Short: Deep Learning Approach to Skeletal Performance Evaluation of Physical Therapy Exercises

Bhanu Garg University of California, San Diego San Diego, California, USA bgarg@ucsd.edu

Pamela Cosman University of California, San Diego San Diego, California, USA pcosman@eng.ucsd.edu

ABSTRACT

At-home exercising strongly predicts physical therapy patient outcomes, underscoring the need for analyzing patient behaviors athome via remote patient monitoring. Contemporary methods for remote patient monitoring rely on specialized sensors, i.e., Inertial Measurement Units, RGB-Depth cameras, motion capture systems, or stereo vision which are costly and not scalable to all physical therapy patients. Here, we observe a lack of literature using only a monocular RGB camera. In this paper, we demonstrate a skeletal feedback model for at-home exercises using only video acquired from a smartphone camera. We propose models for (i) Patient Performance Evaluation - which classifies the correctness of exercises, and (ii) Guidance - which identifies why the exercise went wrong so the patient can correct themselves. We use these models on our dataset of four common physical therapy exercises labeled by a physical therapist. Our results demonstrate the feasibility of using skeletal data from state-of-the-art 3D human pose estimation models for physical rehabilitation exercise evaluation and guidance. Thus, we enable remote patient monitoring and guidance from a single camera - making it highly cost-effective and scalable.

KEYWORDS

PT rehabilitation, Patient Performance Evaluation, posture correction, deep learning, self attention, home monitoring

ACM Reference Format:

Bhanu Garg, Alexander Postlmayr, Pamela Cosman, and Sujit Dey. 2023. Short: Deep Learning Approach to Skeletal Performance Evaluation of Physical Therapy Exercises . In ACM/IEEE International Conference on Connected Health: Applications, Systems and Engineering Technologies (CHASE '23), June 21–23, 2023, Orlando, FL, USA. ACM, New York, NY, USA, 5 pages. https://doi.org/10.1145/3580252.3586984

CHASE '23, June 21-23, 2023, Orlando, FL, USA

© 2023 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-0102-3/23/06.

https://doi.org/10.1145/3580252.3586984

Alexander Postlmayr University of California, San Diego San Diego, California, USA apostlma@ucsd.edu

Sujit Dey University of California, San Diego San Diego, California, USA dey@eng.ucsd.edu

1 INTRODUCTION

The field of human action evaluation (HAE) is broad in its applications, ranging from gait analysis [27] to judging Olympic performance [22]. Physical rehabilitation is one of many applications which can significantly benefit from using HAE technologies [3]. Using skeleton features in HAE was shown to be promising in [9]; however, these studies are limited to low fidelity exercises, i.e., large amplitude movements.

Using specialized sensors, such as Inertial Measurement Units (IMUs) [17], RGB-D cameras [14],[19],[31],[20], or motion capture systems [6] for HAE have shown promising results in displaying accurate assessment and quantification of rehabilitation and strength exercises. While these specialized sensor-based assessment technologies provide high accuracy, they are limited in applicability due to the inherent cost and complex nature of obtaining specialized hardware for action evaluation. Moreover, standardizing skeleton data and deep feature representation methods from the sensors is another key issue in developing reliable quality assessment algorithms for HAE [12]. Thus, relying on specialized hardware limits practical applications of HAE and prevents the creation of standard datasets required to advance the technology.

Recent advances in 3D Human Pose Estimation (HPE) such as [21], [16] have allowed for the feature extraction of skeletal key points from a monocular RGB camera. These 3D HPEs are popularly used to predict key joint positions on the body [11]. In this paper, we show that classical machine learning methods such as Dynamic Time Warping (DTW) are limited for monocular RGB HAE due to the inherent noise associated with predicting the depth dimension from a monocular RGB camera. Using 2D skeletal information in conjunction with deep learning has shown promising results for regression analysis [13]. In HAE for physical therapy (PT), getting 3D skeletons is necessary to clinically assess range of motion in key joints. Furthermore, physical rehabilitation requires corrective feedback for improving patient outcomes [8] [10]. Although monocular RGB 3D HPE is promising, we observe no HAE with high-fidelity feedback in the literature.

