

Skeleton Based Activity Recognition Using GCNs

Göksel Mert Çökmez

ETH Zürich

mcoekmez@student.ethz.ch

Abstract—This report details the pipeline of a Graph Convolutional Network called EfficientGCN and its adaptation to the hand skeleton as captured by Microsoft HoloLens. HoloLens and EfficientGCN is deployed in tandem to perform action recognition from the 3D pose of the human hand. This action recognition is performed in the context of a surgery room for tools used in spine surgeries. This report first discusses the pipeline of EfficientGCN and how it was modified in order to complete the task at hand. It then goes into the certain core ideas and data structures that are employed for this model to work, before benchmarking the action recognition accuracy across different configurations of the model. Lastly, the performance of the model is discussed with potential solutions recommended for its shortcomings.

I. INTRODUCTION

Human pose and action recognition is a field on the rise, with novel methods and use cases coming into existence with an increasing rate. While there are already numerous methods for pose estimation such as conventional RGB based methods. However, skeleton based methods prove to be much more reliable since the location of the joints can be represented as a time series and thus we can make inferences using the dynamic changes in our data.

As a result, this project aims to establish a skeleton based method for robust action recognition using the captured hand pose data. When recording the hand pose data, most of the heavy lifting is done by Microsoft HoloLens, which offers accurate coordinates for the 26 joints defined in the human hand.

Thus, one of the main tasks of this project is choosing an already well functioning Graph Convolutional Network (GCN) based solution and adapting it for use in human hands instead of the whole body. In order to achieve this, the main candidate is MMSkeleton [2], which is a part of the open-mmlab project by the Multimedia Laboratory of The Chinese University of Hong Kong. However, due to compatibility and dependency issues with Windows 10, which the operating system for the workstation at hand, we opted to pursue another implementation of a GCN solution.

Our following candidate is EfficientGCN [3] as its dependencies work flawlessly with Windows 10. Moreover, compared to other State-Of-The-Art (SOTA) methods, EfficientGCN uses fewer parameters to achieve similar accuracy scores. This is especially helpful in reducing training time on the NTU RGB+D 60 dataset [1], which serves as a benchmark for numerous human pose recognition models.

The final step in the implementation of this project is recording various hand poses when using surgical equipment

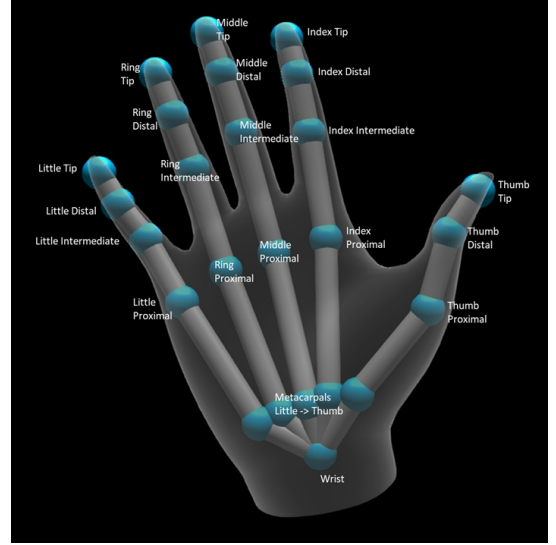


Fig. 1. Hand joints as captured by Microsoft HoloLens

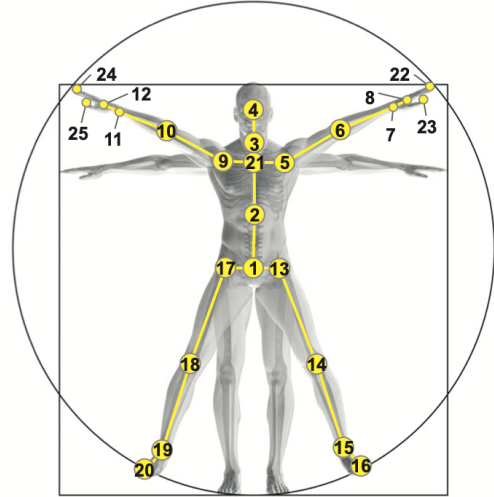


Fig. 2. Human skeleton joints as stated in the NTU RGB+D dataset [1]

related to spine surgery, and converting the HoloLens data to the NTU dataset format, which is the original format used in training EfficientGCN. Once the dataset is converted to the NTU format, the joint indices and their relevant adjacency matrices are modified in order to accommodate the hand model seen in Figure 1, compared to the human skeleton joints in Figure 2.

