

# How Good Is Your Pose?

## Pose Estimation for Weight Lifting Form Correction

Alexander Martin <sup>\*†‡</sup>

Michael Kingsley <sup>\*†</sup>

Vashisth Tiwari <sup>\*†</sup>

### Abstract

1 *Pose estimation is a well-researched field in computer  
2 vision with many published architectures for the task. However,  
3 these methods are not widely used in day-to-day applications due to model computing cost or the general public's  
4 lack of understanding of the methods. Currently, the largest  
5 and most widely advertised deployment of pose estimation  
6 in a non-technical field is the use of automatic offside detection  
7 in the FIFA 2022 World Cup [2].*

8 *Inspired by the use of these methods in a widely spec-  
9 tated sport, we introduce another use of pose estima-  
10 tion in athletics that anyone can use: pose estimation for  
11 weightlifting form. In this paper, we explore this for the  
12 three widely performed exercises: bench presses, deadlifts,  
13 and squats. The framework uses a CNN to classify the lifts  
14 and ViTPose to identify key joints and to make suggestions.  
15 We find that pose estimation can be successfully used for  
16 online coaching without human supervision and can even  
17 be used in competition athletics to flag lifters with illegal  
18 forms.*

### 1. Introduction

20 Form is crucial in athletic performance. Minor changes  
21 in form can increase power, speed, and reduce the risk of  
22 injury for athletes. Form is easy to critique in weight lifting  
23 due to the standards from Olympic and competition lifting.  
24 There is less variability in form for weight lifting than  
25 in other sports where more advanced bio-mechanics are at  
26 play. In order to improve form, athletes must practice and  
27 focus on proper technique. A qualified coach or trainer can  
28 help athletes understand the proper form and provide feed-  
29 back and guidance on how to improve.

31 In this paper, we find that pose estimation can be suc-  
32 cessfully used for online coaching without human supervi-  
33 sion and can even be used in competition athletics to flag

lifters with illegal form. We propose the use of Deep CNN  
34 to classify exercises and then apply a 2D human pose esti-  
35 mation model to provide ways to score the form of lifters in  
36 lieu of a personal trainer.  
37

38 First, the CNN will help identify which exercises the  
39 lifter is doing and can be used as a starting point for the form  
40 analysis. The CNN will be trained on a dataset of images of  
41 people performing either barbell bench press, barbell squat,  
42 or barbell deadlift. Based on this training, we can identify  
43 which of these three labels best predicts the exercise the  
44 lifter is doing.

45 Once the pose has been identified using the CNN, a 2D  
46 human pose estimation model will be used to estimate the  
47 lifter's pose and analyze their form. The pose estimation  
48 model will output 17 key points that represent the lifter's  
49 body position, discussed in section 4.2. These key points  
50 can then be used to compare the lifter's form to the ideal  
51 form for the exercise and score the lifter's form. The scoring  
52 will be based on the lifter's deviation from the ideal form  
53 and can be used to provide feedback on how to improve.

54 Overall, this paper proposes the use of a Deep CNN and  
55 2D human pose estimation to analyze and score the form  
56 of lifters. The proposed system has the potential to provide  
57 feedback and guidance on how to improve form in a more  
58 efficient and cost-effective way than traditional methods.

### 2. Background

#### 2.1. Image Classification

61 Convolutional Neural Networks (CNNs) are a type of  
62 deep learning neural network that are mainly used for im-  
63 age classification. They are composed of multiple layers of  
64 neurons, each of which is responsible for a certain task. The  
65 neurons in the first layer of a CNN are responsible for de-  
66 tecting the lowest level features in an image, such as edges,  
67 corners and basic shapes. As the image is passed through  
68 each layer, more complex features such as parts of objects  
69 are detected. The last layer of the network is used to classify  
70 the image based on the features it has detected.