In our work, we propose a framework for skeleton-based HAE for PT exercises from a single camera, that evaluates patient repetitions as correct or not, and offers explainability for correcting the incorrect movements. For this paper, we narrow the scope to human Patient Performance Evaluation and guidance, assuming 3D skeletal features are obtained from a 3D human pose estimator.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

CHASE '23, June 21-23, 2023, Orlando, FL, USA

2 RELATED WORK

We review various approaches proposed for general HAE and for PT applications, and highlight the major differences between them and our work. Although specialized sensor-based systems for physical rehabilitation are popular in the literature [15], we observe limited approaches based on a single RGB camera. To our knowledge, there are no 3D skeletal feature, monocular-based approaches. The literature highlights the importance of 3D skeletal features as a constituent to HAE in physical therapy; therefore, our work complements existing work that quantifies physiotherapy metrics such as range-of-motion and joint angle-based success criteria implicitly, by learning from data.

2.1 HAE based on complex sensors

Rooted in classical signal processing [4], matching techniques have seen some success in RGB-D applications [29]. Another classical approach [7] was used as a screening tool by performing anomaly detection on several activities of daily living (ADL). This system allows for health service provider intervention for given neuromusculoskeletal conditions but does not address rehabilitation. Moreover, we observe classical analytical methods are insufficient for datasets with low signal-to-noise ratio, i.e., collected using monocular RGB with skeletal features engineered using state-of-the-art HPE models.

Deep learning approaches have been popular for regression analysis [24] [30] on RGB-D, optical tracking system [14], and IMU [26] datasets. These regression models do not provide any patient feedback. The ability to use deep learning to diagnose and track the progression of Alzheimer's Disease was shown in [31].

2.2 HAE based on a single RGB camera

Outside of physical therapy, approaches based on a single RGB camera for HAE have been shown by [13], which utilized 2D skeletal features extracted from monocular RGB to train deep learning regression models on the UNLV Olympic Diving and MIT Olympic Ice Skating Scoring datasets. Pseudo3D was used to extract spatiotemporal features from video on the UNLV Diving dataset [5]. The feature engineering here [23] is limited to pseudo representations of human pose estimation.

For physical rehabilitation, [25] uses a regression-based quantitative scoring model trained on the KIMORE dataset using monocular RGB. A lack of accurate determination of joint angles from 2D skeleton data limits the extent of evaluation of exercises. For example, indicators based on range of motion of joints, a prevalent metric used in PT, cannot be modeled (even implicitly) by 2D joint data. Feedback on static exercises only was shown in [18]; these comprise a small fraction of rehabilitative exercises.

3 METHODOLOGY

We begin by defining key terms and notation in Section 3.1. In Section 3.2, we explain our data collection method, while in Section 3.3, we build the case to show that classical methods such as Dynamic Time Warping are not suitable for evaluating repetitions and providing guidance results for data collected using monocular RGB camera videos. We then describe our proposed methods in Sections 3.4 and 3.5.



Figure 1: Process Flow: (A) Patient records a video performing an exercise, (B) each frame is skeletonized, (C) the skeletons are transformed to an angular domain and post-processed, (D) the resulting vector is fed to PPE and Guidance models, (E) patient receives feedback.

3.1 Terminology

We focus on HAE using a monocular RGB camera which is available in most smartphones and tablet computers. For simplicity, any data is from a monocular RGB camera unless otherwise specified.

For Patient Performance Evaluation (PPE), we formulate and define the following tasks in our processing and data pipeline.

- Skeletonization: Given an RGB frame containing a human, skeletonization refers to the process of extracting the 3D (x,y,z) coordinates of key joint positions and constructing a stick figure-like model of the human. The term skeleton refers to the constructed stick figure. This paper uses a custom skeletonizer that gives results similar to [21].
- **Segmentation:** Given a sequence of skeletons from the video feed, segmentation is the process of extracting exercise repetitions from the entire exercise session. Each exercise repetition has start and stop frame numbers in the video. We use **r** to denote a repetition.
- **Patient Performance Evaluation (PPE):** Given an exercise repetition, PPE assigns the exercise repetition a label from the set {*good*, *bad*}. Ground truth labels are determined by an expert, i.e., a licensed physical therapist.
- **Guidance:** For a repetition that has been tagged *bad*, guidance is the suggestion to correct the form of that repetition. For example, while doing a deadlift, possible guidance could be "keep back straight" if the patient has a rounded back.

Fig. 1 shows the overall flow of our proposed process. Given the abundance of literature on human pose estimation and segmentation, we narrow our attention to PPE and guidance methods.