II. MODELS AND METHODS

The main structure of this project stems from the pipeline of EfficientGCN. As a result, understanding the underlying architecture of EfficientGCN is imperative in order to explore the functionality of this project. Thus, we start with data preprocessing, alongside the explanation of the NTU format and the conversion to it and the form of the graph which diverge from the original implementation of the source code.

A. Data Preprocessing

In this implementation of EfficientGCN, the input is divided into three channels: one for joint positions, one for joint velocities and finally, one for bone features. We start off by normalizing the positions of each joint relative to a predefined center joint with index c . In the original implementation, the joint number 2 of the NTU skeleton is chosen as the center. For our hand pose recognition purposes, we normalize by the wrist index. We define the non-normalized joint coordinates $\mathcal{X} = \{x \in \mathbf{R}^{C_{in} \times T_{in} \times V_{in}}\}$ where C_{in} is the coordinate, T_{in} is the frame index and V_{in} is the joint index. Then, we can define the normalized coordinates $\mathcal{R} = \{r_i | i = 1, 2, 3, \dots, V_{in}\}$ as:

$$r_i = x[:, :, i] - x[:, :, c] \quad (1)$$

The concatenation of the normalized and non-normalized 3D coordinates \mathcal{X} and \mathcal{R} form our first input of joint positions. We then go on to form two sets of velocities $\mathcal{F} = \{f_t | t = 1, 2, 3, \dots, T_{in}\}$ and $\mathcal{S} = \{s_t | t = 1, 2, 3, \dots, T_{in}\}$ by the following equations in order to form our velocity input:

$$\begin{aligned} f_t &= x[:, t+2, :] - x[:, t, :] \\ s_t &= x[:, t+1, :] - x[:, t, :] \end{aligned} \quad (2)$$

Consequently, we form the velocity input by concatenating the matrices \mathcal{F} and \mathcal{S} . Last but not least, we need to form the bone feature input. For that, we define the bone length matrix $\mathcal{L} = \{l_i | i = 1, 2, 3, \dots, V_{in}\}$ and the bone angles $\mathcal{A} = \{a_i | i = 1, 2, 3, \dots, V_{in}\}$. The elements of these matrices are calculated as:

$$\begin{aligned} l_i &= x[:, :, i] - x[:, :, i_{adj}] \\ a_{i,w} &= \arccos\left(\frac{l_{i,w}}{\sqrt{l_{i,x}^2 + l_{i,y}^2 + l_{i,z}^2}}\right) \end{aligned} \quad (3)$$

with i_{adj} denoting the adjacent joint to i and w denoting the 3D coordinates.

B. NTU Dataset Format

The NTU RGB+D is a popular dataset used to benchmark numerous skeleton based human pose estimation algorithms. It is published by the Rapid-Rich Object Search (ROSE) Lab of Nanyang Technological University. The dataset contains 1-4 people in each frame performing certain actions, alongside their associated 3D joint coordinates, RGB image captures and depth maps. Our main focus for our implementation is the files with `.skeleton` extension, which contain the 3D

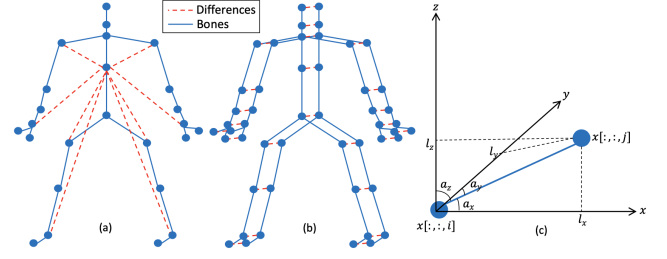


Fig. 3. The above figure indicates how the motion, velocity and bone feature inputs are calculated respectively.

