

Deep convolutional neural network for hand sign language recognition using model E

Yohanssen Pratama, Ester Marbun, Yonatan Parapat, Anastasya Manullang

Faculty of Informatics and Electrical Engineering, Institut Teknologi Del, Indonesia

Article Info

Article history:

Received Nov 3, 2019

Revised Feb 3, 2020

Accepted Mar 15, 2020

Keywords:

Convolutional neural network

Gesture

Hand sign

Hand recognition

Image processing

ABSTRACT

An image processing system that based computer vision has received many attentions from science and technology expert. Research on image processing is needed in the development of human-computer interactions such as hand recognition or gesture recognition for people with hearing impairments and deaf people. In this research we try to collect the hand gesture data and used a simple deep neural network architecture that we called model E to recognize the actual hand gestured. The dataset that we used is collected from kaggle.com and in the form of ASL (American Sign Language) datasets. We doing accuracy comparison with another existing model such as AlexNet to see how robust our model. We find that by adjusting kernel size and number of epoch for each model also give a different result. After comparing with AlexNet model we find that our model E is perform better with 96.82% accuracy.

This is an open access article under the [CC BY-SA](https://creativecommons.org/licenses/by-sa/4.0/) license.



Corresponding Author:

Yohanssen Pratama,

Faculty of Informatics and Electrical Engineering,

Institut Teknologi Del,

Jl. Sisingamangaraja, Sitoluama, Toba Samosir, Sumatera Utara, Indonesia.

Email: yohanssen.pratama@del.ac.id

1. INTRODUCTION

Based on the disability data from the Ministry of Social Affairs in 2012, they figures that in 33 provinces of Indonesia, 223,655 people were deaf (around 10.52%), 151,371 were speech impaired (around 7.12%), and were deaf and speech impaired (dumb deaf) as many as 73,560 (around 3.46%) [1]. If calculated, around 21.10% of Indonesia's population experiences hearing and speech problems. This is very concerning because around 21.10% of the population in Indonesia has difficulty in communicating with normal people, as well as to their fellow deaf and speech impaired people. This difficulty is caused by the inability of the other person to recognize sign language issued by speech impairments. If you want to recognize sign language from the communicator (speech impaired), then the communicator must learn and train themselves to understand the meaning of the sign language [2]. This is the reason why the research team feels the need to make a technology that can help deaf and speechless people communicate with each other as well as normal people.

At present the development of computer technology is rapid, making computing capabilities on computers also increase [3]. Machine learning technology has many benefits for people's lives in different aspects [4-6]. Machine learning technology is used to identify objects in an image, conversations into text, match news, upload products according to user interests, and select relevant results on a search [7]. Some of the algorithms/models used in machine learning technology include neural network and deep learning. Learning representation (representation learning) is a method that allows a machine to process raw data and automatically find the representation needed for detection or classification. Deep learning is a method of learning representation that allows computational models that are composed of many layers of processing to

learn data representation with many levels of abstraction [8]. This method significantly increases performance in the introduction of conversations, visual object recognition, object detection and so on [9]. One of the deep learning methods that are often used for image recognition is deep convolutional neural networks.

Hand gesture recognition technology made based on image processing and computer vision. Computer vision is a system in image processing that obtained an image from electronic cameras and similar to human vision systems where the brain processes images from the eye. To enhance raw image from camera so we can improve the pictorial information we need an image processing [10]. At present, computer vision is a topic that is widely discussed by us originating from electrical engineering, computer scientist, etc. This system is used in many ways, such as checking object size, food quality, detecting faces automatically, recognizing humans through iris, etc. The technology that will be developed by the research team to help speech impaired, deaf and normal people to communicate is a technology that applies computer vision, namely hand recognition using the CNN algorithm. The purpose of hand recognition is to provide natural interactions between humans and computers that can process data and then process it into information. Of course, this is very useful for the speech impaired people and they companion because communicants can understand sign language delivered by communicators (speech impaired). Unlike the hearing aids that have been circulating in the market such as SIGNLY (gloves for detecting movement and finger positioning), Kinect, Albab, etc., the technology that will be developed by the research team is technology that will change sign language from speech impaired people becomes a sentence that can be understood by opponents of communication.

