

Thai Sign Language Recognition: an Application of Deep Neural Network

Anusorn Chaikaew
Computer Science, SCIT, Chiang
Rai Rajabhat University
Chiang Rai, Thailand
anusorn.c@rru.ac.th

Kritsana Somkuan
Information Technology, SCIT
Chiang Rai Rajabhat University
Chiang Rai, Thailand
kritsana@rru.ac.th

Thidalak Yuyen
Computer Science, SCIT Chiang Rai
Rajabhat University
Chiang Rai, Thailand
thidaluk02@gmail.com

Abstract—Thai Sign Language (TSL) is the national sign language for Thai deaf people or hearing impaired in Thailand. These people with disabilities can use sign language to communicate with people with disabilities but face obstacles in communicating daily with ordinary people. Technology should play a key role in helping disadvantaged people achieve a better quality of life. This research aims to find ways to create Thai sign language recognition applications and be developed for real-time sign language translation in the next step. We purpose a simple approach with a MediaPipe framework that helps to extract the hand landmark from video on the preprocessing step and use that landmark to build the model for recognition hand gestures with various Recurrent neural networks (RNN). The result showed that the model builds with LSTM, BLSTM and GRU has an accuracy greater than 90 percent. This approach can produce an accurate close to the traditional approach.

Keywords—Computer Vision, Machine Learning, Neural Network, Thai Sign Language (TSL)

I. INTRODUCTION

A report on the situation of persons with disabilities in Thailand in September 2020 [1] found that there were 2,058,082 people or 3.09% of the entire population. Of these, 388,233 were hearing impaired or interpretive. Human-computer interaction (HCI) could help so many people at all levels of society to upgrade a better life, especially disabled people.

Sign Language is a subset of hand gestures that is part of the human modality to communication among deaf people or the hearing impaired. Thai Sign Language (TSL) was developed from American Sign Language (ASL) but different in gesture detail; for example, in Fig.1, an example of Thai sign language communicates the word love, sick and number three, respectively. TSL is a standard sign language in Thailand, but only some ordinary people understand it and use it to communicate with deaf people that maybe sometimes hearing impaired can not communicate in daily life with ordinary people. Sign Language Recognition (SLR) is the start point to develop sign language translation, which is interdisciplinary research area topics such as computer vision, natural language, artificial neural network processing, etc.

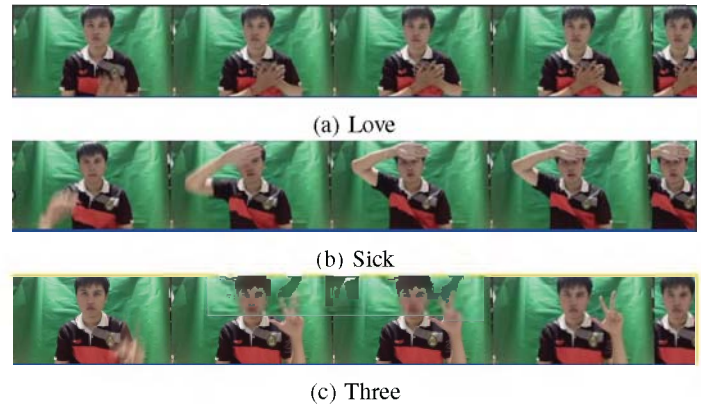


Fig. 1. Example of a Thai Sign Language (TSL).

To allow the machine to recognize the sign language that uses hand movement, we need to prepare the dataset and design the training step with a suitable algorithm. Therefore, we decided to create our own data set to train the model, aiming to build the model that recognizes at least 100 gestures/word capability covering the vocabulary used in daily life in different categories such as Greetings, Numbers, Days, Times and Seasons, etc.

The most state-of-the-art algorithm recognizes body movement using a neural network to extract feature and trained model by decomposing the problem sub unit [2] or using cascade network [5]. This has led to our research aimed at using neural networks as well.

The organization of this paper is as follows. Section II contains some related work. We discuss some algorithm which is adopted or related in our work. Our approach is explained in Section III and evaluate the result in Section IV. Finally we conclude our contribution and discussing future work in Section V.

II. RELATED WORK

One of the most used algorithms today and found to be effective for sign language recognition (SLR) is the deep neural network (DNN) implementation. In this paper, we adopted various algorithms under the DNN area to fit into our criteria. In particular, the RNN (Recurrent Neural Network) algorithm,

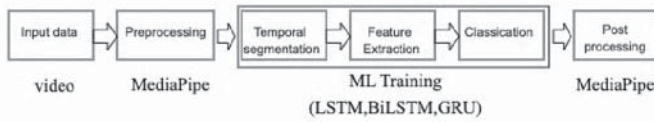


Fig. 2. Our approach.

since this research builds the modeling from extracted the key points on both palms when gesture movements in sign language.