joint coordinates captured by Microsoft Kinect. These files are essentially text files with the total number of frames in the very beginning of the files with the following layout for each frame:

```
S001C001P001R001A001.skeleton - Notepad
File Edit Format View Help
[03
1
72057594037931101 0 1 1 1 0 0.02764709 0.05745083 2
25
0.2181153 0.1725972 3.785547 277.419 191.8218 1036.233 519.1677 -0.2059419 0.05349901 0.9692109 -0.1239193 2
0.232292 0.4326636 3.714767 279.2439 165.8569 1041.918 444.3235 -0.2272637 0.05621852 0.964434 -0.1227094 2
0.2457799 0.6877249 3.633897 281.1529 139.0885 1047.837 367.3966 -0.2456043 0.0660736 0.9535829 -0.1565561 2
0.2128507 0.8079225 3.581995 278.1617 125.6969 1039.476 328.9554 0 0 0 2
0.189304 0.6111551 3.716962 267.2831 148.2511 1007.479 393.5405 0.2179059 0.7353331 -0.6371263 -0.07663021 2
0.180875 0.4286715 3.742593 266.2123 166.5774 1004.187 446.2603 0.03366299 0.7429995 -0.0668001 -0.6650987 1
0.1186992 0.3428923 3.556542 268.5588 173.2219 1011.711 465.4302 -0.5255184 0.8340511 0.1235747 -0.1136005 1
0.1163352 0.3173672 3.550741 268.3543 175.7982 1011.131 472.8558 -0.1148792 0.9877279 -0.0299019 -0.1014999 2
0.3449145 0.5740387 3.588655 291.5822 149.8556 1078.05 398.4594 -0.1106443 0.771432 0.546612 -0.3098302 2
0.4211396 0.3847547 3.627778 298.8716 169.6459 1098.838 455.4754 0.06063642 0.961247 0.2020273 0.1775176 2
0.1918999 0.3189077 3.540756 276.1781 175.5411 1033.77 472.192 0.4641625 0.2649592 -0.3935434 0.7479796 2
0.1298675 0.317039 3.542577 269.7591 175.7559 1015.225 472.7473 0.6194489 0.3464894 -0.3181734 0.6284853 2
0.162766 0.1745395 3.778471 272.0095 191.6048 1020.896 518.4943 -0.05136763 0.6810405 0.7174643 -0.1370778 2
0.1975115 -0.1360554 3.904381 274.8451 221.2371 1028.12 604.0775 -0.1458015 -0.5206363 0.1335629 0.8305664 2
0.2442706 -0.4201954 4.054147 278.3998 246.4252 1037.602 676.7612 -0.0476737 0.1068644 0.244998 0.9624379 2
0.2040377 -0.476396 4.094294 274.5897 251.0867 1026.418 690.1627 0 0 0 2
0.2699234 0.1678491 3.730005 282.8167 192.0363 1052.036 519.8357 -0.2298583 0.6472055 0.644957 -0.3351428 2
0.3105748 -0.1462299 3.818383 286.1002 222.5039 1060.931 607.8004 0.08184162 0.837252 0.1252462 0.5259509 2
0.3428809 -0.433214 3.943258 288.1872 248.7218 1066.199 683.3955 0.2016301 0.9774966 0.05738356 0.02351441 2
0.2951482 -0.5017325 3.968896 283.5814 254.7903 1052.752 700.8542 0 0 0 2
0.2425592 0.6247278 3.650655 280.6673 145.8599 1046.329 386.8308 -0.2489336 0.06207102 0.9568546 -0.1364125 2
0.1200176 0.2955468 3.532409 268.7763 177.8907 1012.419 478.8939 0 0 0 2
0.0989792 0.3459614 3.527539 266.6139 172.6136 1006.227 463.6581 0 0 0 2
0.08710064 0.2988889 3.546119 265.3322 177.6662 1002.409 478.2144 0 0 0 2
0.1268453 0.2813963 3.549525 269.4193 179.4981 1014.19 483.5357 0 0 0 2
```

Fig. 4. NTU dataset format. The first line indicates the number of frames in this recording. The second line is the number of skeletons in that specific frame, with the third line acting as the identifier as the skeleton ID. Then the file goes on to list the number of joints visible in each frame and their coordinates.