3.2 Feature Engineering

We record student volunteers performing exercises on phone cameras. We provide more details on data collection in Section 4.1. We use a custom skeletonization model, which outputs the 3D positions of 17 key joints as in [21]. The videos are then skeletonized and segmented, and the individual repetitions are fed as input to the subsequent steps..

We begin by selecting joint angles relevant to the exercise E as determined by an expert. We note that the selected joint angles a_j could be those that move during the exercise as well as those that



Figure 2: The distributions of similarity scores produced by DTW on good (green) and bad (blue) repetitions are highly overlapping.

might be required to stay stationary. Let $\mathcal{A}_E = \{a_j\}$ be the set of such selected joint angles for exercise *E*.

To formalize, for a repetition **r** belonging to exercise *E* of time length *T* video frames, we skeletonize to obtain a tensor of shape (T, 17, 3), for 17 joints and 3 spatial dimensions. For each frame $i \in [1, ..., T]$, we compute the angles for the joints in \mathcal{A}_E . Let $v_{(j,i)}$ be the angle for joint a_j at time *i*. The 2D tensor $\mathbf{V} = [[v_{(j,i)}]]$ of dimensions $(|\mathcal{A}_E|, T)$ is then smoothed and subsampled as described in the paragraph below to get $\tilde{\mathbf{V}}$. This $\tilde{\mathbf{V}}$ is used as an input to the models described in Sections 3.4.1 and 3.4.2.

Smoothing filter: We use an averaging filter with window size 5 on the time series of individual joint angles to smooth out outliers.

Subsampling: As different subjects perform exercise repetitions at different rates, the V vary in width. This can cause problems in training deep neural networks including difficulties in batch training and data preprocessing. Based on our discussions with the physical therapist, the exercise rate is seldom useful in classifying a repetition as *good* or *bad*. More importance is given to the form, which does not depend on the rate. Hence, we subsample each repetition to 20 equidistant frames.

3.3 Classical Methods: Dynamic Time Warping

Dynamic Time Warping (DTW) is a method for measuring similarity between two temporal sequences. It was used successfully in the context of physical therapy exercises in [29] where data was collected with a Microsoft Kinect Camera. In our experiments, DTW could be applied to patient repetitions to compare them with a "ground truth" coming from our physical therapist consultant. In concept, it could allow the cumulative difference across all key angles to be used as a general threshold for PPE while using anglespecific thresholds to provide feedback on any particularly incorrect angles. Despite its efficacy with RGB-D based datasets, we observed poor results on our 3D HPE skeletal data. Fig. 2 plots the distribution of similarity scores per joint angle from DTW for good and bad repetitions; the substantial overlap between the two distributions makes the thresholding-based approach infeasible. Thus, we turn to deep learning-based methods that are more potent in finding fine patterns in the data.

3.4 Patient Performance Evaluation (PPE)

In this section, we describe our neural network-based approach to PPE. We first describe the data processing steps and then the model architectures.

3.4.1 CNN-based PPE. In this section, we describe the convolutional neural network (CNN) architecture used for classifying \tilde{V}





Figure 3: CNN based PPE: Temporal Convolution followed by Spatial Convolution, Avg Pooling and Classification Head



Figure 4: Guidance Process: Autoencoder trained on good reps learns the intrinsic features of exercises. Reconstruction error is used to get feedback cues.

to the PPE label set $Y = \{good, bad\}$. In the first layer, we apply a convolution filter on the time axis (temporal convolution), i.e., a kernel of shape (3, 1) and output channels 32. Then the next layer is convolution on the angles (spatial convolution), i.e., a kernel of shape $(1, |\mathcal{A}_E|)$ and output channels 64. After each of these two convolutions, a ReLU non-linearity is followed by a batch norm. Finally, the feature map from the last convolution is average-pooled, followed by a binary classification head. The intuition for choosing the above architecture is that the temporal convolution learns the temporal context at time t_i of the angle a_j . Once the context of the angle has been set, the model learns the interdependence relations among the joint angles. To the best of our knowledge, ours is the first method to view human Patient Performance Evaluation as a Spatio-Temporal convolution network.

