
"Follow Through": An Improved Machine Learning Approach to Analyzing Basketball Performance

Sameer Chaturvedi
Boston University
sameerc@bu.edu

Eli Saracino
Boston University
esaracin@bu.edu

Kamoltat Sirivadhna
Boston University
ksirivad@bu.edu

Abstract

1 This paper describes the steps taken and preliminary results of "Follow Through",
2 a one-stop athletic solution for monitoring one's progress and performance through
3 the lens of machine learning. Below, we describe the foundational OpenPose
4 technology, originally developed at Carnegie Mellon University, that we base our
5 framework on, as well as the process of improving that framework for use in a new,
6 more specific setting. This includes the smoothening of the skeletons OpenPose
7 draws when processing a video file, and applying these smoother skeletons to
8 system that converts those skeletons into a standardized vector of values that repre-
9 sent a user's basketball jumpshot. The technology implemented, going forward,
10 could be used in any number of computer vision related applications wherein such
11 movements could be analyzed.

12 1 Introduction

13 Recently, the advent and advancement of Deep Neural Networks have provided myriad new opportu-
14 nities for machine learning. From self-driving cars, to simple image classification, processing images
15 has, in particular, become a cornerstone of this progress. Deep Convolutional Neural Networks are
16 particularly admired for their strengths in this arena, as their ability to apply all the fidelity of a
17 standard multilayer perceptron, but offer regularization in the form of identifying hierarchical patterns
18 in the data, and thus reduce overfitting, is extremely powerful.

19 We wanted to build off of some of that power. Utilizing an already trained Deep CNN, with some key
20 adjustments, such as a post-processing smoothing phase through the application of a Savgol Filter, we
21 could apply the problem of joint recognition to that of athletic improvement. Namely, we can break
22 down the specifics of a basketball player's jumpshot into a series of key moments that allow us to
23 compare the performance of a given player to either a professional athlete, or any other player. In this
24 way, we can apply improved machine learning and computer vision techniques to enter a field left
25 largely untouched by the advent of this technology, and offer the user step-by-step goals to improve
26 their game.

27 2 OpenPose CMU

28 OpenPose is a technology developed at Carnegie Mellon University to apply a VGG-19 Deep
29 Convolutional Neural Network to a particular kind of image processing. Specifically, given a input
30 image or set of images, the trained Neural Network could detect and draw key points on the human
31 subjects (head, neck, hands, etc.) [Cao+18]. The opensource Github Repository can be found below.

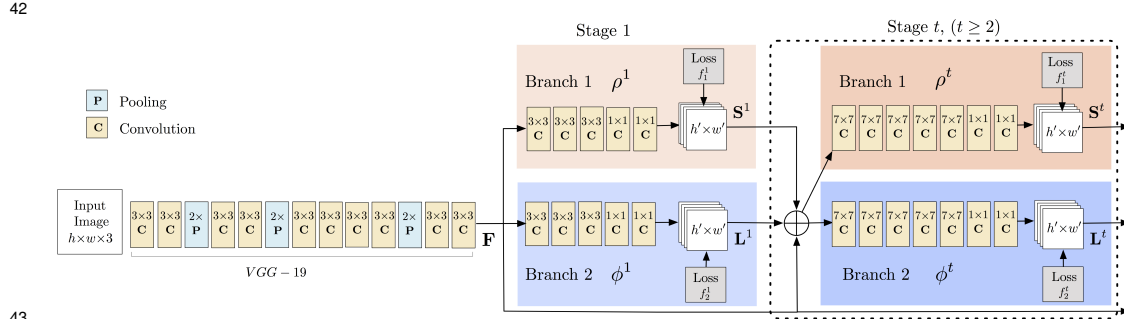
32 <https://github.com/CMU-Perceptual-Computing-Lab/openpose>

The trained weights and structure of the CNN developed at CMU is freely available, and was easily accessible for our own means through Python's OpenCV library.

2.1 VGG-19

The Deep Convolutional Neural Net that OpenPose is based on, VGG-19 (or Visual Geometry Group-19) is a well-known Deep Neural Net trained on more than one million ImageNet images. It's famous for its low error rates, obtaining about a 10 percent error rate on ImageNet, and the fact that its weights and structure are freely available for use [SZ14].

The architecture for the VGG-19 Neural Net is as follows:



Of note here is the three stage process used to process a given image. The first stage involves the first ten layers of the Net, which are used to create feature maps for the input image, to facilitate later processing [Zhe19].

After these feature maps are created, the CNN essentially branches into two distinct parts. The first such branch predicts a set of 2D confidence maps of mpoints on the body (such as the elbow, knee, head, etc.). The second branch predicts a set of 2D vector fields of "part affinities", which effectively encode the degree of association between the parts determined in the first branch (i.e. between the neck and shoulder, or elbow and wrist).