71 ResNet [3] is a type of CNN that uses residual learning  
72 to ease the training of networks that are substantially deeper  
73 than those used previously. In a traditional CNN, each layer

<sup>\*</sup>University of Rochester, Rochester, NY 14627, USA

<sup>†</sup>equal contribution

<sup>‡</sup>corresponding author: amart50@u.rochester.edu

The code for this project can be found here: <https://github.com/alexmartin1722/liftingpose>

74 is responsible for learning a certain set of features. This can  
 75 become difficult when the network is too deep, as the layers  
 76 may not be able to learn the desired features. In a ResNet,  
 77 instead of learning a completely new set of features in each  
 78 layer, each layer is responsible for learning a set of residual  
 79 functions. These residual functions represent the difference  
 80 between the desired output and the predicted output from  
 81 the layer before it. By learning these residual functions, the  
 82 network can better optimize the weights in each layer.

## 83 2.2. Pose Estimation

84 Pose estimation is an important focus in the computer  
 85 vision community due to its large range of real-world ap-  
 86 plications. Pose estimation aims to automatically predict  
 87 and track human posture by localizing joints and defining  
 88 limb orientation. Although numerous methods have been  
 89 developed to address this challenge, most of them rely on  
 90 complex and hand-crafted models, which are expensive to  
 91 train and require a large amount of data. Recent works  
 92 have shown that plain vision transformers [6] can be used to  
 93 achieve excellent performance in visual recognition tasks,  
 94 however, their potential for pose estimation has not been  
 95 fully explored.

## 96 2.3. Pose Datasets



Figure 1. An example annotation from MPII

97 Most datasets annotate 17 key points, while others ex-  
 98 pand their annotations up to 133 (COCO Whole Body [4]).  
 99 The most relevant dataset for our use case is the OCHu-  
 100 man dataset [7]. This dataset includes annotations of bod-  
 101 ies that are occluded by objects or the frame, which is per-  
 102 fect for weightlifting in a commercial setting because you  
 103 cannot guarantee an object-free filming environment. Most  
 104 times, lifters are slightly covered by squat racks, safety bars,  
 105 benches, and large weights. It is important to be able to  
 106 predict and still provide feedback for joints that are not per-  
 107 fectly pictured in the frame, so the OCHuman is a great  
 108 benchmark to test the pose estimation model on. The model  
 109 used for pose estimation in this paper (ViTPose) is trained  
 110 on this dataset for the very reasons discussed above.

## 111 2.4. Models

### 112 2.4.1 ViTPose

113 ViTPose [5] is a pose estimation model consisting of plain  
 114 and non-hierarchical vision transformers, where the back-  
 115 bones are pre-trained with masked image modeling pretext  
 116 tasks.

117 It adopts a simple architecture, with de-convolution and  
 118 prediction layers, and decoders without skip-connections  
 119 and cross-attentions. Given a person instance image, ViT-  
 120 Pose embeds the image into tokens via path embedding.  
 121 After this, the embedded tokens are processed by several  
 122 layers consisting of self-attention layers and feed-forward  
 123 networks. [5]

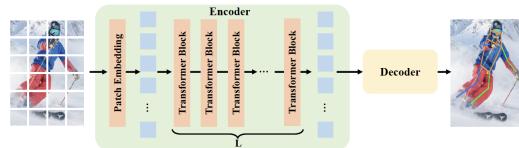


Figure 2. ViTPose Framework

124 ViTPose outperforms representative methods on the MS  
 125 COCO keypoint detection benchmark. The largest model  
 126 sets a new state-of-the-art, 80.9 AP on the MS COCO test  
 127 set as well as high performance on OCHuman with 91.6  
 128 AP.

## 129 3. Dataset

130 For this paper, we collected a new dataset for classify-  
 131 ing and estimating pose on lifts. Using Python's Beau-  
 132 tiful Soup API we web-scraped 100 photos each for barbell  
 133 bench press, barbell squat, and barbell deadlift. After gath-  
 134 ering the data and checking that they are all suitable for clas-  
 135 sification and pose estimation, we used OpenCV to trans-  
 136 form our data. Keeping a copy of the original data, first,  
 137 we flipped the data horizontally, then we proceeded to ran-  
 138 domly rotate our data between 30 degrees and 180 degrees.  
 139 Consequently, we had 900 photos as our final data set, 300  
 140 of each lift.