3.4.2 Attention based PPE. In this section, we describe the architecture of the attention-based transformer classifier used for PPE. Attention networks have gained popularity by learning to attend to values at different timestamps and have been successful in natural language processing and computer vision [28]. The number of heads is chosen to be 2, with a 512 embedding dimension for the encoder. A classification head follows the encoding head.

3.5 Guidance and Exercise Feedback

This section proposes our method for generating guidance and feedback mechanisms for the *bad* repetitions. An autoencoder is a model that encodes the input into a feature map and then uses the encoding to decode it back to the input. It is used to train an

underlying model distribution of the data. With this context, we first train an autoencoder using only the *good* repetitions to make the neural network learn a robust *good* template of the exercise. The mean and variance $(\mu_{a_j}, \sigma_{a_j})$ of the reconstruction error for each joint angle are collected. Then for a given *bad* repetition, the reconstruction error e_{a_j} from the autoencoder is computed. The angles with the highest *z*-score: $z_{a_j} = \frac{e_{a_j} - \mu_{a_j}}{\sigma_{a_j}}$ are selected as the ones responsible for the repetition being *bad*. The process is shown visually in Fig. 4.

AutoEncoder: Following similar arguments on the efficacy of attention-based models on time series data as described in Section 3.4, we propose to use an attention based autoencoder. While implementing a corresponding CNN-based auto-encoder, we found maintaining the convolutional upsampling layers to consistently output the same size vector as the input for different exercises non-trivial, and adding unnecessary complexity to our model pipeline. In the attention-based autoencoder model, we use the default number of dimensions for the embeddings, with the number of heads equal to $|\mathcal{A}_E|$, and mean square error (MSE) as the loss function.

4 EXPERIMENTS

4.1 Data

In this section, we discuss the details of our dataset, which we use to demonstrate the efficacy of our proposed PPE and guidance models. As existing datasets [14] [2] [1] do not have corrective feedback and classification scores, we were required to create our own.

4.1.1 *General Setup.* Based on consultation with a physical therapist, we selected four exercises that are frequently prescribed in physical therapy, involve compound movements, and collectively cover multiple focus areas (shoulder, legs, hips). The exercises are: double leg Romanian dead-lift (DoubleRDL), single leg Romanian dead-lift (SingleRDL), single leg mini squat (SingleMS), and rotator cuff (RotatorCuff).

Students working on the project volunteered to be recorded performing exercises. We selected 10 subjects for our data collection. Before recording each exercise, the subjects were shown a demonstration video created by the physical therapist that shows several repetitions of the exercise with proper form. The subjects were asked to recreate 10 repetitions seen in the video to the best of their ability. For unilateral exercises, subjects were asked to perform 5 repetitions on each side. After the first 10 repetitions, subjects were informed about common mistakes seen in physical rehabilitation. For example, a DoubleRDL is often performed with the subject's spine being too rounded with too much scapular protraction or the subject's knees being too bent or locked out. These subtle mistakes have a profound impact on muscle activation during the exercise. The subjects were instructed to incorporate these incorrect postures into an additional 10 repetitions.

To test our proposed methods' robustness to different viewing angles and camera types, we recorded subjects using five different smartphones (iPhone Models X, 11, 7, 8, and Google Pixel 3A) placed on tripods approximately 4 feet high and at five different angles (30, 60, 90, 120, and 150 degrees) to the subject. To summarize, data

Exercise	Evaluation (F1 Scores)			Guidance (top-2)	
	Baseline	Ours		Baseline	Ours
	DTW	CNN	Attention	DTW	Attention
DoubleRDL	0.197	0.262	0.255	0.846	0.884
SingleRDL	0.052	0.222	0.274	0.652	0.679
SingleMS	0.435	0.544	0.519	0.805	0.858
RotatorCuff	0.569	0.78	0.769	0.652	0.758

Table 1: Results: F1 scores of the PPE models, and top-2 accuracy of the Guidance models. The Baseline used is the DTW method proposed in [29].

consists of 1000 repetitions for each exercise from 10 subjects, and 5 camera angles.

4.1.2 Labeling. For each video recording, each exercise was segmented manually into repetitions which were classified as being good or bad, and the two most erroneous movements that needed correction were noted. Our final dataset includes 19% good repetitions for DoubleRDL, 15% for SingleRDL, 42% for SingleMS, and 50% for RotatorCuff. This variation is due to the fact that subjects were more knowledgeable about proper exercise form for some exercises compared to others.