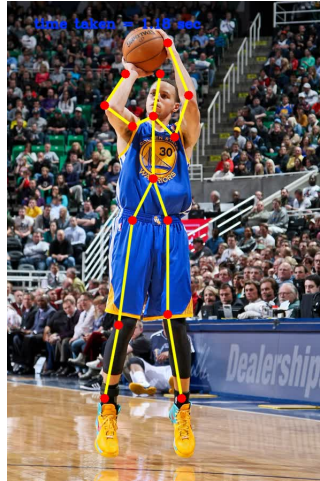
Finally, the confidence and affinity maps are combined to produce the two-dimensional keypoints for the subject of the image, which can then be used to draw the skeletons of our athletes [Cao+18].

2.2 MPII dataset for keypoint detection

The Max Planck Institut Informatik (MPII) dataset includes around 25K images containing over 40K people with annotated body joints. The images were systematically collected using an established taxonomy of every day human activities. Overall the dataset covers 410 human activities and each image is provided with an activity label. Each image was extracted from a YouTube video and provided with preceding and following un-annotated frames. The link to the official website of MPII data set can be found below.

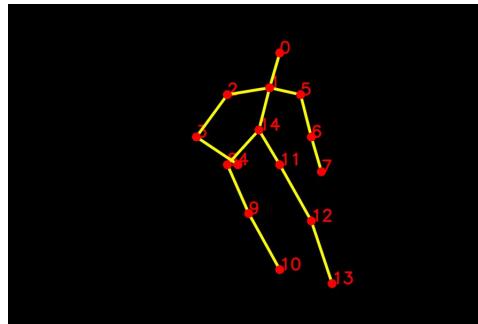
<http://human-pose.mpi-inf.mpg.de/>

This is the dataset that the aforementioned Deep Convolutional Neural Network was trained on by the OpenPose team, and we were able to plot the skeletons of the players based on the two dimensional points returned as the output of the Net. Notably, the MPII dataset highlighted 15 keypoints for human detection, which is less than some other similar datasets (such as COCO, Common Objects in Context, which tries to track as many as 18, including facial features unnecessary for our purpose). A given skeleton might look like the following, after running the image through the Neural Network, and connecting the points computed as the output of the process:



72

73 After the skeleton is computed, it's often easier to process just the skeleton, as the colors and objects
 74 present in the background of the image serve mostly to add noise to our later processing of the
 75 skeleton. A given blank skeleton might look like the following:



76

77 This format allows us to focus solely on corresponding keypoints between two images (or, expanding
 78 this further, two videos), and process them purely analytically.

79 2.3 Issues with OpenPose and MPII

80 While OpenPose is clearly a powerful technology that could offer an effective way to analyze a given
 81 player's basketball shot, it is not without fault. For instance, in processing videos, or noisy images
 82 (potentially of multiple subjects), it's very easy for OpenPose to induce jitter into its outputs, as can
 83 be seen below:



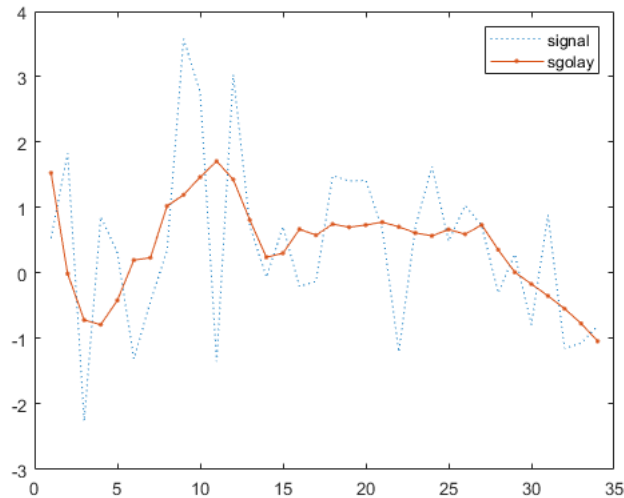
84

85 When processing a video with such jitter, it makes it much hard to capture the specifics of a player's
 86 performance. The first step in developing our technology, then, became smoothening the skeleton's
 87 attachment to the human subject- essentially, taking the skeleton output by OpenPose, and apply
 88 some technique that would increase its fidelity to the actual movements taken by the player.

3 Smoothing the OpenPose skeleton

Although OpenPose is arguably good enough to create an accurate skeleton on the player who is in the video, we found out that it is actually not consistent, since in some cases the model causes the body parts of the players to bounce around the true value. With this we later found out that the data is distributed normally around the true value, so a filter can be used to smooth out the data, and generate accurate values from noisy data. Amongst many filtering techniques, we found Savitzky-Golay filter (Savgol) to be the most effective in smoothing out the data and creating an accurate value.