141 To score a person's lift against a lifter with good form,  
 142 we annotate gold data that we've found of lifting coaches  
 143 and athletes performing the exercise. We then annotate each  
 144 lift in the same style as COCO body annotations. We've  
 145 provided three images (one at each angle [front, left, right])  
 146 for each lift that we classify in our CNN in section 4.1.

## 147 4. Framework

148 To be able to score a lifter's form, we perform three  
 149 tasks. First, we classify the lift into three categories, squat,  
 150 bench, and deadlift (section 4.1). Second, we perform pose

151 estimation on the lift to get the joint locations of the lifter  
 152 [4.2](#). After we have the lifter’s keypoint data, we then score  
 153 the lifter’s form against the three gold annotations for that  
 154 lift (section [4.3](#)).

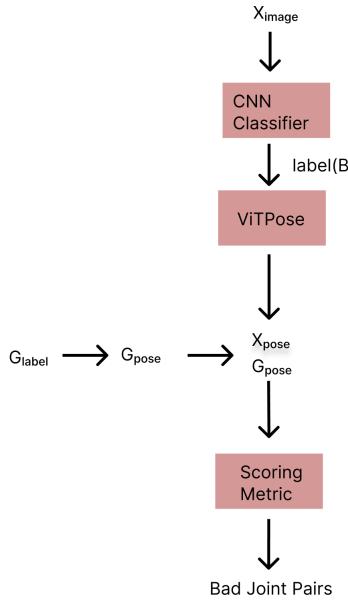


Figure 3. Framework  
 (the input image is fed to the CNN to decide the type of the lift, then annotated images are obtained using ViTPose which are compared against gold standards to provide feedback)

## 155 4.1. CNN

156 To classify the lift types, we attempt to implement two  
 157 architectures for a CNN. One is a CNN based off the ar-  
 158 chitecture from project one, and the other is based on the  
 159 ResNet [3] architecture. The CNN used in project one con-  
 160 sists of three convolutional layers with RELU activations.  
 161 The loss criteria were changed from BCE (binary cross-  
 162 entropy loss) to cross-entropy loss, given we have three lifts  
 163 for the classification.

## 164 4.2. ViTPose

165 For pose estimation, we rely on the newest state-of-the-  
 166 art model for human pose estimation, ViTPose [5]. ViTPose  
 167 is the perfect model for this because of its high throughput  
 168 and AP on the OCHuman dataset.

169 To use this model, we configure it to annotate 17 key  
 170 points on a person: nose, eye, ear, shoulder, elbow, wrist,  
 171 hip, knee, ankle<sup>1</sup>. For the full pose estimation, we employ  
 172 a top-down pose estimation method, utilizing the ViT back-  
 173 bone from the ViTAE transformer [6]. A top-down process  
 174 is a two-step method, first object detection is used to predict

<sup>1</sup>All are annotated left/right except for nose

175 where the target is and a bounding box is annotated around  
 176 the target. Once the target is predicted, the pose estimation  
 177 is performed on the target.

## 178 4.3. Scoring Metric

179 The scoring metric used in this paper utilizes a technique  
 180 from natural language processing called inter-annotator  
 181 agreement (IAA). One of the limitations of the CNN classi-  
 182 fication is that if we were to filter each lift by its angle (or  
 183 range of angles) human annotators would have to spend lots  
 184 of time classifying the lift in a range of angles. This also  
 185 would increase the complexity of the CNN and reduce its  
 186 overall accuracy in classifying the lift. To solve this prob-  
 187 lem we use a scoring metric on each of the gold annotations  
 188 and utilize annotator agreement to decide if the feedback  
 189 should be passed along to the user as seen in Algorithm 1.  
 190 One should note that if a joint is not annotated by ViTPose,  
 191 those joints are ignored in the scoring metric.

