

# Assessing Pull-Up ROM with Deep Learning Pose Estimation

Cameron Boydd

## Abstract

The goal of this project is to estimate if a pull-up taken from video is a full range-of-motion or "perfect" pull-up. This project utilizes a pre-trained deep learning human pose estimation model to extract key points and assess the form of the pull-up exercise. By analyzing the 2-dimensional spatial relationships between these key points and using vector arithmetic and trigonometry, the program calculates the angle of arm extension and flexion. To ascertain the completeness of the pull-up movement, the program compares the calculated arm extension or flexion angle against a predefined threshold indicative of full extension or flexion. A "perfect" pull-up is completed when the subject is in the bottom state in full-extension, moves to the top state at the predefined threshold for flexion, and back to the bottom state in full-extension.

## 1. Introduction

Strength training is important for enhancing muscular strength and overall well-being. Among exercises, the pull-up stands out as a fundamental compound movement that targets multiple muscle groups, including the back, shoulders, and arms. However, doing a pull-up with proper form, characterized by a full range of motion, can be challenging.

This program addresses the challenge of assessing pull-up form using human pose estimation techniques. The objective is to accurately detect key body points, such as the shoulders, elbows, and wrists to determine the completeness of each pull-up.

To achieve the objective, the program utilizes the pre-trained deep learning model "Caffe," trained on the COCO dataset. The COCO dataset offers a large number of annotated images providing detailed keypoints for human pose estimation.

The experimentation started with custom data collection and annotation with a small dataset of people performing pull-ups in varying conditions. However, due to limitations in model training the approach switched to utilize a pre-trained model for improved performance.

As a result, the project demonstrates success with the system successfully processing input frames and generating

heatmaps indicating body keypoints. By analyzing the spatial relationships between these keypoints, the system calculates the angles of arm extension and flexion, enabling real-time assessment of pull-up range-of-motion.

In the subsequent sections of this report, it delves into the technology, experimental setup, results, and conclusions.

## 2. Related Work

This project used a couple of open-source materials. The open-source materials used are OpenCV, Numpy and a pre-trained deep learning model called "pose-iter-440000.caffemodel."

### 2.1. OpenCV

OpenCV (Open Source Computer Vision Library) is an open-source computer vision and machine learning software library. It provides a wide range of functionalities for processing images and videos. OpenCV is written in C++ and includes interfaces for Python, Java, and MATLAB/Octave, among others. OpenCV was used for many of its functions. In the project it loads the deep-learning model, opens a camera and/or loads a video, reads a frame, sets the input image for the neural network and performs forward pass inference on the input image. It also draws on the output frame for the displayed video. It draws circles representing detected key points, adds text annotations to the frame, draws lines connecting key points to form a pose skeleton and displays the resulting frames.

### 2.2. NumPy

NumPy, short for Numerical Python, is a powerful library in Python used primarily for numerical computing. In this project, NumPy is primarily used for numerical computations involving vectors and matrices. The coordinates of key points are initially extracted as tuples and then converted into NumPy arrays. This conversion allows easy manipulation and computation of vectors. NumPy arrays support element-wise subtraction, resulting in the calculation of vectors representing the upper arm and forearm. The dot product between the upper arm and forearm vectors is computed using NumPy's np.dot() function. The magnitudes of the upper arm and forearm vectors are computed using the np.linalg.norm() function. NumPy functions np.arccos()

080 and `np.degrees()` are used to calculate the angle between the  
081 upper arm and forearm vectors in radians and then convert  
082 it to degrees.

### 083 2.3. Caffe

084 Caffe (Convolutional Architecture for Fast Feature Embed-  
085 ding) is an open-source deep learning framework developed  
086 by Berkeley AI Research (BAIR) and community contribu-  
087 tors. In this project, the deep learning Caffe model is used  
088 for human pose estimation. Included is the `.caffemodel`  
089 file that contains the trained weights of the neural network  
090 and the `.prototxt` file that describes the network architecture.  
091 The input frame is turned into a blob suitable for the model,  
092 then a forward pass is applied and saved as the output which  
093 contains the predicted key point heat-maps.

## 094 3. Data

095 This project uses the pre-trained deep learning model  
096 "Caffe." This model was trained on the COCO dataset. The  
097 COCO dataset contains a large collection of images with  
098 complex scenes and diverse backgrounds. The images are  
099 sourced from everyday contexts covering a wide range of  
100 scenarios. Each image in the COCO dataset is densely an-  
101notated with bounding boxes around object instances, seg-  
102mentation masks delineating object boundaries, keypoints  
103 indicating specific body joints or landmarks, and captions  
104 describing the image content. COCO annotations include  
105 keypoints for human pose estimation, specifying the coor-  
106 dinates of keypoints such as joints in each annotated person  
107 instance. The COCO 2020 Keypoint Detection Task sup-  
108 plied the COCO train, validation, and test sets containing  
109 more than 200,000 images and 250,000 person instances la-  
beled with keypoints.