As it can be observed in the NTU format, there are 11 coordinates associated with each joint. These stem from the recording format of Kinect v2, and they correspond to (in order):

- 3D x coordinate
- 3D y coordinate
- 3D z coordinate
- RGB camera 2D x coordinate
- RGB camera 2D y coordinate
- Depth camera x coordinate
- Depth camera y coordinate
- Quaternion w coordinate
- Quaternion x coordinate
- Quaternion y coordinate
- Quaternion z coordinate

However, a large part of these coordinates are redundant. Our script in HoloLens captures uniquely the 3D x , y and z coordinates, which is very useful since these three coordinates are the only ones used by EfficientGCN. Thus, when converting

the JSON file generated by the HoloLens to the NTU format, we end up with a similar layout, albeit with 26 joints and uniquely the 3D coordinates.

Fig. 5. Hand joint coordinates in the NTU format. Notice that the joint coordinates only list the 3D coordinates.

C. Graphs

Another important modification in the code is the inclusion of the aforementioned hand skeleton. In order to convey information regarding the skeletons, EfficientGCN stores a kind of adjacency matrix in which the positions of the joints are described relative to the center joint, which is the wrist joint in our case.

```
if self.dataset == 'hand':
    num_node = 26
    neighbor_lbase = [(6, 5), (5, 4), (4, 3), (3, 1),
                     (11, 10), (10, 9), (9, 8), (8, 7), (7, 1),
                     (16, 15), (15, 14), (14, 13), (13, 12), (12, 1),
                     (21, 20), (20, 19), (19, 18), (18, 17), (17, 1),
                     (26, 25), (25, 24), (24, 23), (23, 22), (22, 1)]
    neighbor_link = [(i - 1, j - 1) for (i, j) in neighbor_lbase]
    connect_joint = np.array([5, 4, 3, 1, 10, 9, 8, 7, 1, 15, 14, 13, 12, 1, 20, 19, 18, 17, 1, 25, 24, 23, 22, 1]) - 1
    parts = [
        np.array([6, 5, 4]) - 1,      #thumb
        np.array([11, 10, 9, 8]) - 1, #index
        np.array([16, 15, 14, 13]) - 1, #middle
        np.array([21, 20, 19, 18]) - 1, #ring
        np.array([26, 25, 24, 23]) - 1, #pinky
        np.array([1, 2, 3, 7, 12, 17, 22]) - 1 #palm
    ]
```

Fig. 6. Data structure detailing the bones and the extremity of the bone closer to the center.

Notice the data structure outlined in Figure 6 that is found in the file *graph.py*. The *neighbor_lbase* matrix contains the bones by specifying end joints, and *connect_joints* specifies which extremity of the bone is closer to the center, defining its orientation.

Moreover, similar to limbs, the hand skeleton is divided into six separate sections. There, we can see that joints are grouped under the following parts:

- Thumb
- Index
- Middle finger
- Ring finger
- Pinky
- Palm

III. RESULTS

One thing that is very important to note regarding the benchmarks is the vast difference between accuracies of the 3 action and 4 actions datasets. The activities in question are the following:

- Piercing with the Lenki
- Poking with the Knopfsonde
- Tightening with the screwdriver
- Cutting with the scalpel

The main difference between our datasets is that the 3 action dataset only contains the first three action, whereas the 4 action dataset contains all actions. In the following tables, results from normalized and non normalized coordinates are compared alongside left hand only training compared to both hand training.

TABLE I
3 ACTION DATASET ACCURACIES

	<i>Normalized</i>	<i>Non-normalized</i>
<i>Left hand training</i>	$\mu = 0.83$ $\sigma = 0.0714$	$\mu = 0.82$ $\sigma = 0.1156$
<i>Left and right hand training</i>	$\mu = 0.78$ $\sigma = 0.0905$	$\mu = 0.89$ $\sigma = 0.0275$

TABLE II
4 ACTION DATASET ACCURACIES

	<i>Normalized</i>	<i>Non-normalized</i>
<i>Left hand training</i>	$\mu = 0.61$ $\sigma = 0.0917$	$\mu = 0.71$ $\sigma = 0.1067$
<i>Left and right hand training</i>	$\mu = 0.81$ $\sigma = 0.0249$	$\mu = 0.70$ $\sigma = 0.0265$

As we can see, the 4 action dataset performs significantly worse than the 3 action one. The main reason behind that is when data is well separated spatially, the model performs well. However, when we introduce another action that has similar temporal data to another existing action but differs temporally, the model struggles to differentiate between the two actions, leading to a significantly lower accuracy.

