00	0.99	1.00	136
Q	1.00	1.00	1.00	137
R	1.00	1.00	1.00	132
S	0.99	0.98	0.99	117
T	0.99	1.00	1.00	132
U	1.00	0.90	0.95	126
V	0.90	0.98	0.94	130
W	0.95	0.97	0.96	125
X	1.00	1.00	1.00	127
Y	0.99	0.99	0.99	122
Z	1.00	1.00	1.00	127
del	1.00	1.00	1.00	133
nothing	1.00	1.00	1.00	112
space	1.00	1.00	1.00	121
accuracy			0.99	3625
macro avg	0.99	0.99	0.99	3625
weighted avg	0.99	0.99	0.99	3625

Fig. 8. Performance of The Testing Model

Fig. 8 shows that the performance of the testing model achieves the best accuracy of 99%. and the Sensitivity, Recall, and f1-score levels are 99% respectively. This is because we apply resizing and background correction to the training and test images. These results are relatively better than previous similar studies as shown in Table 2.

Table 2. Comparison of Results with previous related studies

Study	Model	Dataset	Accuracy
Chong et al <sup>[5]</sup>	SVM, DNN	American Sign Language (ASL)	80.30% and 93.81%
Pigou et al <sup>[8]</sup>	CNN	Italian gestures	91.7%
Rastgoo et al <sup>[10]</sup>	RBM	Massey University Gesture Dataset 2012	99.31%
Rathi et al <sup>[13]</sup>	ResNet50 based	American Sign Language Hand gesture	99.03%
Daroya et al <sup>[14]</sup>	CNN	American Finger Spelling format	90.03%
Abdulhusein et al <sup>[17]</sup>	CNN	American Sign Language (ASL)	99.3%
Sabeenian et al <sup>[18]</sup>	CNN	NIST ASL dataset	93 %
Al-Hammadi et al <sup>[20]</sup>	CNN	King Saud University Saudi Sign Language (KSU-SSL) dataset	87.69%
Our Methode	CNN	ASL Alphabet repository	99%

The findings of this study imply that this hand sign language recognition model can be a hand sign language independent learning tool with a relatively better level of recognition accuracy than previous similar studies. The strength of this study is that the proposed hand sign language recognition model can perform hand sign language recognition from the alphabet in real-time. While the limitation of this model is that the performance of model is strongly influenced by the specifications of the web camera and lighting system.

#### 4. CONCLUSION

In this study, a hand sign recognition model from letters of the alphabet using the CNN model has been successfully created, with significant results compared to previous related studies. Our contribution is

the addition of preprocessing to the background correction which results in good accuracy in the proposed model. Our future work is to add to the sign language dataset basic words in addition to the letters of the alphabet, and also to [increase the accuracy of the model](#) by adding other hyperparameters.

# Deep Learning Approach For Sign Language Recognition

## ORIGINALITY REPORT

24%

SIMILARITY INDEX

## PRIMARY SOURCES

1	<a href="https://ims.onnocenter.or.id">ims.onnocenter.or.id</a> Internet	164 words — 5%
2	<a href="https://eprints.uad.ac.id">eprints.uad.ac.id</a> Internet	157 words — 4%
3	<a href="https://www.researchgate.net">www.researchgate.net</a> Internet	60 words — 2%
4	Muneer Al-Hammadi, Mohamed A. Bencherif, Mansour Alsulaiman, Ghulam Muhammad et al. "Spatial Attention-Based 3D Graph Convolutional Neural Network for Sign Language Recognition", Sensors, 2022 Crossref	53 words — 2%
5	<a href="https://journal.uad.ac.id">journal.uad.ac.id</a> Internet	45 words — 1%
6	Razieh Rastgoo, Kourosh Kiani, Sergio Escalera. "Sign language recognition: A deep survey", Expert Systems with Applications, 2020 Crossref	41 words — 1%
7	<a href="https://link.springer.com">link.springer.com</a> Internet	31 words — 1%
8	"Applied Computer Sciences in Engineering", Springer Science and Business Media LLC, 2021 Crossref	29 words — 1%

