

This is a self-archived version of an original article. This version may differ from the original in pagination and typographic details.

Author(s): Cronin, Neil; Rantalainen, Timo; Ahtiainen, Juha; Hynynen, Esa; Waller, Benjamin

Title: Markerless 2D kinematic analysis of underwater running : A deep learning approach

Year: 2019

Version: Accepted version (Final draft)

Copyright: © 2019 Elsevier Ltd.

Rights: CC BY-NC-ND 4.0

Rights url: <https://creativecommons.org/licenses/by-nc-nd/4.0/>

Please cite the original version:

Cronin, N., Rantalainen, T., Ahtiainen, J., Hynynen, E., & Waller, B. (2019). Markerless 2D kinematic analysis of underwater running : A deep learning approach. *Journal of Biomechanics*, 87, 75-82. <https://doi.org/10.1016/j.jbiomech.2019.02.021>

1 **Markerless 2D Kinematic Analysis of Underwater Running: A Deep Learning Approach**

2 Neil Cronin¹, Timo Rantalainen¹, Juha Ahtiainen¹, Esa Hynynen², Ben Waller^{1,3}

3 ¹ Faculty of Sport and Health Sciences, University of Jyväskylä, Finland

4 ² KIHU- Research Institute for Olympic Sports, Jyväskylä, Finland

5 ³ Physical Activity, Physical Education, Sport and Health Research Centre (PAPESH), Sports

6 Science Department, School of Science and Engineering, Reykjavik University, Reykjavik, Iceland

7

8 Original Article

9

10 Correspondence:

11 Neil Cronin

12 University of Jyväskylä, Neuromuscular Research Center, Faculty of Sport and Health Sciences, P.

13 O. Box 35, FI-40014, University of Jyväskylä, Finland

14 tel: +358 40 805 3735

15 e-mail: neil.j.cronin@jyu.fi

16

17 Keywords: deepwater running, kinematics, deep learning, artificial intelligence, motion analysis

18

19 Word count: 2339

20

21

22

23

24

25

26 Abstract

27

28 Kinematic analysis is often performed with a camera system combined with reflective markers
29 placed over bony landmarks. This method is restrictive (and often expensive), and limits the ability
30 to perform analyses outside of the lab. In the present study, we used a markerless deep learning-
31 based method to perform 2D kinematic analysis of deepwater running, a task that poses several
32 challenges to image processing methods. A single GoPro camera recorded sagittal plane lower limb
33 motion. A deep neural network was trained using data from 17 individuals, and then used to predict
34 the locations of markers that approximated joint centres. We found that 300-400 labelled images
35 were sufficient to train the network to be able to position joint markers with an accuracy similar to
36 that of a human labeler (mean difference <3 pixels, around 1cm). This level of accuracy is sufficient
37 for many 2D applications, such as sports biomechanics, coaching/training, and rehabilitation. The
38 method was sensitive enough to differentiate between closely-spaced running cadences (45-85
39 strides per minute in increments of 5). We also found high test-retest reliability of mean stride data,
40 with between-session correlation coefficients of 0.90-0.97. Our approach represents a low-cost,
41 adaptable solution for kinematic analysis, and could easily be modified for use in other movements
42 and settings. Using additional cameras, this approach could also be used to perform 3D analyses.
43 The method presented here may have broad applications in different fields, for example by enabling
44 markerless motion analysis to be performed during rehabilitation, training or even competition
45 environments.

46

47 Introduction

48

49 Kinematic analysis is used to characterise changes in joint angles during human movement. This
50 information can be combined with other sources, e.g. force data, to build a more complete picture of

51 how a movement is performed (Winter, 1991), and thus has important implications for various
52 fields such as sports biomechanics, injury risk assessment and rehabilitation (see Colyer et al. 2018
53 for a review). Kinematic analysis is often performed with a camera system combined with a set of
54 reflective markers placed over bony landmarks, allowing a digital model of the moving person to be
55 reconstructed (van der Kruk and Reijne, 2018). However, the use of reflective markers can restrict
56 the settings in which data can realistically be collected, and many existing camera-based methods
57 still rely on expensive hardware and software. Moreover, in an aquatic environment, the use of
58 markers is impractical because they impede normal movement and are prone to significant motion
59 artifact.