Apart from that, we can clearly observe that normalization leads to higher accuracies throughout. The only exception is the left hand normalized training with 4 action dataset, where the network fits on the location of the hand rather than its pose. However, this is a suboptimal fit. The location of the hand when performing the action should be irrelevant since we would like to focus on the pose instead. Since the scalpel and Knopfsonde are held in similar poses, this suboptimal fit increases the accuracy nonetheless. Lastly, while the use of both hands in training leads to a more accurate model, left hand training is not that far behind and offers a comparable accuracy.

IV. DISCUSSION

A very evident observation from our results is that our model performs well with distinct spatial configurations, while

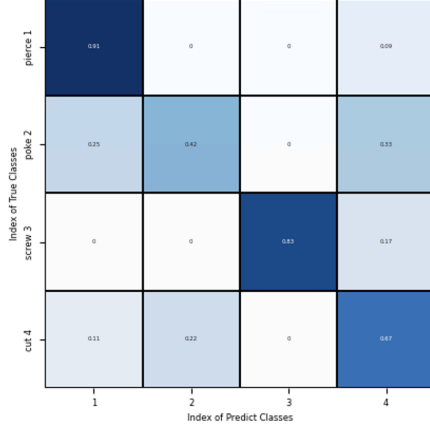


Fig. 7. The confusion matrix for the 4 action dataset. Notice the mix-ups between poking with Knopfsonde and cutting with scalpel. While these two actions differ temporally, they are very close spatially, leading to this state of confusion.

it suffers when the difference is limited to the temporal domain. A possible explanation for that is the limited frame window. In EfficientGCN, the datasets recorded with Kinect have around 60-100 frames while our dataset contains around 350-1000 frames due to the better capabilities of HoloLens. Since the maximum number of frames allowed by our GPU’s memory is around 500 frames, the window cannot cover the whole action.

A direct approach to this issue would be simply increasing the number of frames the model is trained on, necessitating a larger memory for our GPU. However, since HoloLens records with such a high speed, we have around 60 frames covering one second in our recordings, which leads to bloated frame numbers. We could simply include one in five frames, effectively reducing the memory requirements of our model with very little information loss since a decrease from 60 FPS to 12 FPS should still retain the necessary temporal information for model fitting. As a result, the 500 frame window would be able to cover the whole duration of the action without the need of introducing more powerful hardware.

It should also be noted that our dataset is relatively small. In order to have more concrete information regarding the accuracy in the operating room, a larger dataset where the surgery tools are used by experts will yield a more accurate benchmark since our model is geared towards a very specific context that is the surgery room.

V. SUMMARY

As a result of our benchmarks, we can conclude that GCNs are indeed reliable methods for skeleton based activity recognition. Even though the model was not extensively trained and it was trained on a small dataset, the resulting accuracies are already promising. This enables it to be potentially used for more niche tasks of which there are few datasets about.

The model can be trained and deployed quickly for accurate results, and it is especially well suited for spatially distinct tasks. However, it should be noted that these benchmarks are also done in a very niche setting. In order to claim a wider use area, the model should be trained on a much larger dataset with more quotidian actions having different spatial and temporal configurations. A benchmark performed with this trained model would be more representative of the model’s capabilities in the bigger picture.

REFERENCES

- [1] Amir Shahroudy, Jun Liu, Tian-Tsong Ng, and Gang Wang. Ntu rgb+d: A large scale dataset for 3d human activity analysis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1010–1019, 2016.
- [2] Jingbo Wang, Dahua Lin, Sijie Yan, Yuanjun Xiong. Mmskeleton. <https://github.com/open-mmlab/mmskeleton>, 2019.
- [3] Yi-Fan Song, Zhang Zhang, Caifeng Shan, and Liang Wang. Efficientgc: Constructing stronger and faster baselines for skeleton-based action recognition. *arXiv:2106.15125*, 2021.