---

### Algorithm 1 Scoring Metric

---

```

for each critical pair in keypoints do
  Get the Euclid distance between two joints
   $distance = \sqrt{(\frac{x_2 - x_1}{width})^2 + (\frac{y_2 - y_1}{height})^2}$ 
  for each gold annotation do
    Compare each ratio with gold ratios
    if outside of size threshold then
      lift is bad form
    end if
  end for
  if Lift agrees with front then
    consider lift a front lift and good form
  end if
  if left angle and right angle agree then
    consider lift a side lift good form
  end if
  if  $\frac{2}{3}$  lifts call it bad form then
    report those joints as bad form
  end if
end for
  
```

---

192 The critical pairs of joints vary for each lift. In bench  
 193 press, the critical joint relations are (wrist, wrist), (wrist,  
 194 elbow), (elbow, elbow), and (elbow, shoulder). For dead-  
 195 lifts, the critical joints are (wrist, wrist), (hip, hip), (ankle,  
 196 ankle), and (knee, knee). For squat, the critical joints are  
 197 (wrist, wrist), (knee, knee), (ankle, ankle), and (hip, knee).

## 198 5. Results

### 199 5.1. CNN

200 For the CNN pose classification, we find that our CNN  
 201 framework performs well, while our ResNet implemen-

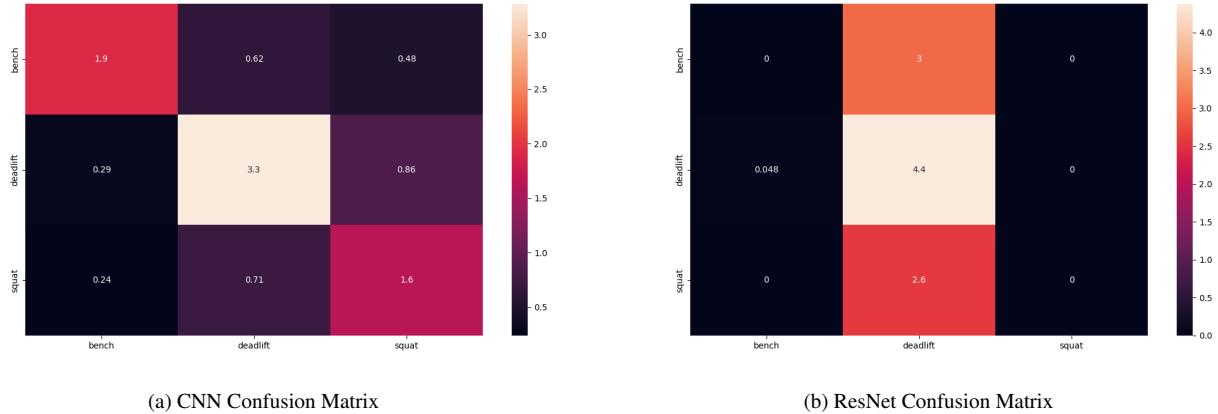


Figure 4. Confusion Matrix for CNN Implementations on the validation set (true labels on the x-axis, predicted labels on the y-axis)

Model	Bench	Squat	Deadlift
CNN	<b>63.5%</b>	63.0%	<b>74.2%</b>
ResNet	0%	<b>100%</b>	0%

Table 1. CNN Classification vs ResNet Classification Accuracy  
 (The table shows the percentage of lifts classified correctly for each of the three lifts)

202 tion struggles. In Figure 4a and Table 1 we see that our clas-  
203 sification method is successful in differentiating the types  
204 of lifts and provides good results. Our ResNet implemen-  
205 tation is obviously incorrect as it can only classify the deadlift  
206 pose as seen in Figure 4b.

207 In our implementation, it can be seen in Figure 4 that  
208 the classification between deadlift and squat encounters the  
209 most confusion due to the visual similarity and setup of the  
210 two lifts.

It should also be noted that we tried many implementations of CNN architectures, as well as prebuilt networks and networks with pre-trained weights that we fine-tuned on our task. Even with these methods, we were not able to achieve high classification scores for the squat lift and most models tend to classify it as a deadlift.

## 5.2. ViTPose

218 The pose estimation results we get are very good. Since  
219 we do not have a large amount of gold data, we are not able  
220 to quantify the accuracy. However, to ensure satisfactory  
221 performance, we went through each image by hand. In an  
222 ideal experiment, we would supply this to Mechanical Turk-  
223 ers to note the points that are out of position, however, we  
224 lack the funding and time for this. Instead, we went through  
225 each image and noted the points that were off.