60
61 Recently, several attempts have been made to develop markerless methods, which in theory could
62 be used outside of the laboratory and allow movement to be analysed in more natural, unconstrained
63 conditions (see Drory, Li, and Hartley 2017 for a comprehensive overview). In particular, methods
64 that rely on artificial intelligence have demonstrated promising results (see Colyer et al., 2018 for
65 review), and have the potential to revolutionise the way movement analysis is performed due to
66 their powerful ability to ‘learn’ patterns in data. In the present study, we used DeepLabCut
67 (Insafutdinov et al., 2016; Mathis et al., 2018; Pishchulin et al., 2015) to track the locations of
68 (approximated) lower limb joint centres and used this information to perform 2D kinematic analysis
69 of deepwater running, a task that poses several challenges to image processing methods, such as
70 poor contrast and changes in light intensity. DeepLabCut is an open-source method that combines a
71 residual neural network (ResNet-50) pretrained on ImageNet with deep convolutional and
72 deconvolutional neural network layers (Insafutdinov et al., 2016) to predict the ‘learned’ locations
73 of individual points in an image using feature detectors (He et al. 2015). The network ‘learns’
74 marker locations by being trained on labeled data, which consists of individual images accompanied
75 by a human-defined label of the ‘correct’ marker location. During training, the weights are adjusted

76 iteratively so that for each image, the network assigns high probabilities to target marker locations
77 and low probabilities to all other regions. Training thus allows the network to ‘learn’ feature
78 detectors for each user-defined marker, rather than relying on hard-coded, pre-defined features.
79
80 In this study we demonstrate that a modified version of the DeepLabCut method can be used for
81 accurate 2D kinematic analysis of deepwater running filmed using a single GoPro camera. We used
82 this method to determine lower limb segment lengths and joint angles, and we present various other
83 parameters that could be useful in motion analysis applications.

84

85 Methods

86

87 **Participants.** A total of 21 individuals (age: 24 ± 4 years, height: 177 ± 10 cm, mass 67 ± 9 ; 13 males
88 and 8 females) volunteered to participate and provided written informed consent. The study was
89 approved by the University’s ethics committee, and testing was conducted in accordance with the
90 most recent Helsinki declaration.

91

92 **Experimental protocol.** Participants performed bouts of deepwater running whilst immersed to
93 shoulder level, and were tethered to the edge of the pool by a non-elastic cable attached to a
94 buoyancy aid (Aquawallgym©, Hungary). A single GoPro camera (Hero 3 model) was enclosed in
95 a waterproof case and positioned underwater in the sagittal plane to the participants’ left side at a
96 distance of approximately 5 m. A custom-made calibration frame (2m x 2m) was used to calibrate
97 the field of view for each participant and test. The camera was then set to record at 60Hz whilst
98 participants ‘ran’ at different cadences controlled by a metronome (increased by 5 strides per
99 minutes (spm) from 45 to 90 spm). A subset of participants were tested a second time
100 approximately 1 week after the first test, to enable test-retest comparisons to be performed. A deep

neural network was trained and then used to predict the locations of several markers that approximated joint centres. Predicted joint coordinates were used to determine lower limb segment lengths and joint angles from the left leg, which was closest to the camera.

Deep Neural Network. The method used here largely followed the method described by Mathis et al. (Mathis et al., 2018; v1). We first trained the network using 500 images from 17 randomly chosen participants (i.e. 28-30 images per participant), leaving aside data from the remaining 4 participants (see below). The training images were randomly selected using a custom-written script in Matlab (Mathworks, v2016b). These images were cropped (dimensions: 580 x 480 pixels) and then manually labelled, with markers placed on the lateral side of the trunk (approximately mid-way between the shoulder and hip), greater trochanter, lateral femoral condyle, lateral malleolus, and 5th meta-tarsal head. The labelled images were used to train a deep neural network with a 90% training, 10% test split. The ResNet model was initialised with weights trained on ImageNet (He et al., 2015), and the cross-entropy loss between the predicted score-map and the ground truth score-map was minimised using stochastic gradient descent (Insafutdinov et al., 2016). The network was trained for 200,000 iterations using a single Tesla K80 GPU via Microsoft Azure's cloud platform running Python (Python Software Foundation; v.3.5) and Tensorflow (Abadi et al. 2018; v.1.2.1). The training process was repeated with smaller training sets (400, 300, 200 and 100 images respectively), to determine the minimum number of images required to reach satisfactory predictive performance for this task. The number of frames used for training was selected based on previous work using a similar method (Mathis et al., 2018), and for each trained model, frames were randomly assigned to the test or training set.

Evaluation of deep neural network performance. To compare between joint coordinates labelled by a human and those labelled by the network, pairwise Euclidean distances were computed for

each marker location (root mean square error: RMSE). In the results section the RMSE values are shown for individual joints or as the average across all joints, as appropriate. To quantify the evolution of the training error, and to enable training to be resumed later if needed, the Tensorflow weights were stored every 10,000 iterations. As noted above, data from 4 randomly chosen participants were excluded completely from the training set. After training of the neural networks was complete, videos from these 4 participants were evaluated by each neural network model, thereby serving as additional test data. This approach was chosen to enable out of sample predictions that were completely independent of the training process, thus giving some indication of the generalisability of our trained models.

