

---

# Development of a hand pose recognition system on an embedded computer using CNNs

---

**Dennis Núñez Fernández**  
 Universidad Nacional de Ingeniería  
 Lima, Peru  
 dnunezf@uni.pe

## Abstract

Demand of hand pose recognition systems are growing in the last years in technologies like human-machine interfaces. This work suggests an approach for hand pose recognition in embedded computers using hand tracking and CNNs. Results show a fast time response with an accuracy of 94.50% and low power consumption.

## 1 Introduction

Hand gesture recognition is one obvious strategy to build user-friendly interfaces between machines and users. In the near future, hand posture recognition technology would allow for the operation of machines through only series of hand postures, eliminating the need for physical contact. However, hand gesture recognition is a difficult problem because occlusions, variations of appearance, etc. Despite these difficulties, several approaches to gesture recognition on images has been proposed [9].

In recent years, convolutional neural networks (ConvNets) have become the state-of-the-art for object recognition [5]. In spite of the high potential of CNNs in object detection problems [1, 6] and image segmentation [5], only a few papers report successful results. A recent survey on hand gesture recognition [9] reports only one important work [11]. Some obstacles to wider use of CNNs are high computational costs, lack of sufficiently large datasets, as well as lack of appropriate hand detectors.

## 2 Methodology

The proposed system works with images captured from a CMOS camera and runs on embedded computers without GPU support such as the Raspberry Pi, BeagleBone, Intel Galileo among others. Therefore, the goals of the proposed system are as follows: high accuracy, fast response time and low power consumption. The system was implemented in C++ in order to obtain the best performance.

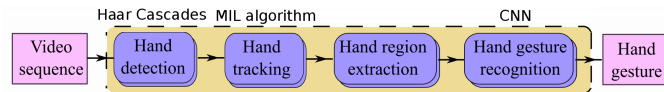


Figure 1: Diagram for the proposed system

Haar cascade classifier allows better detection for objects with static features such as balloons, faces, eyes, etc. But a hand in motion has few static features because shape changes over time. So, this classifier is not suitable to recognize a hand poses in motion. However, its deficiency could be compensated with a hand tracker based on wrist region, which features keep invariant over time. Furthermore, tracking reduces the processing time since it requires less computational resources than detection. We use the MIL (Multiple Instance Learning) tracking algorithm [2]. It avoids the drift

problem for a robust tracking and consumes less memory and computational resources than Haar cascade classifier. In addition, due skin color is a powerful feature for fast hand detection, a model in RGB-YCbCr color spaces have been constructed on the basis of a training dataset. Then, the hand region was obtained by thresholding and morphological operations. The dataset for hand gesture classification was taken from AGH University of Science and Technology [8]. It has 73,124 grayscale images of 48x48 pixels divided into 10 hand gestures. The proposed CNN takes as input a binary image of 48x48 pixels. The architecture is: C(5x5)-S(2x2)-C(3x3)-S(2x2)-FC(120)-FC(84)-FC(10), where C: Conv. layer, S: Sub sampling, FC: Full connection. We used Caffe framework [4].

### 3 Results

The performance of the proposed Convolutional Neural Network for hand poses classification was evaluated using different metrics such as confusion matrix and accuracy. Fig. 2 depicts the hand pose for each class in grayscale format. The confusion matrix of our model is shown in Fig. 3 and discloses which hand poses are misclassified. These errors happen because of similarities between the classes. Furthermore, our architecture shows an outstanding accuracy of 94.50%.