2. RESEARCH METHOD

Convolutional neural network (CNN) is an approach that builds the invariance properties into the structure of a neural network and it was based a neocognitron model which is also hierarchical and multilayered structured [11-13]. CNN have succeeded in the problem of image recognition and classification, and have been successfully implemented for the introduction of human body movements in recent years. In particular, experiments have been carried out in the field of introduction of sign language using CNN, with the condition of inputs that are sensitive to background change. By using a camera that can detect the depth and contours, the process of recognizing sign language is made easier through the development of depth characteristics and movement of profiles. For each sign language introduction, the use of sensing technology is rapidly gaining popularity, and other tools have been incorporated into processes that have proven successful. Developments such as specially designed color gloves have been used for hand recognition processes and make feature extraction steps more efficient by making certain gesture units easier to identify and classify [14, 15].

However, until now, the method of introducing automatic sign language cannot utilize sensing technology that is usually carried out daily. Previous works used camera technology that was very basic to produce simple image data sets, without information on depth or contours available, only pixels exist. Efforts to deal with the task of classifying several ASL letter movement images have been successful but using the AlexNet architecture that has been previously trained.

In this paper we proposed a simple model E architecture to detect the hand gesture and varying the filter size for the input. The Model is based on traditional convolutional neural network which being used for many applications such as house numbers digit classification and has a good performance [16]. The architecture could be seen in Figure 1, that from the first convolution layer until the fourth we modified the filter size (1x1, 3x3, 9x9, and 11x11) and also we did the experiment that using 5x5 filter sizes for all convolution layer except the fourth layer that using 4x4 filter size. We have model E with different filter there are: model E which have 1x1 filter size, model E with 3x3 filter size, model E with 9x9 filter size, and model E with 11x11 filter size. Also we have model E with 5x5 filter size except its fourth layer with 4x4 filter size. In this research, we will see which filter gives the best accuracy result by doing some test experiments. After that we compared our model with previous AlexNet architecture to get accuracy. Estimating the accuracy of our model is central for this research, so we can measure our system performance [17].

In this research we use a collection of SIBI (Indonesian Language Signal System) dataset that taken from online datasets sourced kaggle.com and in the form of ASL (American Sign Language) datasets which consisting of 29 objects, 26 letters (a-z), nothing, delete, and space [18]. This is because SIBI has adopted ASL to become a standard sign language and standard in Indonesia. In carrying out this research, we conducted several steps as follows:

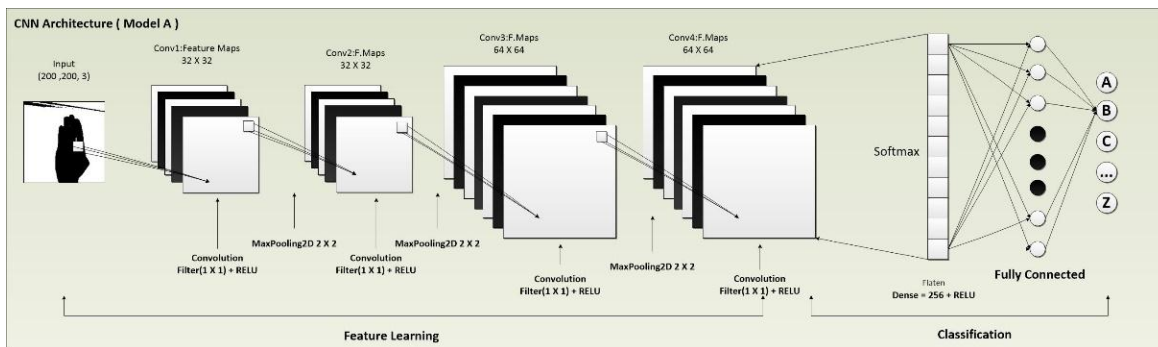


Figure 1. Model E architecture with 1x1 filter size

2.1. Collecting a dataset

The SIBI dataset, shown in Figure 2(a) that has been obtained divided into two categories with a ratio of 80:20 using holdout method in cross validation. The data is consist of 80% training dataset and 20% testing dataset. Furthermore, the data were analyzed for physical forms such as file size, pixel number, image clarity (noisy) to facilitate the training process and data validation test. After that in Figure 2(b), the colored dataset (RGB/red green blue) undergoes a color transformation stage to grayscale (gray) and then undergoes transforms again to black and white (binary image) to separate between background and foreground image using a background subtraction [19, 20]. Here is the SIBI dataset that we used:

