# Learning Hand Articulations by Hallucinating Heat Distribution - Supplementary Material

Chiho Choi     Sangpil Kim     Karthik Ramani

Purdue University

West Lafayette, IN 47907, USA

{chihochoi, kim2030, ramani}@purdue.edu

This document serves as supplementary materials to our paper *Learning Hand Articulations by Hallucination Heat Distribution*. We present details of our localization network and the system specifications.

## 1. Architecture of the Localization Network

We visualize the graph of our localization network in Figure 1. The proposed network solves two sub-tasks for hand localization: hand segmentation and hand center regression. The segmentation stream identifies pixel-wise class labels through a series of the convolution process, so the resulting probability map can effectively reconstruct the detailed hand segment. Also, the regression stream robustly estimates the centroid of the hand and thus enable us to draw the bounding box around the segmented hand.

## 2. Implementation Specifications

The proposed system was trained with GPUs using the Caffe framework [1]. We trained the localization network by setting the learning rate to 0.0005 and the number of epochs to 250 using the Adam optimizer[1] with $\beta_1 = 0.9$ and $\beta_2 = 0.99$. In addition, our depth network was converged after 80 epochs with the learning rate 0.01 for 60 epochs and 0.001 afterward. Our heat distribution network converged after 250 epochs. Here, we used the learning rate 0.01 for 200 epochs and then dropped it by a factor of 0.1. The hallucination network converged after 100 epochs with the learning rate 0.01. Lastly, the refinement network converged after 170 epochs by dropping the initial learning rate (0.01) in every 60 epochs by a factor of 0.1. At runtime, the computation time for each frame is split as 2 $ms$ for processing data (*i.e.*, depth normalization, resizing, and bounding box cropping), 0.7 $ms$ for hand localization, 1.8 $ms$ to estimate the joint angle parameters from the proposed system. The additional hardware specifications are as follows: Intel's Core i5-4690K, 32GBs RAM, NVIDIA's Geforce GTX 1070, and Intel's RealSense SR300.
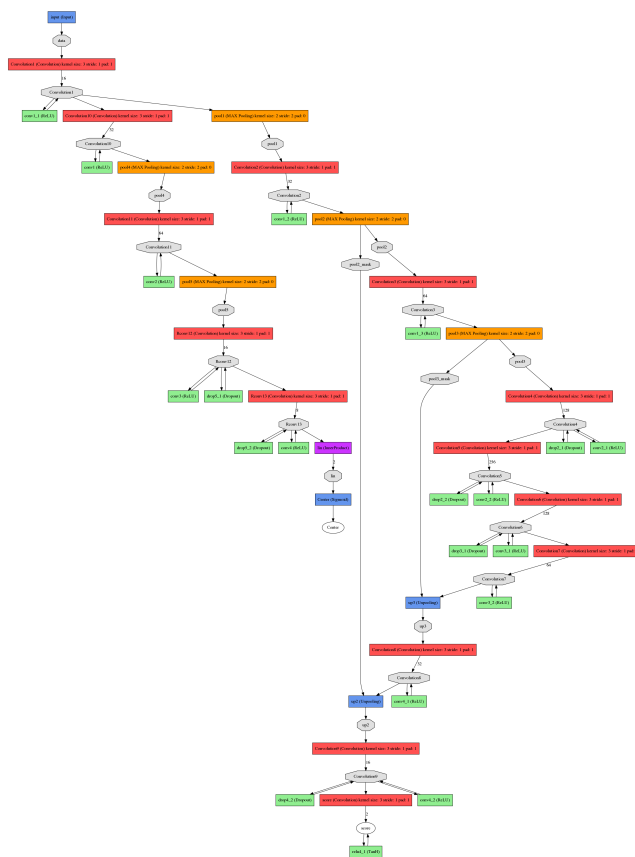


Figure 1: The graph of the proposed localization network architecture. There are two streams: the segmentation stream for pixel-wise hand segmentation, and the regression stream for hand center regression.

## References

[1] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 22nd ACM international conference on Multimedia*, pages 675–678. ACM, 2014. 1

---

[1]The rest of the networks used the SGD optimizer.