4.2 Results

We split the 10 subjects into 3 folds (with 4, 3, and 3 subjects), and train the model on 2 of the folds, and evaluate the model on the remaining fold. We show the mean F1 score of the 3-fold cross validation.

Patient Performance Evaluation For PPE, we use negative log likelihood loss for training. We use the *Adam* optimizer for training with learning rate 1e-4 and default weight parameters $\beta = (0.9, 0.98)$, and batch size 16. We report the F1 score of the binary classification model. We see that both the CNN and Attention-based models consistently outperform DTW (see Table 1). Between the CNN and Attention-based models, the Attention-based model performs slightly better than the CNN-based model.

Our results show promise that our novel solutions can classify exercises based on single camera input.

Guidance We map each joint angle to an action item that the subject could do to correct that angle based on the output of the autoencoder. For guidance models, we report the top - 2 accuracy. Specifically, it is considered correct if either of the two most erroneous predicted angles match the PT judgment. Human body movements have constraints and the joint angles do not operate in isolation from one another. Correcting one angle could affect the correctness of the other angle - for example, in DoubleRDL, making the back straight would still be conducive to correcting locked knees. Hence, the top-2 accuracy metric is justified.

For training the Attention autoencoder, we use the Adam optimizer with learning rate 1e - 3 and default weight parameters $\beta = (0.9, 0.98)$ with batch size 32. In Table 1, we see that the Attention-based autoencoder method outperforms DTW. Ours is the first deep learning-based approach to PT guidance.

5 CONCLUSION AND FUTURE WORK

We propose deep-learning Patient Performance Evaluation and guidance models for PT rehabilitation and at-home monitoring using only a smartphone camera. Our CNN and Attention networks Deep Learning for Physical Therapy Evaluation

Exercise	Guidance Criteria	
DoubleRDL	Knees Too Bent, Knees Locked, Back Too Round,	
	Feet Too Far Apart	
SingleRDL	Knees Too Bent, Knee Locked, Back Too Round,	
	Leg Not In-Line	
SingleMS	Hips Not Level, Squat Too Low, Twisting Torso	
RotatorCuff	Twisting Torso, Arm Too Extended, Lifting Arm	
	Too High/Low	

Table 2: Guidance cues defined by PT

show improved results on PPE and guidance for all four exercises. We observe the importance of a robust dataset as a priority for future work. In contrast, DTW would not benefit from such a dataset. To the best of our knowledge, this is the first deep learning-based approach for PT guidance. In the future, we plan to run clinical trials on the effectiveness of our method on patient outcomes and expand our exercise library to many more exercises.