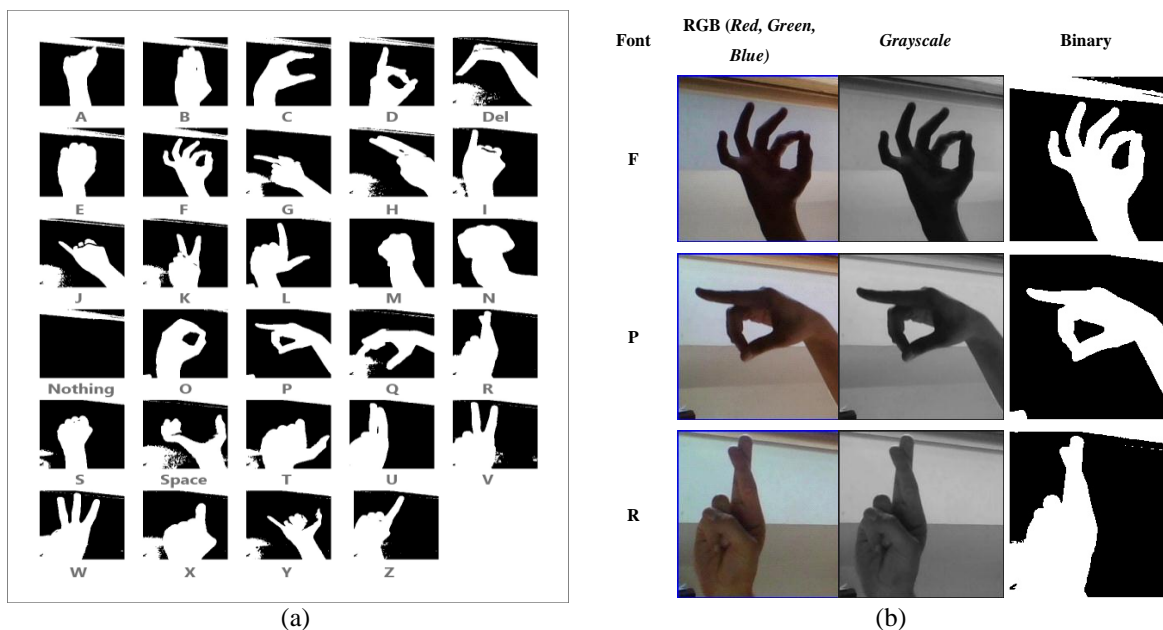


Figure 2. Dataset, (a) SIBI dataset, (b) Conversion of RGB images to binary images

2.2. Image preprocessing

In this stage we try to convert images from RGB to binary color space. Image preprocessing is done to improve the computer accuracy when recognizing images and the time efficiency to recognize images. In the image preprocessing stage, several steps are carried out such as:

- Segmentation: good segmentation process leads to perfect feature extraction process and the later play an important role in a successful recognition process [21, 22]. So for segmentation we use thresholding method to find the pattern and the background image. The value for the background image is 0 and foreground (hand) image is 1 as shown in Figure 3.

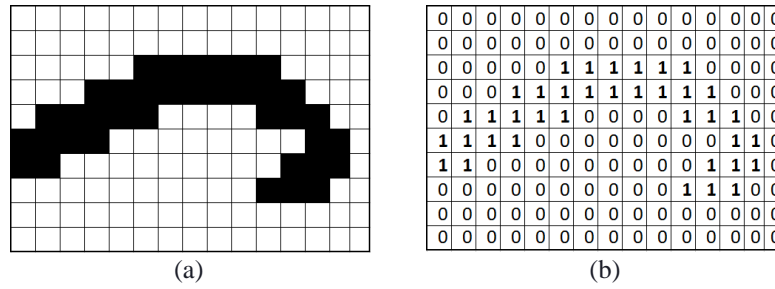


Figure 3. Thresholding process, (a) Thresholding output, (b) Output binary representation

- Resizing: resizing to help image processing performance where the image size is adjusted to the standard height and width of the source image input. Each hand dataset could have a different resolution, so to maintain the input size we need to standardize the image resolution for 100x100 pixels.
- Rescaling: re-scaling to prevent image distortion.
- Noise removal: use an efficient method of removing noise from image, so we will get the better analysis result from the system. In this research we use BM3D filter [23].