A. Convolutional Neural Network

Convolutional neural networks (CNN), also called ConvNets, were first introduced in the 1980s by Yann LeCun is a class of deep neural networks most commonly applied to analyzing visual imagery or video. Sign language research area adopted CNN such as [2], [7], [8], [10] and so on for feature extraction and classification. Especially [10] focus on TSL like our paper showed accuracy up to 97.7 percent but based on the Kinect device for data acquisition step, unlike our research, the process of collecting the dataset does not require any special equipment.

B. Recurrent Neural Network

Recurrent Neural Network(RNN) is a class of artificial neural networks that mostly focus on temporal sequence data such as NLP, text, voice, video, or time-series data. So many research paper in machine translation using RNN or improved algorithm from RNN such as LSTM or GRU [4], and in this paper, we focus on using RNN since the extracted hand key point's position is saved in a text file format (CSV) corresponding to the reference position, it is more suitable for choosing RNN algorithms than CNN.

C. MediaPipe framework

MediaPipe [6] is a framework for building a pipeline mainly for rapid prototype base on C++. In this paper, we use MediaPipe for extracted hand key point that MediaPipe already provide by customize some source code and write the absolute position of hand landmark to text file base on guideline [11], [12].

III. OUR APPROACH

Sign Language mainly focus on two ways the feature extraction from images or extract hand key point and concatenated their location [9]. Our approach based on pipeline as Fig.2. It can be separated into four large steps: creating a video dataset. Preparing data to train the model; Model training and using models to create applications for predict Thai sign language; details are as follows.

A. Input Data

The very first step is to prepare a dataset to use for model training. We shoot 100 videos per word. However, in this experiment, we recorded five videos of five gestures, a total of



(a) 21 points on each hand detected by MediaPipe Framework

```

1.075075 0.983447 0.635131 0.808252 0.395826 0.80169 0.572331 0.81088 0.551083 0.844567
1.888058 0.473008 0.981095 0.58708 0.874218 0.545845 0.879116 0.527843 0.806792 0.512226
1.559852 0.873458 0.577882 0.814722 0.356732 0.803278 0.555497 0.82553 0.515284 0.826792
1.808856 0.551443 0.871188 0.50714 0.876355 0.551288 0.877202 0.808802 0.808833 0.584866
1.576872 0.80166 0.532418 0.507855 0.160172 0.945481 0.577284 0.802551 0.572536 0.867436
1.541895 0.570343 0.925518 0.524551 0.976239 0.551688 0.906191 0.559224 0.571718 0.557998
1.55453 0.907741 0.561852 0.808627 0.562536 0.805158 0.525932 0.903439 0.555939 0.934603
1.617825 0.540275 0.971571 0.540283 0.95595 0.54028 0.93804 0.403210 0.376631 0.453862
1.54404 0.870346 0.516546 0.942335 0.553646 0.895037 0.558519 0.851663 0.554508 0.544781
1.532220 0.555080 0.915587 0.463101 0.949175 0.476977 0.80442 0.486431 0.810773 0.515353
1.659323 0.871734 0.558233 0.079042 0.575956 0.808003 0.567113 0.922761 0.559327 0.91851
1.535347 0.460722 0.00007 0.406790 0.040973 0.516784 0.023108 0.523568 0.03077 0.538445
1.629717 0.833489 0.552845 0.943692 0.593709 0.915043 0.613754 0.504837 0.628132 0.87845
1.651163 0.505063 0.832654 0.525689 0.80588 0.507159 0.907434 0.558242 0.852117 0.576607
1.571647 0.908545 0.58373 0.564959 0.406066 0.830752 0.537273 0.908561 0.563439 0.907892
1.401805 0.535136 0.911318 0.557348 0.865727 0.501808 0.83382 0.602705 0.808493 0.527024
1.507146 0.825953 0.535838 0.027576 0.558936 0.807506 0.565579 0.852756 0.566474 0.83307
1.831138 0.576184 0.707152 0.593549 0.468444 0.533858 0.899705 0.557288 0.847097 0.5848
1.558936 0.86552 0.570801 0.934715 0.375255 0.813818 0.425582 0.717214 0.407616 0.638094
1.555576 0.369167 0.513108 0.363078 0.520199 0.378009 0.552002 0.52311 0.579111 0.549204
1.552427 0.672361 0.387573 0.70677 0.409816 0.645584 0.412886 0.585389 0.404336 0.533659
1.558945 0.572297 0.364477 0.546855 0.363308 0.565176 0.305334 0.556871 0.546706 0.54633
1.468573 0.408089 0.465155 0.413227 0.302865 0.408708 0.537568 0.4001 0.511285 0.359201
1.373443 0.560073 0.352764 0.508757 0.347856 0.540182 0.356643 0.558572 0.361507 0.574947
1.414632 0.594797 0.404781 0.545851 0.401530 0.518578 0.394249 0.520554 0.333942 0.509805
1.298227 0.550244 0.346254 0.543718 0.355121 0.575587 0.360843 0.50895 0.312368 0.395559
1.0 0.0 0.0 0.0 0.370889 0.710334 0.404781 0.656977 0.418583 0.595915 0.399853 0.3
1.264215 0.533885 0.374348 0.566602 0.379481 0.59159 0.343574 0.553940 0.558172 0.55775
1.761889 0.404952 0.640322 0.418562 0.588913 0.396572 0.543615 0.38883 0.521162 0.390815
1.373736 0.574182 0.340889 0.562709 0.347683 0.549456 0.357173 0.571564 0.363134 0.591518