Determining joint angles and segment lengths. Segment lengths were initially computed in terms of pixels, using the coordinate data of each point exported during the analysis of each video. Segment lengths were calculated based on the distance formula: $d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$, where d = segment length in pixels, and x and y values denote the coordinates of the two points that make up a segment. A scaling factor was calculated for each participant and trial based on the corresponding calibration frame video, and used to scale segment lengths. Joint angles at the hip, knee and ankle were determined using the atan2 function in Matlab. In several cases, the neural network did not attempt to place a marker because the target joint location was blocked by the hand or moved beyond the image field of view. To overcome the effect of missed (and misplaced) markers on the resulting kinematic and segment length data, raw data were first filtered with a median filter (10-20 data points generally yielded good results) followed by a Butterworth 4th order low-pass filter (Figure 3). In some cases, e.g. when a marker was missing for several consecutive frames, it was necessary to experiment with different filtering procedures.

Results

151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175

Deep Neural Network Performance

Using the full training set of 500 labelled images, the mean training error across all images was 1.4 pixels. The mean test error was 2.92 pixels (approximately 1 cm). This model represents the best performance achieved out of all of the tested models. As seen in Figure 1A, training performance was similar between all of the tested models after 200,000 iterations. Test performance, i.e. how well the network predicts marker locations on images it has not ‘seen’ during training, was similar between models trained on 300-500 images, suggesting that 300 training images was sufficient for this task. However, test performance clearly decreased with training datasets of 100-200 images, indicating overfitting of these models during the training stage. For all models, training time varied between approximately 9-12 hours.

*** FIGURE 1 HERE ***

For the better performing models, both training and test errors were largely independent of which marker was being tracked, whereas for the poorer performing models, the disparity between different markers was much larger (Figure 1B and C). It should be noted that the trunk marker was not placed by the network in around 20% of frames due to the hand blocking the target area. In some images, the 5th metatarsal marker also was not placed because the foot moved beyond the camera field of view. However, these issues did not substantially affect kinematic tracking (Figures 3 and 4).

*** FIGURE 2 HERE ***

176 Figure 2 shows some examples of the same images labelled by the 100 and 500 models. In some
177 cases, the models make similar predictions, compared to each other and to a human labeller (e.g.
178 middle image in Figure 2). In other images, the 100 model consistently makes larger errors, and in
179 the most extreme cases, identifies an ostensibly correct location but on the wrong limb (left image
180 in Figure 2), which largely accounts for the bigger test errors of the poorer performing models (for
181 further model comparisons see Supplementary Video 1). Segment length calculations yielded
182 consistent traces across consecutive stride cycles, particularly for models trained on more images.
183 Segment lengths varied somewhat throughout a stride (see Figure S1), due to minor fluctuations in
184 marker locations, as well as the inevitable 3D rotation of the lower limb that cannot be quantified
185 with this method.

186

187 Kinematics of deepwater running

188

189 Using this method we obtained consistent joint angle traces over several consecutive stride cycles.
190 Figure 3 shows examples of data computed from three different 10s videos obtained from different
191 individuals whose data were not seen by the neural network during training. These videos were
192 processed entirely by a trained neural network, and did not require a human labeler at any stage (see
193 also Supplementary Videos 2-4).

194

195 *** FIGURE 3 HERE ***

196

197 Based on visual identification from the videos, it was possible to approximate the start of individual
198 stride cycles. Figure 4 shows the results of this segmentation for a single 20s trial from a participant
199 whose data were not seen by the neural network during training.

200

201 *** FIGURE 4 HERE ***

202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226

For the data shown in Figure 4, the mean range of motion at the hip joint was $96.7 \pm 5.4^\circ$. Corresponding values for the knee and ankle joint were $124.0 \pm 8.2^\circ$ and $59.0 \pm 6.3^\circ$ respectively. Similar range of motion values were also obtained for the 3 participants' data in Figure 3 (hip: 102.2 - 121.7° ; knee: 102.0 - 133.0° ; ankle: 67.2 - 78.2°). To demonstrate some additional applications of our method, we used it to examine kinematics at a range of different cadences, to ensure that the method was sufficiently robust to small changes in movement velocity and the resulting kinematics. These results are shown in Figure 5, demonstrating that small changes in cadence of 5 spm can be distinguished reliably based on the kinematic traces.