2.3. Modeling and training

At this stage, we built a simple model called Model E and the processing is done by using NVIDIA Geforce 930mx software that is operating under the Windows 10 Pro operating system. This model consists of 4 convolutional layers, 4 max pooling 2D pieces, and 1 fully connected layer. We use 4 max pooling because it can greatly improve the statistical efficiency of the network [24]. The training process uses a 5x5 matrix filter on the first layer to the third layer and uses a 4x4 filter in the fourth layer, shown in Figure 4. Furthermore, the number of epochs used is 150 with a step count of 1000 and this is done to produce high data accuracy. We take conventional approach to look for similar problems and deep learning architectures, which have already been shown to work, because there is no general answer to determine hyperparameter such as kernel size, output maps, and CNN layers. Then a suitable architecture can be developed by experimentation. The model E is a simple model of conventional CNN that improve based on experiment. If we compare a smaller filter size to larger filter size, the larger one extracted quite generic features that spread across the image and capture the basic component of image. Model E use larger filter size from 1st to 3rd layer because the hand sign does not have a complex features, although the amount information that extracted is lesser but it will use lesser memory for backpropagation. In the 4th layer we used smaller filter because we will have a better weight sharing.

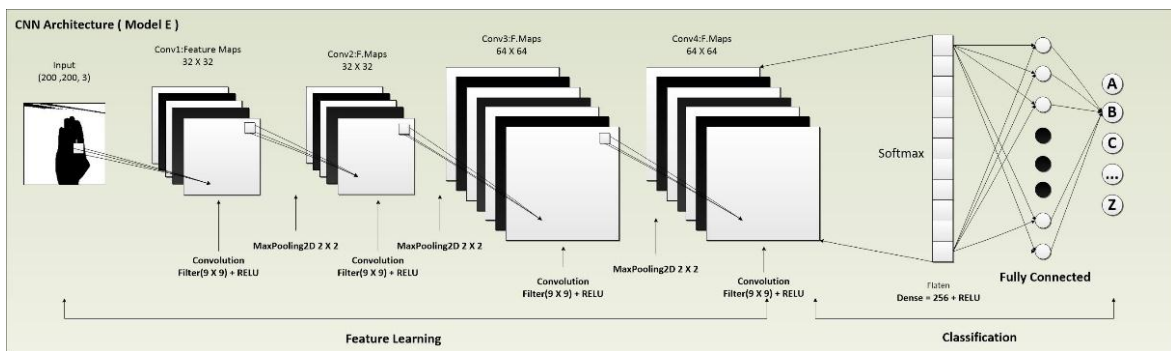


Figure 4. Model E architecture with 4x4 filter size in 4th convolutional layer

2.4. Testing

After the model has been built and the data is already passed the training phase, then the next step is testing. At this stage we tried to classifying test images using 2 models with different filters. A common problem that found during training neural networks is the choice of number training epochs that will used. Too many epochs can lead to overfitting of the training dataset, whereas too few may result in an underfit model [25]. So we try various number of epochs and see when the model performance stops improving. We choose the number of epoch based on the best model performance.

3. RESULTS AND DISCUSSION

The results of this research will consist of 3 main parts, there are data preparation stages, modeling stages, and testing stages.

3.1. Data preparation

The dataset is consist of 3000x29 imagery of hand gesture language. The image size is converted to 100x100 pixels with the binary format and stored in .png format.

3.2. Modeling

The experiment is conducted by using model E which uses several types of filter sizes. Model E implement the stochastic gradient descent (SGD) method with learning rate=0.01 and 50 epochs. Then the highest result of the fifth filters will be compared to AlexNet model. The graph of the training process can be seen in the Figure 5. From Figure 5 within the five filters, filter with combination 4x4 and 5x5 kernel size has the best accuracy when the epoch reaches 50.