For most lifts, ViTPose does a great job with body coordinates and tends to fail on the ear and nose estimations.

### 5.3. Scoring Metric

We find the scoring metric successful on the 15 images we took ourselves. It is able to note the regions of poor form with high success. The only time an annotation is scored poorly when it might be good is in the case of a ViTPose annotation with low confidence. This is covered in more depth in Appendix B.

## 6. Conclusion

In conclusion, our CNN and ViTPose implementations are successful in pose classification and pose estimation, respectively. Our CNN framework is able to differentiate between different types of lifts and our ResNet implementation struggles due to the visual similarity and setup of the two lifts. Our scoring metric is able to accurately note the regions of poor form with high success. The ViTPose annotations are found to have a high success rate. However, there are some inaccuracies in cases of partial limb visibility or when a full limb is not visible. Overall, our implementations are successful and show promise for future projects.

## 7. Future Work

To further improve the performance of this model, it would be important to classify lifts by the angle of the camera in a more critical way. The current method ignores this constraint and tries to work around it with a unique scoring metric.

259 It would also be beneficial to annotate a dataset of lifting  
260 poses for other use cases and expand the model beyond the  
261 bench, squat, and deadlift. When any exercise does not involve  
262 advanced bio-mechanics, it would be possible to implement this  
263 method on that exercise. Additionally, more work needs to be done on  
264 making a more rigorous classifier because our classifier does poorly in differentiating between  
265 squats and deadlifts, given their visual similarities.

## 267 A. Use Cases

### 268 A.1. Qualifying Olympic Lifts

269 With the new bench press rules being introduced by the  
270 International Powerlifting Federation (IPF) in 2023, qualifying  
271 bench press has gotten much more difficult [1]. The  
272 rule introduced mandates lifters to get their elbows parallel  
273 with the top of the shoulder joint to prevent lifters from  
274 getting away with smaller ranges of motion.

## 275 B. Annotated Pose Estimations

276 In Figures 5-7, we can see that ViTPose is successful on  
277 our own lift examples. It has high success on the limbs in  
278 the frame and even limbs occluded by our example lifter.  
279 However, It struggles with the face annotation and in Figure 5,  
280 you can see a poor annotation of the right wrist and elbow  
281 because the full limb is out of frame. It also is less  
282 successful on the ankles in Figure 6 that are cut out of the  
283 frame.