- 
- 9 Razieh Rastgoo, Kourosh Kiani, Sergio Escalera. "Multi-Modal Deep Hand Sign Language Recognition in Still Images Using Restricted Boltzmann Machine", *Entropy*, 2018  
29 words — 1%  
Crossref
- 
- 10 Kirti Aggarwal, Anuja Arora. "Chapter 8 Hand Gesture Recognition for Real-Time Game Play Using Background Elimination and Deep Convolution Neural Network", Springer Science and Business Media LLC, 2022  
27 words — 1%  
Crossref
- 
- 11 "Machine Learning and Autonomous Systems", Springer Science and Business Media LLC, 2022  
22 words — 1%  
Crossref
- 
- 12 [www.chordkuncigitar.com](http://www.chordkuncigitar.com)  
17 words — < 1%  
Internet
- 
- 13 Sonam Chhikara, Rajeev Kumar. "Information theoretic steganalysis of processed image LSB steganography", *Multimedia Tools and Applications*, 2022  
16 words — < 1%  
Crossref
- 
- 14 E. Rajalakshmi, R. Elakkiya, Alexey L. Prikhodko, M. G. Grif et al. "Static and Dynamic Isolated Indian and Russian Sign Language Recognition with Spatial and Temporal Feature Detection Using Hybrid Neural Network", *ACM Transactions on Asian and Low-Resource Language Information Processing*, 2022  
15 words — < 1%  
Crossref
- 
- 15 Sakshi Sharma, Sukhwinder Singh. "Vision-based sign language recognition system: A Comprehensive Review", 2020 International Conference on Inventive Computation Technologies (ICICT), 2020  
14 words — < 1%  
Crossref

- 
- 16 Lecture Notes in Computer Science, 2007. 11 words — < 1%  
Crossref
- 
- 17 S Putrawangsa, U Hasanah. "Mathematics Education in Digital Era: Utilizing Spatialized Instrumentation in Digital Learning Tools to Promote Conceptual Understanding", Journal of Physics: Conference Series, 2020 10 words — < 1%  
Crossref
- 
- 18 [www.ijrte.org](http://www.ijrte.org) 10 words — < 1%  
Internet
- 
- 19 "Micro-Electronics and Telecommunication Engineering", Springer Science and Business Media LLC, 2021 8 words — < 1%  
Crossref
- 
- 20 "Telematics and Computing", Springer Science and Business Media LLC, 2019 8 words — < 1%  
Crossref
- 
- 21 Abu Saleh Musa Miah, Jungpil Shin, Md Al Mehedi Hasan, Md Abdur Rahim. "BenSignNet: Bengali Sign Language Alphabet Recognition Using Concatenated Segmentation and Convolutional Neural Network", Applied Sciences, 2022 8 words — < 1%  
Crossref
- 
- 22 Arturs Lavrenovs, Roman Graf, Kimmo Heinaaro. "Towards Classifying Devices on the Internet Using Artificial Intelligence", 2020 12th International Conference on Cyber Conflict (CyCon), 2020 8 words — < 1%  
Crossref
- 
- 23 I NYM Yoga Saputra, Siti Saadah, Prasti Eko Yunanto. "Analysis of Random Forest, Multiple 8 words — < 1%



Regression, and Backpropagation Methods in Predicting Apartment Price Index in Indonesia", Jurnal Ilmiah Teknik Elektro Komputer dan Informatika, 2021

Crossref

---

24 Lecture Notes in Computer Science, 2010. 8 words — < 1%

Crossref

---

25 mdpi-res.com 8 words — < 1%

Internet

---

26 Visual Analysis of Humans, 2011. 7 words — < 1%

Crossref

---

27 Abdul Mannan, Ahmed Abbasi, Abdul Rehman Javed, Anam Ahsan, Thippa Reddy Gadekallu, Qin Xin. "Hypertuned Deep Convolutional Neural Network for Sign Language Recognition", Computational Intelligence and Neuroscience, 2022 6 words — < 1%

Crossref

---

28 Chengcheng Wei, Wengang Zhou, Junfu Pu, Houqiang Li. "Deep Grammatical Multi-classifier for Continuous Sign Language Recognition", 2019 IEEE Fifth International Conference on Multimedia Big Data (BigMM), 2019 6 words — < 1%

Crossref

---

29 Kartikasari Kusuma Agustiningsih, Ema Utami, Omar Muhammad Altoumi Alsyabani. "Sentiment Analysis and Topic Modelling of The COVID-19 Vaccine in Indonesia on Twitter Social Media Using Word Embedding", Jurnal Ilmiah Teknik Elektro Komputer dan Informatika, 2022 6 words — < 1%

Crossref

---

30 Razieh Rastgoo, Kourosh Kiani, Sergio Escalera. "Hand sign language recognition using multi-view hand skeleton", Expert Systems with Applications, 2020 6 words — < 1%  
Crossref

---

EXCLUDE QUOTES OFF  
EXCLUDE BIBLIOGRAPHY OFF

EXCLUDE SOURCES OFF  
EXCLUDE MATCHES OFF