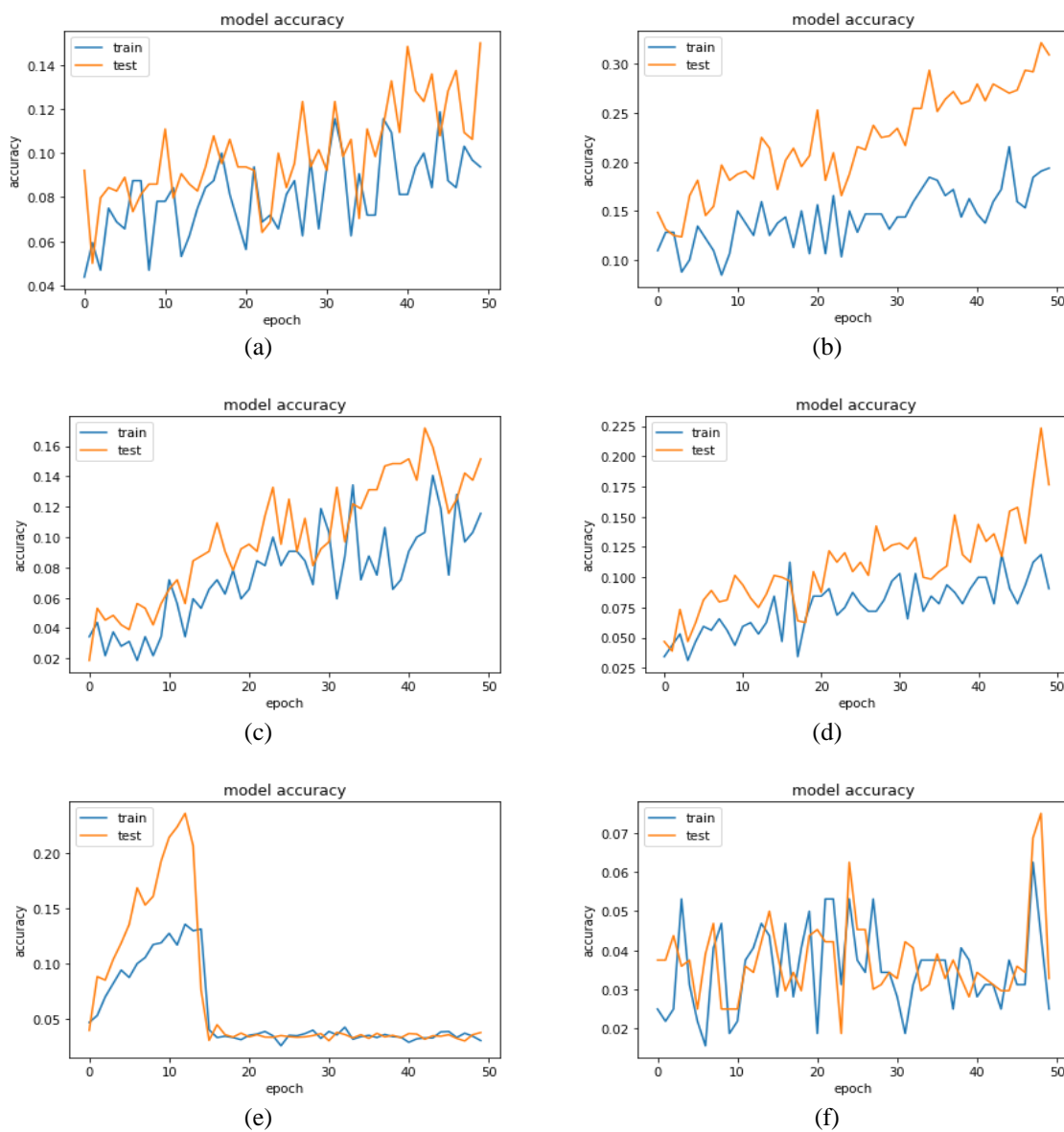


Figure 5. Experiment results based on the dataset type, (a) Model E with filter 1x1, (b) Model E with filter combination 4x4 and 5x5, (c) Model E with filter 3x3, (d) Model E with filter 11x11, (e) Model E with filter 9x9, (f) Model E AlexNet with filter 9x9

3.3. Testing

In our experiment, we have succeeded in classifying images that taken from the dataset. The image data was taken from ASL (American Sign Language) hand motion images.

3.3.1. Dataset test

The results of our experiment can be seen in Table 1 which is the experimental results based on the different dataset type. In Table 1 can be seen that there are three types of color spaces that are compared, there are RGB (red, green, blue), grayscale (gray), and binary (black and white). In this comparison, the parameters that we used has the same value, such as the kernels/filters having a matrix of size 1x1, epoch numbers are 100, and the total steps are 1000. If we sorted by the accuracy of the training data, the highest accuracy is in the RGB type with values 0.8759 and the lowest accuracy is in the Binary type which is 0.6776. However, at this initial comparison stage, we only take the lowest time/duration during the computational process to produce fast output without considering the accuracy of each type of dataset.