```

(b) After Extracted 42 key points and write to CSV files

Fig. 3. key point position detection on each palm and writing key point to text files (CSV)

500 videos, to test our approach. Each video contains 50 FPS with format H.264, and this process takes the longest time but does not rely on any additional equipment except a mobile camera.

B. Preprocessing

In this step, we customize source code from the MediaPipe framework to extracted the hand key point 3 (a) and write it to a text file (CSV) as Fig. 3 (b). At this stage, a total of 21 key locations on both palms, a total of 42 key points, will be recorded in a text file; all 42 landmarks will pull from each frame of the video. Therefore, we need each video to have the same length as possible to create an equivalent data set.

C. Model Training

At this stage, the hand's key point we extracted from the video is suitable for RNN, and we decide to use the three different variations of architecture related to RNN are LSTM (Long short-term memory), Bi-LSTM (Bidirectional LSTM), and GRU (Gated recurrent unit). LSTMs were developed to deal with the vanishing gradient problem that can be encountered when training traditional RNNs. Bi-LSTM is a compound of Long Short-Term Memory (LSTM) and Bi-directional Recurrent Networks (BiRNN). The GRU is like an LSTM with a forget gate but has fewer parameters than LSTM. Our approach in this step is that using a text file training set takes less time to train a model than the CNN approach that relies on the large image or video datasets, which can take a longer time than text file such as [10]. Moreover, this step also may be necessary to repeat the process many times by

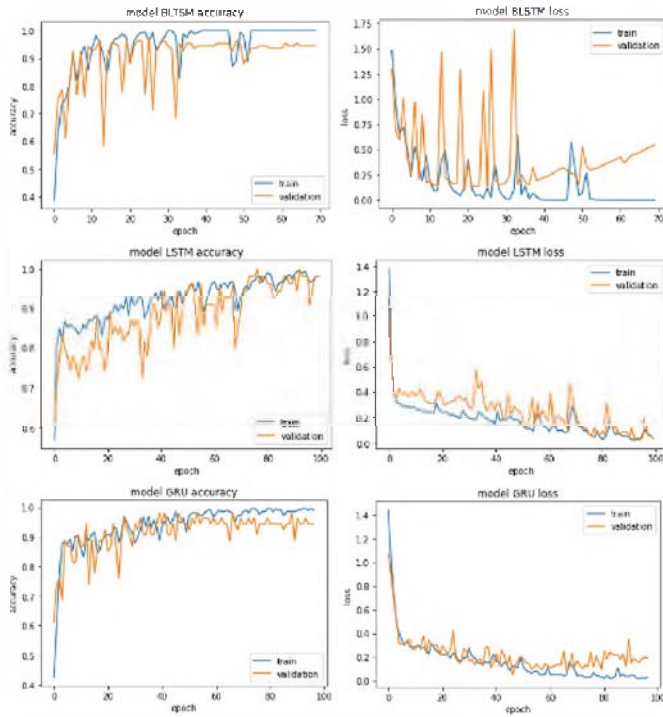


Fig. 4. Accuracy/Loss

TABLE I. Hyperband parameter for tuning each algorithm

parameter	min	max
input node	64	256
hidden layer	1	3
optimizer	'Adagrad','Adamax','Adam','RMSprop'	

adjusting the parameters to get the best accuracy. Thus text files might be easier to train model than other data set.