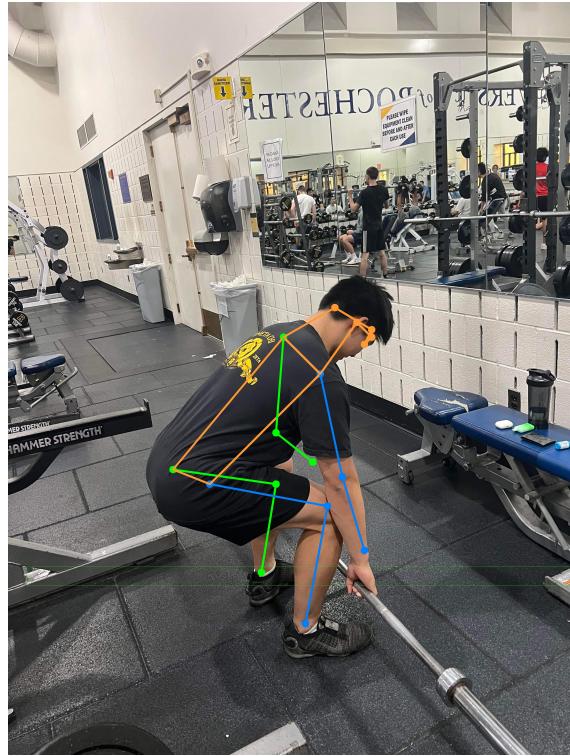


Figure 5. Deadlift Pose Estimation

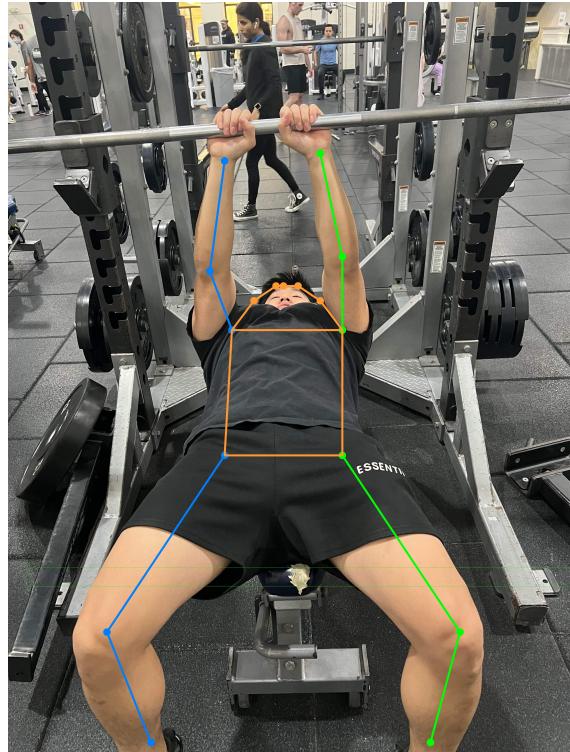


Figure 6. Bench Pose Estimation

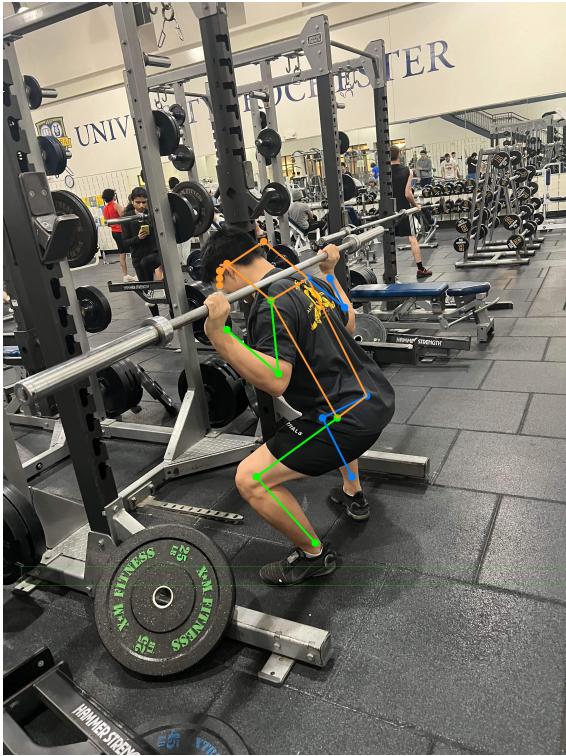


Figure 7. Squat Pose Estimation

## 284 References

- 285 [1] Jake Dickson. International powerlifting federation unveils  
286 bench press rule change for 2023. 5
- 287 [2] FIFA. Semi-automated offside technology to be used at fifa  
288 world cup 2022. 1
- 289 [3] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.  
290 Deep residual learning for image recognition, 2015. 1, 3
- 291 [4] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays,  
292 Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence  
293 Zitnick. Microsoft coco: Common objects in context. In  
294 *European conference on computer vision*, pages 740–755.  
295 Springer, 2014. 2
- 296 [5] Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao. Vit-  
297 pose: Simple vision transformer baselines for human pose es-  
298 timation. *arXiv preprint arXiv:2204.12484*, 2022. 2, 3
- 299 [6] Yufei Xu, Qiming Zhang, Jing Zhang, and Dacheng Tao. Vi-  
300 tiae: Vision transformer advanced by exploring intrinsic induc-  
301 tive bias, 2021. 2, 3
- 302 [7] Song-Hai Zhang, RUILONG Li, Xin Dong, Paul L. Rosin, Zixi  
303 Cai, Han Xi, Dingcheng Yang, Hao-Zhi Huang, and Shi-Min  
304 Hu. Pose2seg: Detection free human instance segmentation.  
305 2018. 2