3.1 Savitzky-Golay Filter



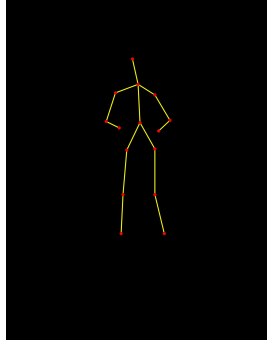
As can be seen, at a high level, the filter is limiting the number and severity of the peaks and valleys of its function of input. In the context of OpenPose, what this allows is a much less frequent occurrence of jitter, or moments in the processed video or image where the skeleton tends to have trouble matching the keypoints it's supposed to. The result is a skeleton that fits the player much better, which allows for more accurate measurements to be taken in building a user's model, described in more detail below.

4 Applying the smoother skeletons: "Follow Through"

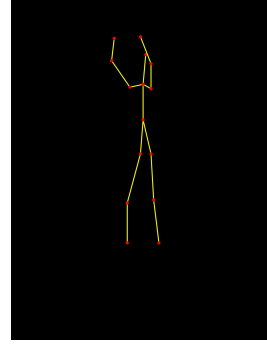
After the application of our Savgol filters, the results of applying OpenPose's Neural Network drastically improved. It was then that we could shift our focus to implementing the core of the technology that would actually use the skeletons to analyze basketball performance. The general process would be as follows: build a model for each user of our application, load in .mp4 video files of that user's jumpshot, and use the gathered information about that player to compare his/her performance with other players, or even professionals.

4.1 Processing an input video through Computer Vision and pattern matching

The framework used to process the input video was, of course, OpenPose's CNN. Given the output skeleton for every frame of an input video, we could compare that skeleton to two of our average template skeletons. These average templates, shown below, were computed by combining the output skeletons from each of several different images hand-labeled as the base and release of a jumpshot, respectively. By creating these averages, and comparing an input video's frames to them, we are able to capture the point in the input video when the user reaches the base and release of their shot through simple pattern matching of the user's skeleton to the average.



(a) Average Skeleton for Base of Shot



(b) Average Skeleton for Release of Shot

Figure 1: Average Skeletons used for Pattern Matching

Measurements used in the pattern matching include a specific range of angles for the elbows to denote the user's position, as well as a computation of the average distance of each joint from its pairwise position in the template. Fine tuning these values was non-trivial, and required much tweaking by hand. However, in the end, we came up with a range that seemed to provide a good result for the pattern matching. For the base template, the degree of elbow angles, summed together, could not be more than 150 degrees different from that of the template, and the average Euclidean distance of each pairwise keypoint couldn't be greater than 210. If both of these conditions were met, and the base was not already found, we labeled that frame as the base. For the release, the elbow degree difference could not exceed 60 degrees, and the pairwise distance from the template could not exceed 230.

As can be inferred from the above, it became easier to restrain the overall joint difference on the base template, and easier to restrain the angle of the elbows when comparing to the release template. However, once these key points were determined in a given input video, we could compute a series of measurements that could quantify that shot. Namely, we recorded the average right and left elbow angle at the base and release of the shot, the average length of the rise of the shot and the total shot overall, in number of frames, and the average time in the shot when the base and release were met. These six initial measures were computed (or averaged, if a user supplied multiple shot videos), and used to build a given user's model.

4.2 Preliminary results of model comparisons

While the technology is currently limited, due to the time constraints placed on us by virtue of this project only lasting a single semester, the preliminary results of our system are promising. As opposed to processing any human users at this time, we decided an efficient way to test the accuracy of our model would be to build models based off the jumpshots of pro basketball players, and compare these models instead. Ideally, the behavior available when analyzing a pros model should be the same as any given user, and the functionalities tested could easily be extended to input videos of that user's jumpshot.

4.2.1 Creating pro models

To that end, we used the well-known basketball simulation game, NBA 2K19, to create data to use for each pro's shot. It is known that the character models present in the game are realistic representations of their real-life counterparts. They are even created by bringing professional athletes in, equipping them with sensors, and letting them play the game [Spo18].

What this allowed, then, was a perfect setting for us to create the videos we needed to build a professional player's profile. We could, using a freeplay mode in the game, have the players take shots, specify the exact camera angle and distance, and record the entire shot, for later use in building a model for that player. In this way, we created very simple models from each player, using only one to four videos for each athlete (which, admittedly, builds a model still very prone to error).