Additionally, In this step, we also are tuning the hyperparameter with Hyperband [3] that focuses on speeding up random search through adaptive hyper-parameter allocation the importance Hyperband parameter list on Table I. From the table, we can see that we have the default number of nodes at 64 nodes and can be increased by a maximum of 256 nodes, and there are from 1 to 3 hidden layers and 4 types of optimizer for Hyperband to choose from. Fig. 5 show each RNN algorithm's summary after tuning with Hyperband. We found that the GRU had the least number of parameters while BiLSTM had the largest number of parameters, meaning BiLSTM was the largest model while the GRU was the smallest model in this experiment.

After tuning with Hyperband the next step is to build the model with training data and the output result of each algorithm after train shown in Table II. The data set of training data has 100 videos per word divided into 60 videos for training, 30 for the validation test, and the last 10 for testing the accuracy after trained models. All algorithms can surpass accuracy more than 90 percent with train data set and surpass more than 90 percent with test data set as Table III. During the model's practice with different epoch numbers ranging from

Layer (type)	Output Shape	Param #
lstm (LSTM)	(None, 70, 128)	109056
lstm_1 (LSTM)	(None, 70, 192)	246528
lstm_2 (LSTM)	(None, 70, 64)	65792
dropout (Dropout)	(None, 70, 64)	0
lstm_3 (LSTM)	(None, 128)	98816
dense (Dense)	(None, 6)	774

Total params: 520,966
Trainable params: 520,966
Non-trainable params: 0

(a) LSTM

Layer (type)	Output Shape	Param #
bidirectional (Bidirectional)	(None, 70, 256)	218112
bidirectional_1 (Bidirectional)	(None, 70, 256)	394240
bidirectional_2 (Bidirectional)	(None, 70, 256)	394240
bidirectional_3 (Bidirectional)	(None, 128)	164352
dense_1 (Dense)	(None, 6)	774

Total params: 1,171,718
Trainable params: 1,171,718
Non-trainable params: 0

(b) BiLSTM

Layer (type)	Output Shape	Param #
gru (GRU)	(None, 70, 192)	160128
gru_1 (GRU)	(None, 70, 192)	222336
dropout_1 (Dropout)	(None, 70, 192)	0
gru_2 (GRU)	(None, 32)	21696
dense_2 (Dense)	(None, 6)	198

Total params: 404,358
Trainable params: 404,358
Non-trainable params: 0

(c) GRU

Fig. 5. Model Summary of each model after tuning with Hyperband

TABLE II. Model accuracy comparison while trained

Network	Accuracy	Loss
LSTM	0.94	0.16
BLSTM	0.97	0.39
GRU	0.93	0.20

70 to 100, we found that the graphs obtained while practicing Fig. Ref fig: acc-lost showed significant fluctuations in the BiLSTM, LSTM and slightly up and down on GRU due to the number of videos required for the practice (60 videos), it may not have been enough for BiLSTM or LSTM with parameters larger than GRU.

D. Post-processing

In this step, we developed both desktop and mobile applications, as shown in the Fig. 6 both of them have used different

TABLE III. Model Accuracy (Predicted with test data)

Model	Accuracy	Loss
LSTM	0.97	0.06
BLSTM	0.94	0.23
GRU	0.94	0.14

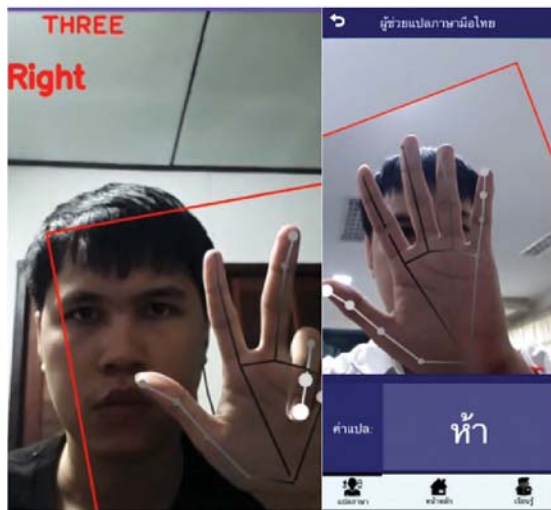


Fig. 6. Desktop Application (left) and Mobile Application (right)

models in the application development, choosing the model with the highest precision from the table. III, that is, LSTM has 97 percent accuracy for desktop. And the GRU model was chosen for mobile use because of its small size but acceptable accuracy.