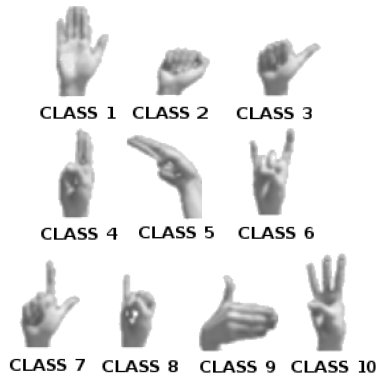


Figure 2: Hand poses

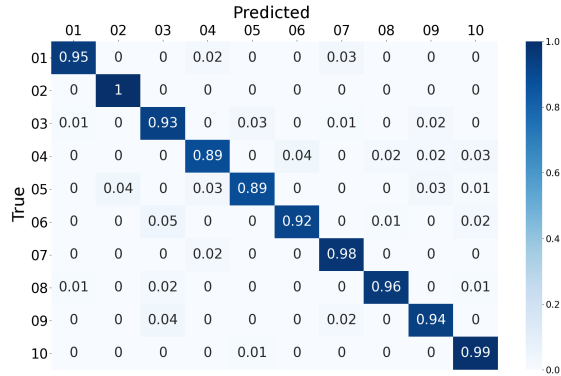


Figure 3: Confusion matrix

The implementation of the proposed recognition system on a desktop PC has no issues due to its high computational resources. However, when a recognition system is implemented on embedded computers like the Raspberry Pi 3 we have two major obstacles working against us: limited RAM memory (only 1 GB) and limited processor speed (4 ARM Cortex-A53 @ 1.2 GHz). The Table 1 shows the performance of CNNs on the Raspberry Pi 3 platform. As you can see, the proposed CNN achieves the fastest response time and the lowest power consumption.

Table 1: Response time and power consumption for different CNNs on a Raspberry Pi 3 using Caffe

Model	Proposed CNN	VGG_F [3]	NiN [7]	AlexNet [5]	GoogLeNet [10]
Layers	9	13	16	11	27
Power (W.)	0.690	0.760	0.840	0.750	0.790
Time (s.)	0.351	0.857	0.553	1.803	1.175

### 4 Conclusions

In this work we demonstrated that our system is capable to recognize 10 hand gestures with an accuracy of 94.50% on images captured from a single RGB camera, and using low power consumption, which is about 0.690 W. In addition, we show that the average time to process each image on the Raspberry Pi 3 is about 351.2 ms. The trained CNN models (Caffe models) as well as a version of the source code are fully available at: <https://github.com/dennishnf/cnn-hand-gesture-interface>. The results explained before show that our hand pose recognition system can be used for controlling robots, for virtual reality interaction, for human-machine interfaces among others.

## References

- [1] Itamar Arel, Derek Rose, and Thomas Karnowski. Deep machine learning - a new frontier in artificial intelligence research [research frontier]. *IEEE Comp. Int. Mag.*, 5:13–18, 01 2010.
- [2] B. Babenko, M. Yang, and S. Belongie. Visual tracking with online multiple instance learning. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 983–990, June 2009.
- [3] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman. Return of the devil in the details: Delving deep into convolutional nets. In *British Machine Vision Conference*, 2014.
- [4] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 22Nd ACM International Conference on Multimedia*, MM '14, pages 675–678, New York, NY, USA, 2014. ACM.
- [5] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1*, NIPS' 12, pages 1097–1105, USA, 2012. Curran Associates Inc.
- [6] Bogdan Kwolek. Face detection using convolutional neural networks and gabor filters. In Włodzisław Duch, Janusz Kacprzyk, Erkki Oja, and Sławomir Zadrozny, editors, *Artificial Neural Networks: Biological Inspirations – ICANN 2005*, pages 551–556, Berlin, Heidelberg, 2005. Springer Berlin Heidelberg.
- [7] Min Lin, Qiang Chen, and Shuicheng Yan. Network in network. *CoRR*, abs/1312.4400, 2013.
- [8] Dennis Núñez Fernández and Bogdan Kwolek. Hand posture recognition using convolutional neural network. In Marcelo Mendoza and Sergio Velastín, editors, *Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications*, pages 441–449, Cham, 2018. Springer International Publishing.
- [9] Oyebade Oyedotun and Adnan Khashman. Deep learning in vision-based static hand gesture recognition. *Neural Computing and Applications*, 28, 04 2016.
- [10] C. Szegedy, Wei Liu, Yangqing Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–9, June 2015.
- [11] Jonathan Tompson, Murphy Stein, Yann Lecun, and Ken Perlin. Real-time continuous pose recovery of human hands using convolutional networks. volume 33, 08 2014.