REFERENCES

- Massimo Camplani, Adeline Paiement, L Tao, Sion Hannuna, Dima Damen (Aldamen), Majid Mirmehdi, and Tilo Burghardt. 2014. Depth video and skeleton of people walking up stairs.
- [2] Marianna Capecci, Maria Gabriella Ceravolo, Francesco Ferracuti, Sabrina Iarlori, Andrea Monteriu, Luca Romeo, and Federica Verdini. 2019. The KIMORE Dataset: KInematic Assessment of MOvement and Clinical Scores for Remote Monitoring of Physical REhabilitation. *IEEE transactions on neural systems and rehabilitation engineering: a publication of the IEEE Engineering in Medicine and Biology Society* 27, 7 (July 2019), 1436–1448. https://doi.org/10.1109/TNSRE.2019.2923060
- [3] Oya Celiktutan, Ceyhun Burak Akgul, Christian Wolf, and Bülent Sankur. 2013. Graph-Based Analysis of Physical Exercise Actions. In Proceedings of the 1st ACM International Workshop on Multimedia Indexing and Information Retrieval for Healthcare (MIIRH '13). Association for Computing Machinery, New York, NY, USA, 23–32. https://doi.org/10.1145/2505323.2505330 event-place: Barcelona, Spain.
- [4] Yinpeng Chen, Margaret Duff, Nicole Lehrer, Hari Sundaram, Jiping He, Steven Wolf, and Thanassis Rikakis. 2011. A Computational Framework for Quantitative Evaluation of Movement during Rehabilitation. *AIP Conference Proceedings* 1371 (06 2011). https://doi.org/10.1063/1.3596656
- [5] Li-Jia Dong, Hong-Bo Zhang, Qinghongya Shi, Qing Lei, Ji-Xiang Du, and Shangce Gao. 2021. Learning and fusing multiple hidden substages for action quality assessment. *Knowledge-Based Systems* 229 (2021), 107388. https://doi.org/10. 1016/j.knosys.2021.107388
- [6] Chen Du, Sarah Graham, Colin Depp, and Truong Nguyen. 2021. Assessing Physical Rehabilitation Exercises using Graph Convolutional Network with Self-supervised regularization. In 2021 43rd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC). 281–285. https: //doi.org/10.1109/EMBC46164.2021.9629569
- [7] Amr Elkholy, Mohamed E. Hussein, Walid Gomaa, Dima Damen, and Emmanuel Saba. 2020. Efficient and Robust Skeleton-Based Quality Assessment and Abnormality Detection in Human Action Performance. *IEEE Journal of Biomedical* and Health Informatics 24, 1 (2020), 280–291. https://doi.org/10.1109/JBHI.2019. 2904321
- [8] Marc P. Gruner, Nathan Hogaboom, Ike Hasley, Jared Hoffman, Karina Gonzalez-Carta, Andrea L. Cheville, Zhuo Li, and Jacob L. Sellon. 2021. Prospective, Singleblind, Randomized Controlled Trial to Evaluate the Effectiveness of a Digital Exercise Therapy Application Compared With Conventional Physical Therapy for the Treatment of Nonoperative Knee Conditions. Archives of Rehabilitation Research and Clinical Translation 3, 4 (2021), 100151. https://doi.org/10.1016/j. arrct.2021.100151
- [9] Reza Haghighi Osgouei, David Soulsby, and Fernando Bello. 2020. Rehabilitation Exergames: Use of Motion Sensing and Machine Learning to Quantify Exercise Performance in Healthy Volunteers. *JMIR Rehabil Assist Technol* 7, 2 (Aug. 2020), e17289. https://doi.org/10.2196/17289
- [10] Stephanie Hewitt, Ruth Sephton, and Gillian Yeowell. 2020. The Effectiveness of Digital Health Interventions in the Management of Musculoskeletal Conditions: Systematic Literature Review. *Journal of medical Internet research* 22 (June 2020), e15617. https://doi.org/10.2196/15617
- [11] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. 2014. Human3.6M: Large Scale Datasets and Predictive Methods for 3D Human Sensing in Natural Environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36, 7 (jul 2014), 1325–1339.
- [12] Qing Lei, Ji-Xiang Du, Hong-Bo Zhang, Shuang Ye, and Duan-Sheng Chen. 2019. A Survey of Vision-Based Human Action Evaluation Methods. Sensors (Basel) 19,

19 (Sept. 2019).