Source : [COCO 2020 Keypoint Detection Task](#)

110  
111 The formatting of the data includes `x` number of image  
112 files with an equal number of annotation files that correlate  
113 to each image file. The annotations in the case of keypoint  
114 detection would be a list of `x` and `y` coordinates (Ex: "key-  
115 points": `[x1, y1, x2, y2, ..., xn, yn]`). Using the preprocessed  
116 COCO dataset, the model learns to predict the locations of  
117 keypoints based on input images, aiming to minimize a loss

118 function that measures the discrepancy between predicted  
119 keypoints and ground truth annotations.

## 120 4. Experiments

121 The first experiment started with custom data collection and  
122 annotation. Experimenting started with 10 images of people  
123 wearing different colors, doing pull-ups in different angles  
124 and in different ranges of the repetitions. For each image,  
125 a correlating .txt file containing `x` and `y` coordinates of the  
126 keypoints located at the shoulder joints, elbow joints, and  
127 wrists was created. The tool used for annotating keypoints  
128 was CVAT, an open-source annotation tool for labeling data.  
129 The closest attempt to training a model with custom data  
130 was using the YOLO-pose model and training it on this  
131 data. This ended up failing. Downloading other models to  
132 train and downloading public data sets also ended up failing.  
133 This is why a pre-trained model was used. This project  
134 required a well-trained model to capture the keypoints of  
135 the input and reach the goal of the project.

## 136 5. Results

137 The results of the project align with the goals. The program  
138 successfully passes input frames as blobs through the model  
139 producing heatmaps as outputs.

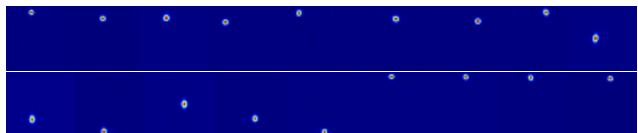


Figure 1. 18 heatmaps indicating 18 different keypoints

140 Using these heatmaps, the keypoints necessary for pull-  
141 up detection are extracted. By analyzing the spatial rela-  
142 tionships between these keypoints using vector arithmetic  
143 and trigonometry, the program calculates the angle of arm  
144 extension and flexion. These angles determine whether the  
145 person in the video is in the top state or bottom state of the  
146 pull-up. The pull-ups are incremented when the person en-  
147 ters the bottom state from the top state. If the person cuts  
148 the repetition short by not moving with full range-of-motion  
149 it will not increment.

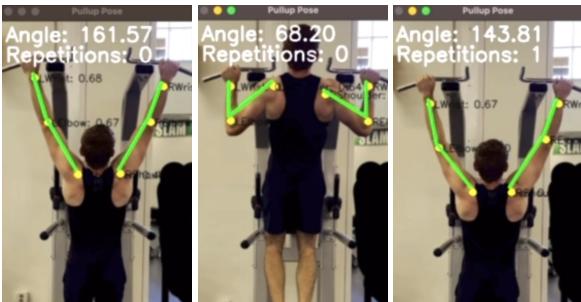


Figure 2. Bottom state, to top, to bottom +1 rep increment

## 150 6. Conclusion

151 The proposed project aimed to leverage human pose estima-  
 152 tion techniques to analyze pull-up range-of-motion with the  
 153 goal of assessing the completeness of each repetition based  
 154 on the angles of the elbows. By utilizing a pre-trained deep  
 155 learning model trained on the COCO dataset, the project  
 156 successfully extracted key joint points from input frames  
 157 and calculated the angles of arm extension and flexion to  
 158 determine if the pull-ups were done with a full range-of-  
 159 motion.

160 Through the generation of heatmaps and extraction of  
 161 relevant keypoints, the program calculated the spatial re-  
 162 lationships between shoulder, elbow, and wrist joints. By  
 163 implementing vector arithmetic and trigonometric calcula-  
 164 tions, the program identified the transition between the top  
 165 and bottom states of the pull-up.

166 While initial attempts at custom data collection and  
 167 model training proved challenging, the decision to utilize  
 168 a pre-trained model helped achieve the project's objectives.  
 169 By building upon these open-source resources, the project  
 170 successfully demonstrated the feasibility of using pose esti-  
 171 mation for real-time evaluation of exercise form.

## 172 References

- 173 [1] Coco dataset: All you need to know to get started. V7.  
<https://www.v7labs.com/blog/coco-dataset-guide>
- 174 [2] Dmitry Golovkin, "CNN Caffe Models," Available at: <https://github.com/foss-for-synopsys-dwc-arc-processors/synopsys-caffe-models/blob/master/README.md>
- 175 [3] Gholamalinezhad, H., Khosravi, H. (2020, September 16). Pooling methods in Deep Neural Networks, a review. arXiv.org. <https://arxiv.org/abs/2009.07485>
- 176 [4] Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., Darrell, T. (2014, June 20). Caffe: Convolutional Architecture for fast feature embedding. arXiv.org. <https://arxiv.org/abs/1408.5093>