Table 1. Experiment results based on the dataset type

Dataset type	Computing time duration (hour)	Kernel	Epoch	Step	Accuracy	Loss
Grayscale	49.52	1x1	100	1000	0.8538	0.3961
Binary	30.49	1x1	100	1000	0.6776	0.9380
RGB	42.84	1x1	100	1000	0.8759	0.4316

In this comparison process, the kernel was intentionally implemented following the SqueezeNet model measuring 1x1 to speed up the computing [26]. After the training phase is done, it turns out that the type of dataset that has the fastest computing duration is the “binary type” with an execution time for 30.49 hours and the duration of computing is longer for the grayscale type which is for 49.52 hours. Based on the results that already obtained, we decided to use binary data types to speed up the process of training and testing of sign language data.

3.3.2. Filter size test

Next, the dataset in the form of binary is compared again to determine the most effective filters to apply in the process of recognizing sign language. In Table 2 we could see an accuracy comparison based on filter size. From Table 2 can be seen that the number of the epoch, step, and convolutional layers that we used has to be in the same amounts. The epochs numbers are 50, steps are 10, and convolutional layers are 4. This experiment was done to test the accuracy of the model with a different number of filters. The difference in filter size can affect the accuracy of training data and testing data. This is proofed by the difference in accuracy that we acquired. The highest level of accuracy is in the 5x5 filter combined with 4x4 filters, which are equal to 0.19, while the lowest level of accuracy is in the 11x11 filter, which is equal to 0.09. Based on these results, we decided to use a 5x5 filter combined with a 4x4 matrix.

Table 2. Experiment results based on filter size

Conv. layer	Kernel	Epoch	Learning rate	Step	Accuracy	Loss
4	1x1	50	0.01	10	0.0938	3.2911
4	5x5 (3 layer) +4x4	50	0.01	10	0.1938	2.9334
4	3x3	50	0.01	10	0.1156	3.2564
4	11x11	50	0.01	10	0.0906	3.2578
4	9x9	50	0.01	10	0.0305	3.3674

3.3.3. Epoch size test

After do some experiment for dataset type and filter size, we also conducted a comparison for the number of epochs to obtain the highest accuracy from the training and testing data. The results of this experiment can be seen in Table 3.

Table 3. Experiment results based on epoch amount

Conv. layer	Kernel	Epoch	Learning rate	Step	Accuracy	Loss
4	5x5 (3 layer) +4x4	50	0.01	150	0.7388	0.7651
4	5x5 (3 layer) +4x4	100	0.01	150	0.8596	0.4024
4	5x5 (3 layer) +4x4	150	0.01	150	0.9379	0.1777
4	5x5 (3 layer) +4x4	200	0.01	150	0.8360	0.4860
4	5x5 (3 layer) +4x4	250	0.01	150	0.8981	0.2879

In Table 3, it can be seen that the number of convolutional layers, steps, and filters are in the same size, the convolution consists of four layers, 1000 total steps, and 5x5 filters size for the first into third convolutional layer and 4x4 sized matrices in the fourth convolutional layer. Based on these results, the highest accuracy is at the 150th epoch with an accuracy rate of 0.9379 and the lowest accuracy is at the 50th epoch with an accuracy rate of 0.7388. The higher the number of epochs that are applied during the training and testing process, the higher the level of accuracy of the data, so that we decide to use epoch which amounts to 150 to be used as parameters when performing the recognition process.

3.4. Discussion of differences in 2 models

In determining whether model E is a suitable and relatively efficient model for performing recognition processes on SIBI, we compare model E with a pre-existing model and are generally used for recognition, the AlexNet model. The AlexNet model is often used in the process of hand gesture recognition for sign language. This is due to several advantages of AlexNet such as simple, using the CNN architecture that is basic, and effective. An accuracy and loss comparison of model E with the trained Alexnet model can be seen from the following TensorBoard graph in Table 4.

Table 4. Comparison of model E and AlexNet

Aspect	Model E	Model AlexNet
Accuracy		
Loss		

Both models use the same number of epochs as 50, step 150, and learning rate 0.01, but use different numbers of convolutional layers, five layers in the AlexNet model and four layers in model E. The different layers make the AlexNet model require longer computation time to do the training process compared to model E, where AlexNet requires 1.72 hours while model E takes 0.71 hours.