IV. RESULT AND DISCUSSION

By looking at the results, we found that we could build and deploy the model in the development phase of both desktop and mobile applications. It shows that this approach can be applied to the development of both desktop and mobile sign language recognition. However, there are some limitations in binding to the MediaPipe Framework, but it may be acceptable in terms of open-source code.

V. CONCLUSION AND FUTURE WORK

In this paper, we presented Thai Sign Languages recognition that builds from the MediaPipe framework by extracting key points from hand landmarks. We use just a simple smartphone camera and train it with CPU on the laptop and even faster when training with GPU. Our experiments demonstrate how to build Thai Sign Language recognition with the MediaPipe framework and will be the next step to make a Sign Language Translation and the next future.

ACKNOWLEDGMENT

This research has been supported by the Research and Development Institute of Chiang Rai Rajabhat University

REFERENCES

- [1] Department of Empowerment of Persons with Disabilities. (2020, October 24). A report on the situation of persons with disabilities in Thailand [Online]. Available: <http://www.dep.go.th>
- [2] N. C. Camgoz, S. Hadfield, O. Koller, and R. Bowden, "SubUNets: End-to-End Hand Shape and Continuous Sign Language Recognition," in 2017 IEEE International Conference on Computer Vision (ICCV), Venice, Oct. 2017, pp. 3075–3084, doi: 10.1109/ICCV.2017.332.

- [3] L. Li, K. Jamieson, G. DeSalvo, A. Rostamizadeh, and A. Talwalkar, "Hyperband: A Novel Bandit-Based Approach to Hyperparameter Optimization," arXiv:1603.06560 [cs, stat], Jun. 2018, Accessed: Nov. 02, 2020. [Online]. Available: <http://arxiv.org/abs/1603.06560>.
- [4] K. Cho et al., "Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation," arXiv:1406.1078 [cs, stat], Sep. 2014, Accessed: Nov. 02, 2020. [Online]. Available: <http://arxiv.org/abs/1406.1078>.
- [5] D. Chen, G. Hua, F. Wen, and J. Sun, "Supervised Transformer Network for Efficient Face Detection," arXiv:1607.05477 [cs], Jul. 2016, Accessed: Nov. 02, 2020. [Online]. Available: <http://arxiv.org/abs/1607.05477>.
- [6] C. Lugaresi et al., "MediaPipe: A Framework for Building Perception Pipelines," arXiv:1906.08172 [cs], Jun. 2019, Accessed: Nov. 02, 2020. [Online]. Available: <http://arxiv.org/abs/1906.08172>.
- [7] P. Kishore, G. A. Rao, E. K. Kumar, M. T. K. Kumar, and D. A. Kumar, Selfie sign language recognition with convolutional neural networks. *International Journal of Intelligent Systems and Applications*, 10(10):63, 2018.
- [8] H. Shin, W. J. Kim, and K.-a. Jang. Korean sign language recognition based on image and convolution neural network. In *Proceedings of the 2nd International Conference on Image and Graphics Processing*, pages 52–55. ACM, 2019.
- [9] D. Li, C. R. Opazo, X. Yu, and H. Li, "Word-level Deep Sign Language Recognition from Video: A New Large-scale Dataset and Methods Comparison," arXiv:1910.11006 [cs], Jan. 2020, Accessed: Oct. 25, 2020. [Online]. Available: <http://arxiv.org/abs/1910.11006>.
- [10] N. Sriparojthikoon and J. Harnsomburana, "Thai Sign Language Recognition Using 3D Convolutional Neural Networks," in *Proceedings of the 2019 7th International Conference on Computer and Communications Management*, Bangkok Thailand, Jul. 2019, pp. 186–189, doi: 10.1145/3348445.3348452.
- [11] M. Tachionstrahl, "Tachionstrahl/SignLanguageRecognition," 27-Oct-2020. [Online]. Available: <https://github.com/Tachionstrahl/Sign-LanguageRecognition>. [Accessed: 11-Nov-2020].
- [12] A. Kim, "rabBit64/Sign-language-recognition-with-RNN-and-Mediapipe," Apr-2020. [Online]. Available: <https://github.com/rabBit64/Sign-language-recognition-with-RNN-and-Mediapipe>. [Accessed: 01-Oct-2020].