- [13] Qing Lei, Hong-Bo Zhang, Ji-Xiang Du, Tsung-Chih Hsiao, and Chih-Cheng Chen. 2020. Learning Effective Skeletal Representations on RGB Video for Fine-Grained Human Action Quality Assessment. *Electronics* 9, 4 (2020). https: //doi.org/10.3390/electronics9040568
- [14] Yalin Liao, Aleksandar Vakanski, and Min Xian. 2020. A Deep Learning Framework for Assessing Physical Rehabilitation Exercises. *IEEE Transactions on Neural Systems and Rehabilitation Engineering* 28, 2 (2020), 468–477. https: //doi.org/10.1109/TNSRE.2020.2966249
- [15] Yalin Liao, Aleksandar Vakanski, Min Xian, David Paul, and Russell Baker. 2020. A review of computational approaches for evaluation of rehabilitation exercises. *Computers in Biology and Medicine* 119 (2020), 103687. https://doi.org/10.1016/j. compbiomed.2020.103687
- [16] Dushyant Mehta, Srinath Sridhar, Oleksandr Sotnychenko, Helge Rhodin, Mohammad Shafiei, Hans-Peter Seidel, Weipeng Xu, Dan Casas, and Christian Theobalt. 2017. VNect: Real-time 3D Human Pose Estimation with a Single RGB Camera. In ACM Transactions on Graphics, Vol. 36. https://doi.org/10.1145/3072959.3073596 Issue: 4.
- [17] Slobodan Milanko and Shubham Jain. 2020. LiftRight: Quantifying strength training performance using a wearable sensor. *Smart Health* 16 (2020), 100115. https://doi.org/10.1016/j.smhl.2020.100115
- [18] Cristian Militaru, Maria-Denisa Militaru, and Kuderna-Iulian Benta. 2020. Physical Exercise Form Correction Using Neural Networks. In Companion Publication of the 2020 International Conference on Multimodal Interaction (Virtual Event, Netherlands) (ICMI '20 Companion). Association for Computing Machinery, New York, NY, USA, 240–244. https://doi.org/10.1145/3395035.3425302
- [19] Ferda Ofli, Gregorij Kurillo, Štěpán Obdržálek, Ruzena Bajcsy, Holly Brugge Jimison, and Misha Pavel. 2015. Design and Evaluation of an Interactive Exercise Coaching System for Older Adults: Lessons Learned. IEEE J Biomed Health Inform 20, 1 (Jan. 2015), 201–212.
- [20] Paritosh Parmar and Brendan Tran Morris. 2016. Measuring the quality of exercises. In 2016 38th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC). 2241–2244. https://doi.org/10.1109/EMBC. 2016.7591175
- [21] Dario Pavllo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 2019. 3D Human Pose Estimation in Video With Temporal Convolutions and Semi-Supervised Training. 7745–7754. https://doi.org/10.1109/CVPR.2019.00794
- [22] Hamed Pirsiavash, Antonio Torralba, and Carl Vondrick. 2014. Assessing the Quality of Actions. https://doi.org/10.1007/978-3-319-10599-4_36
- [23] Zhaofan Qiu, Ting Yao, and Tao Mei. 2017. Learning Spatio-Temporal Representation with Pseudo-3D Residual Networks. In *IEEE International Conference on Computer Vision (ICCV)* (ieee international conference on computer vision (iccv) ed.). https://www.microsoft.com/en-us/research/publication/learning-spatiotemporal-representation-pseudo-3d-residual-networks/
- [24] Noureddin Sadawi, Alina Miron, Waidah Ismail, Hafez Hussain, and Crina Grosan. 2019. Gesture Correctness Estimation with Deep Neural Networks and Rough Path Descriptors. In 2019 International Conference on Data Mining Workshops (ICDMW). 595–602. https://doi.org/10.1109/ICDMW.2019.00090
- [25] Faegheh Sardari, Adeline Paiement, Sion Hannuna, and Majid Mirmehdi. 2020. VI-Net–View-Invariant Quality of Human Movement Assessment. Sensors 20, 18 (2020). https://doi.org/10.3390/s20185258
- [26] Dapeng Tang. 2020. Hybridized Hierarchical Deep Convolutional Neural Network for Sports Rehabilitation Exercises. *IEEE Access* 8 (2020), 118969–118977. https: //doi.org/10.1109/ACCESS.2020.3005189
- [27] Lili Tao, Adeline Paiement, Dima Damen, Majid Mirmehdi, Sion Hannuna, Massimo Camplani, Tilo Burghardt, and Ian Craddock. 2016. A comparative study of pose representation and dynamics modelling for online motion quality assessment. *Computer Vision and Image Understanding* 148 (2016), 136–152. https://doi.org/10.1016/j.cviu.2015.11.016
- [28] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. CoRR abs/1706.03762 (2017). arXiv:1706.03762 http://arxiv.org/abs/ 1706.03762
- [29] Wenchuan Wei, Yao Lu, Catherine D. Printz, and Sujit Dey. 2015. Motion Data Alignment and Real-Time Guidance in Cloud-Based Virtual Training System. In Proceedings of the Conference on Wireless Health (Bethesda, Maryland) (WH '15). Association for Computing Machinery, New York, NY, USA, Article 13, 8 pages. https://doi.org/10.1145/2811780.2811952
- [30] Bruce X. B. Yu, Yan Liu, and Keith C. C. Chan. 2020. Skeleton-Based Detection of Abnormalities in Human Actions Using Graph Convolutional Networks. In 2020 Second International Conference on Transdisciplinary AI (TransAI). 131–137. https://doi.org/10.1109/TransAI49837.2020.00030
- [31] Bruce X. B. Yu, Yan Liu, Keith C. C. Chan, Qintai Yang, and Xiaoying Wang. 2021. Skeleton-based human action evaluation using graph convolutional network for monitoring Alzheimer's progression. *Pattern Recognition* 119 (2021), 108095. https://doi.org/10.1016/j.patcog.2021.108095