Comparison of model E with the AlexNet model can be seen that the performance of model E is more efficient than the AlexNet model because in model E, accuracy tends to rise from the lower left to the lower right, while in the AlexNet model accuracy tends to increase and decrease in succession. In the AlexNet model, the accuracy of the 4th epoch increase in training reaches accuracy at 0.0531 and decreases at 0.02 at the 5th epoch, and at the 7th epoch decreases to 0.01. Then, at the 48th epoch, it increases to 0.0625, drops back to the 49th epoch to 0.0438, and drops again to 0.0250 at the 50th epoch.

In model E, the training accuracy is 19.38% and validation accuracy is 30.94%. Then, training loss is 2.9334 and the validation loss is 2.4456. In the AlexNet model, the training accuracy is 2.50% and validation accuracy is 3.28%. Then, the training loss is 3.3700 and the validation loss is 3.3669. Based on these data, the model E has better prediction compare to AlexNet Model because the total loss in model E is smaller than the AlexNet model.

Based on the results of a comparison between model E and AlexNet, we decided to use model E in classifying and predicting sign language (SIBI). Then, we optimizes the model again by increasing the number of epochs to 100, the stepper epoch being 960, and the validation step being 2600, resulting in

a 96.83% accuracy. This quite high result indicates that the model has a good learning ability so that it has a good level of prediction. The model E that has been improved is then used by us to recognize sign language, and then combine the results into a sentence.

4. CONCLUSION

The model E is the model that based on simple CNN architecture and uses a number variation of filter size. The model E is proven to have a better performance with larger filter size from 1st to 3rd layer because the hand sign does not have a complex features, although the amount information that extracted is lesser but it will use lesser memory for backpropagation. In the 4th layer the smaller filter is being used here because it will have a better weight sharing. For the epoch variable already tried a various number of epochs and seen when the model performance stops improving, for model E the number epoch that give the best performance is 150. As the result we suggest to use 150 number of epoch for hand gesture experiment with model E. Moreover, CNN can be managed to recognize or classify hand gestures especially in the Indonesian Language Signal System with accuracy of 96.82% by using model E. With the same number of epoch for AlexNet model we found our model E was perform slightly better with delta 16.88%. For further research, we will try a model that can recognize sign language by observing static to dynamic gestures. So the gesture that observed from the camera could directly recognize into a language.

ACKNOWLEDGEMENTS

This work was supported in part by Institut Teknologi Del.