*** FIGURE 5 HERE ***

We also performed test-retest comparisons on data collected from the same individual one week apart (Figure 6). Figure 6 shows a segment of data (~30s), as well as individual strides from each session. Based on the mean strides, the range of motion values for tests 1 and 2 were 107.9° and 102.1° (hip), 115.2° and 121.9° (knee), and 35.5° and 33.8° (ankle) respectively. Corresponding mean differentials of these traces were 0.004 and 0.011 (hip), -0.014 and -0.011 (knee), and -0.062 and 0.003 (ankle). Correlation coefficients computed on the pairwise mean stride data showed values of 0.97, 0.90 and 0.93 using the raw data, and 0.93, 0.78 and 0.79 when computed on the differential of the mean stride data, for the hip, knee and ankle respectively.

*** FIGURE 6 HERE ***

Discussion

227 In this paper we demonstrate the ability to perform markerless 2D kinematic tracking using a deep
228 residual neural network trained on human-labelled data. Our results show that 300-400 labelled
229 images were sufficient to train the network to be able to position joint markers with an accuracy
230 similar to that of a human labeler (with a mean difference of around 1cm). This level of accuracy is
231 sufficient for many 2D applications, such as sports biomechanics and coaching, and
232 rehabilitation/training scenarios. Moreover, it is likely that network performance could be further
233 improved, for example by using a deeper pre-trained network or by modifying model
234 hyperparameters (Mathis et al., 2018). In addition to assessing joint kinematics, we also
235 demonstrate the ability to compute relevant parameters such as joint range of motion and cadence
236 on a stride by stride basis, and show strong test-retest reliability of kinematics measured with this
237 method.

238

239 The kinematic results obtained in this study are largely comparable to those of the few previous
240 studies conducted in this area. For example, our joint range of motion results (Figures 3 and 4) are
241 similar to values reported by Kato et al. (2001) and Kilding et al. (2007). For some participants we
242 observed larger hip range of motion than in the Kato study, but this is likely due to the
243 unconstrained nature of deepwater running, compared to running on a treadmill in Kato's study.
244 Similarly, we observed less peak knee flexion than Kilding et al., likely due to differences in
245 deepwater running technique (high knees versus cross-country technique). At all joints, the
246 kinematic traces in our study were qualitatively similar to those observed in overground running
247 (e.g. Voloshina and Ferris 2015).

248

249 We applied Deep Learning to a task that is very challenging from an image processing perspective.
250 For example, the light intensity of the image background varied between (and even within) tests due
251 to the fact that data were collected in a public swimming pool. The camera used in the present study

252 had automatic shutter speed, and due to the low amount of light, motion blur was evident in the
253 videos, particularly in the distal portion of the image, which may have contributed to the larger
254 RMSE at the ankle than at other joints for some models (Figure 1). This could conceivably have
255 increased errors in marker placement by both the human and the neural networks. As light is filtered
256 by the water, the contrast of the videos also seemed to be low. With these constraints, image quality
257 was arguably quite low, further exacerbated by isolating and cropping individual video frames
258 during training. It seems likely that using a more advanced camera could have improved overall
259 image quality and thereby minimised tracking errors (for human and neural network labelling).
260 However, we see this as a strength of the present approach, since it highlights the robustness of the
261 method in spite of the factors mentioned above.

262

263 Aside from issues related to filming underwater, we also encountered difficulties common to gait
264 analysis such as an arm blocking a marker's position. Additional cameras were not necessary to
265 overcome this issue, and a simple filtering procedure combined with a robust deep neural network
266 was sufficient to produce consistent kinematic results. Nonetheless, implementing this method in
267 3D may help to reduce the effect of marker occlusion, due to the redundancy provided by additional
268 cameras (see Drory et al., 2017 for a similar approach based on single images). Other difficulties
269 included the occasional placement of a marker on the wrong limb by the neural network. To
270 overcome this issue, other studies have used information about spatial relations between markers
271 (e.g. the hip is always an approximately constant distance from the knee) to better inform
272 predictions (Drory et al., 2017), and these techniques could have helped to improve accuracy in the
273 present study. It should also be noted that camera-based methods are not the only possible solution
274 for kinematic analysis. Some studies have used inertial measurement units (Dadashi et al., 2012),
275 with the advantage that cameras are not needed, and so the issue of placing markers is avoided
276 completely. Finally, we only compared neural network performance to that of a human labeler, and

could not evaluate our method against traditional systems that use reflective markers. However, reflective markers in the image would influence neural network performance during training, with a high risk that the network would simply learn to identify the reflective markers, and subsequently fail when used to predict marker positions in images where reflective markers are not present. Overcoming this issue thus requires an alternative approach.