154 4.2.2 Comparing pro models

155 Even with the limited number of samples used, however, we could already see some preliminary
156 results. Using the average body measurements at different points in the shot, as mentioned above in
157 section 4.1, we could effectively treat each player's model as a single, six-dimensional vector, and
158 compare them in any of several ways.

159 Our preliminary comparison was the following: we held out one of our videos of Michael Jordan
160 when we built his model. Then, after all models had been made for each of our pro athletes, we added
161 this held-out sample as a new, unidentified user. Finally, we simply computed this new user's distance
162 to each professional model built, and, sure enough, Michael Jordan's model was the closest:

```
jordan's jumpshot is closest to...  
[122.11437783004688, 'jordan']  
[146.3981404622901, 'lebron']  
[153.40539578518727, 'wade']
```

163

164 Other simple tests that could be run include, for every pro player, find out which other pro player's
165 jumpshot is most like their own (shown in part below):

```
shaq is closest to...  
[34.736617092689514, 'noah']  
[49.82063139929438, 'dirk']  
[50.291850196503226, 'marion']  
  
dirk is closest to...  
[49.82063139929438, 'shaq']  
[55.60234256461159, 'noah']  
[57.99514601599389, 'jordan']  
  
lebron is closest to...  
[21.85778795836751, 'wade']  
[38.158792923137334, 'steph']  
[41.30439788216181, 'barry']
```

166

167 5 Conclusion and future steps

168 The codebase pertaining to our experiments and results is here:

169 <https://github.com/esaracin/cs542-follow-through-complete>

170 The README.txt provides explicit instructions for reviewing and replicating our results.

171 While these results are promising, there is obviously still much left to do. Pairwise comparisons are
172 interesting, but they don't, by themselves, answer the original problem in depth. Going forward, we
173 plan to use the results of these model comparisons in several different ways: we can perform a kind
174 of classification, essentially answering the question of which pro athlete's jumpshot is most like a
175 given "Follow Through" user's. At this point, we could go even further, and be more specific: given a
176 specific pro athlete, we can specify where a given user's own shot differs, and provide concrete steps
177 for them to improve.

178 The implications of this technology also supersede just the analysis of basketball play, as well.
179 Any sport or activity where computer vision and machine learning could be applied, the use of our
180 improved OpenPose skeletons and template matching could offer insight into one's performance.

181 5.1 Improvements to the pattern matching and model creation

182 In addition to simply widening the breadth of services offered by "Follow Through," we'd also like
183 to make improvements to the base framework of the technology. Namely, we'd like to further tune

184 the template matching to be more accurate, as more accurately matching the specific key moments
185 in a shot makes the model built from their information much more useful. This could further be
186 remedied by the measurement of more angles, such as the bend of the knee, which could provide
187 useful information, not just for pattern matching, but for actually building the user's model. The more
188 measurements we have, the more specific the model will be to a given user, and the more accurate
189 our report of the user's performance can be.

190 We could also improve the templates used for matching. Right now, due to the limited scope of this
191 project, our average templates essentially consist of ten to fifteen hand-picked "base" and "release"
192 images. Given many more samples, our averages could be adjusted to better account for variations in
193 the base and release of the basketball shot across different players.

194 5.2 Training the CNN on a basketball specific dataset

195 Going forward, we might also look into adjusting the CNN used for the pose estimation portion of
196 the technology. Having been trained on the MPII dataset, there are four-hundred and nine activities
197 present in the myriad images there that are completely irrelevant to the game of basketball. If we
198 removed those images from the set, and trained the Neural Net solely on basketball images, we could
199 potentially further improve the accuracy of the pose estimation by limiting the noise induced by the
200 poses present in these images.

201 References

- 202 [SZ14] Karen Simonyan and Andrew Zisserman. "Very Deep Convolutional Networks for
203 Large-Scale Image Recognition". In: *arXiv e-prints*, arXiv:1409.1556 (Sept. 2014),
204 arXiv:1409.1556. arXiv: 1409.1556 [cs.CV].
- 205 [Cao+18] Zhe Cao et al. "OpenPose: Realtime Multi-Person 2D Pose Estimation using Part Affinity
206 Fields". In: *arXiv e-prints*, arXiv:1812.08008 (Dec. 2018), arXiv:1812.08008. arXiv:
207 1812.08008 [cs.CV].
- 208 [Spo18] 2K Sports. *NBA 2K19*. 2018.
- 209 [Zhe19] ZheC. *Realtime_{Multi}—Person_{Pose}Estimation*. 2019. URL: [https://github.com/
210 ZheC/Realtime_Multi-Person_Pose_Estimation](https://github.com/ZheC/Realtime_Multi-Person_Pose_Estimation) (visited on 05/10/2019).