REFERENCES

- [1] Ministry of Social Affairs Republic of Indonesia, "Ministry of social affairs in social welfare development numbers," in Bahasa "Kementrian sosial dalam angka pembangunan kesejahteraan sosial," Jakarta, 2012.
- [2] V. Priyadharshni, M. S. Anand, and N. M. Kumar, "Hand gesture recognition system using hybrid technology for hard of hearing community," *Int. J. of Eng. Mathematics & Computer Science*, vol. 1, no. 2, pp. 50-63, 2013.
- [3] A. Chaudhary, J. L. Raheja, K. Das, and S. Raheja, "Intelligent approaches to interact with machines using hand gesture recognition in natural way: A survey," *Int. J. of Comp. Sci. & Eng. Survey*, vol. 2, no. 1, pp. 122-133, 2011.
- [4] P. Yohanssen, I. G. B. B. Nugraha, and E. T. Samosir, "Vehicle counting and classification for traffic data acquisition," *Jurnal Teknologi*, vol.78, no. 6-3, pp. 77-82, 2016.
- [5] P. Yohanssen and P. R. Puspoko, "The addition symptoms parameter on sentiment analysis to measure public health concerns," *TELKOMNIKA Telecommunication Computing Electronics and Control*, vol.15, no. 3, pp. 1301-1309, 2017.
- [6] P. Yohanssen and S.S. Hadi, "Cassava quality classification for tapioca flour ingredients by using ID3 algorithm," *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 9, no. 3, pp. 799-805, 2018.
- [7] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp.436-444, 2015.
- [8] M. Tanaka and M. Okutomi, "A novel inference of a restricted boltzman machine," *22nd International Conference on Pattern Recognition*, pp.1526-1531, 2014.
- [9] A. Deshpande, "A beginner's guide to understanding convolutional neural networks," *Retrieved March*, vol. 31, 2016.
- [10] B. Chitradevi and P. Srimathi, "An overview on image processing techniques," *International Journal of Innovative Research in Computer*, vol. 2, no. 11, pp. 6466-6472, 2014.
- [11] C. M. Bishop, "Pattern recognition and machine learning," *Springer*, pp. 227-229, 2006.
- [12] L. V. Fausett, "Fundamental of neural networks: Architecture, algorithm, and application," Prentice-Hall Inc., 1994.
- [13] K. Fukushima, "Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position," *Biol. Cybernetics*, vol. 36, pp. 192-202, 1980.
- [14] J. F. Paul, D. Seth, C. Paul, and J. G. Dastidar, "Hand gesture recognition library," *Int. J. of Science and Applied Information Technology*, vol. 3, no. 2, pp. 44-50, 2014.
- [15] Y. Xu and Y. Dai, "Review of hand gesture recognition study and application," *Contemporary Engineering Science*, vol. 10, no. 8, pp. 375-384, 2017.
- [16] P. Sernamet, S. Chintala, and Y. LeCun, "Convolutional neural networks applied to house numbers digit classification," *Proc. Of the 21st International Conference on Pattern Recognition*, pp. 3288-3291, 2012.
- [17] E. Platanios, H. Poon, T. M. Mitchell, and E. J. Horvitz, "Estimating accuracy from unlabeled data: A probabilistic approach," *31st Conference on Neural Processing Systems, (NIPS 2017)*, pp. 4361-4370, 2017.
- [18] J. Davis and M. Shah., "Recognizing hand gestures," *European Conference on Computer Vision*, pp. 331-340, 1994.
- [19] R. Suguna and P. S. Neethu, "Hand gesture recognition using shape features," *Int. J. of Pure and Applied Mathematics*, vol. 117, no. 8, pp. 51-54, 2017.
- [20] S. D. Badgujar, G. Talukdar, O. Gondhalekar, and S. Y. Kulkarni, "Hand gesture recognition system," *Int. J. of Scientific and Research Publications*, vol. 4, no. 2, pp. 1-5, 2014.
- [21] R. Z. Khan and N. A. Ibraheem, "Hand gesture recognition: A literature review," *Int. J. of Artificial Intelligence & Application*, vol. 3, no. 4, pp. 161-174, 2012.

- [22] X. Li, "Gesture recognition based on fuzzy c-means clustering algorithm," Department of Computer Science, University of Tennessee Knoxville, 2003.
- [23] M. Verma and D. J. Ali, "A comparative study of various types of image noise and efficient noise removal techniques," *International Journal of Advanced Research in Computer Science and Software Engineering*, vol. 3, no.10, pp. 617-622, 2013.
- [24] I. Goodfellow, Y. Bengio, and A. Courville, "Deep learning," MIT Press, pp. 204-210, 2016.
- [25] F. Chollet, "Deep learning with python," Manning Publications, 2018.
- [26] W. Rawat and Z. Wang, "Deep convolutional neural networks for image classification: A comprehensive review," *Neural Computation*, vol. 29, no. 9, pp. 2532-2449, 2017.

BIOGRAPHIES OF AUTHORS



Yohanssen Pratama, S.Si., M.T.

Current Faculty Members & Researcher in Del Institute of Technology. 4+ years experience specializing in back-end/infrastructure, analytical tools development and computer programming. Teach academic and vocational subjects to undergraduate also pursue my own research to contribute to the wider research activities of my department.



Ester Enjela Marbun, S.Tr.Kom

Currently, work as a Software Engineer at PT. Citra Global Dinamika (Esecurity), Batam Center. Have experience at image processing using neural network, automation testing, web & desktop developing, and digital designing.



Yonatan Vikario Resha Parapat, S.Tr.Kom

Currently, work as a IT Consultant Analyst at Mitra Integrasi Informatika(subsidiary of PT. Metrodata Electronics Tbk.) APL Tower 37th Floor Suite 3 Jl. Letjen S. Parman Kav. 28, RT.13/RW.7, Jelambar Baru, Grogol petamburan, RT.10, RT.12/RW.6, Tj. Duren Sel., Kec. Grogol petamburan, West Jakarta. Have experience at image processing using neural network, Networking, Web PenTest, and Mobile developing.



Anastasya Pehulisa Manullang, S.Tr.Kom

Currently, work as a Product Owner at PT. Moonlay Technologies, Sudirman Central Business District, South Jakarta. Have experience at image processing using neural network, automation testing, web & desktop developing, digital designing and system analysis.