The approach used here offers a very low-cost, adaptable solution for simple kinematic analysis. The method only requires a small amount of manual labelling of image frames, and in the best case, this training process only needs to be performed once. The successfully trained network can then be used to label new videos quickly (45s for a 10s video on a standard CPU), and near real-time tracking is also possible with GPU support (Nath et al., 2018). Given the challenges associated with imaging deepwater running, it is likely that this approach could easily be modified to analyse kinematics in other human movements and measurement settings, simply by re-training the network using a suitable dataset. Moreover, using additional cameras, this approach could be used to perform 3D analyses (Nath et al., 2018). As stated by Coyler et al. (2018), the development of methods aided by artificial intelligence could revolutionise sports biomechanics and rehabilitation by broadening the applications of motion analysis to training or even competition environments.

Acknowledgements

The authors gratefully acknowledge the technical assistance of Markku Ruuskanen during data collection.

References

Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G.S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S.,

305 Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker,
 306 P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M.,
 307 Wicke, M., Yu, Y., Zheng, X., Research, G., 2015. TensorFlow: Large-Scale Machine
 308 Learning on Heterogeneous Distributed Systems.
 309 Colyer, S.L., Evans, M., Cosker, D.P., Salo, A.I.T., 2018. A Review of the Evolution of Vision-
 310 Based Motion Analysis and the Integration of Advanced Computer Vision Methods Towards
 311 Developing a Markerless System. *Sport. Med. - Open* 4, 24. doi:10.1186/s40798-018-0139-y
 312 Dadashi, F., Crettenand, F., Millet, G.P., Aminian, K., 2012. Front-Crawl Instantaneous Velocity
 313 Estimation Using a Wearable Inertial Measurement Unit. *Sensors* 12, 12927–12939.
 314 doi:10.3390/s121012927
 315 Drory, A., Li, H., Hartley, R., 2017. A learning-based markerless approach for full-body kinematics
 316 estimation in-natura from a single image. *J. Biomech.* 55, 1–10.
 317 doi:10.1016/J.JBIOMECH.2017.01.028
 318 He, K., Zhang, X., Ren, S., Sun, J., 2015. Deep Residual Learning for Image Recognition. *arXiv*.
 319 Insafutdinov, E., Pishchulin, L., Andres, B., Andriluka, M., Schiele, B., 2016. DeeperCut: A
 320 Deeper, Stronger, and Faster Multi-Person Pose Estimation Model. *arXiv*.
 321 Kato, T., Onishi, S., Kitagawa, K., 2001. Kinematical Analysis of Underwater Walking and
 322 Running. *Sport. Med. Train. Rehabil.* 10, 165–182. doi:10.1080/10578310210396
 323 Kilding, A.E., Scott, M.A., Mullineaux, D.R., 2007. A Kinematic Comparison of Deep Water
 324 Running and Overground Running in Endurance Runners. *J. Strength Cond. Res.* 21, 476.
 325 doi:10.1519/R-17975.1
 326 Mathis, A., Mamidanna, P., Cury, K.M., Abe, T., Murthy, V.N., Mathis, M.W., Bethge, M., 2018.
 327 DeepLabCut: markerless pose estimation of user-defined body parts with deep learning. *Nat.*
 328 *Neurosci.* 21, 1281–1289. doi:10.1038/s41593-018-0209-y
 329 Nath, T., Mathis, A., Chen, A.C., Patel, A., Bethge, M., Mathis, M.W., 2018. Using DeepLabCut
 330 for 3D markerless pose estimation across species and behaviors. *bioRxiv* 476531.
 331 doi:10.1101/476531
 332 Pishchulin, L., Insafutdinov, E., Tang, S., Andres, B., Andriluka, M., Gehler, P., Schiele, B., 2015.
 333 DeepCut: Joint Subset Partition and Labeling for Multi Person Pose Estimation. *arXiv*.
 334 van der Kruk, E., Reijne, M.M., 2018. Accuracy of human motion capture systems for sport
 335 applications; state-of-the-art review. *Eur. J. Sport Sci.* 18, 806–819.
 336 doi:10.1080/17461391.2018.1463397
 337 Voloshina, A.S., Ferris, D.P., 2015. Biomechanics and energetics of running on uneven terrain. *J.*
 338 *Exp. Biol.* 218, 711–719. doi:10.1242/jeb.106518
 339 Winter, D.A., 1991. The biomechanics and motor control of human gait : normal, elderly and
 340 pathological. University of Waterloo Press.
 341
 342