# Skeleton-based Approaches based on Machine Vision: A Survey

Jie Li, *Member, IEEE,* Binglin Li, and Min Gao

*Abstract*—**Recently, skeleton-based approaches have achieved rapid progress on the basis of great success in skeleton representation. Plenty of researches focus on solving specific problems according to skeleton features. Some skeleton-based approaches have been mentioned in several overviews on object detection as a non-essential part. Nevertheless, there has not been any thorough analysis of skeleton-based approaches attentively. Instead of describing these techniques in terms of theoretical constructs, we devote to summarizing skeleton-based approaches with regard to application fields and given tasks as comprehensively as possible. This paper is conducive to further understanding of skeleton-based application and dealing with particular issues.**

*Index Terms*—**Skeleton-based, Action recognition, Pose estimation.**

## I. INTRODUCTION

Skeleton-based approaches, also known as kinematic techniques, cover a set of joints and a group of limbs based on physiological body structure. Typically, the number of joints is determined by computational complexity, which is between ten to thirty. In recent years, some significant techniques have worked successfully on discovering the representation of skeletons with dots and edges, which can be categorized as top-down methods (e.g., cascade pyramid network [1] and PoseFix [2]) and bottom-up methods (e.g., Openpose [3]).

In order to analyze skeleton-based approaches deeply, we observe these researches from two aspects(i.e., single-frame and multi-frame). Fig. 1 describes the summarization. In regard to single-frame type, tasks are handled through investigating independent images within videos. Correspondingly, multi-frame type requires a series of sequential images among videos to explore the inherence.

### A. Single-frame Approaches

*1) Multi-view Pose Estimation:* Although skeletons with 3D structures enhance the accuracy of pose estimation, the training operation needs plenty of 3D ground-truth data which is costly. Thus, multiple 2D skeletons from different views are integrated into 3D structures under a given strategy to determine the pose. EpipolarPose [4] implements epipolar geometry to combine the 2D skeletons and trains a 3D pose estimator with camera geometry information. RepNet [5] finds the mapping from 2D to 3D skeletons by designing an adversarial training approach and adding feedback projection from 3D to 2D skeletons. Based on Wasserstein generative

J. Li and M. Gao are with Chongqing University of Science and Technology, Chongqing, China, 401331 e-mail: 2014008@cqust.edu.cn.

B. Lin is with Chongqing University, Chongqing, China, 400001

adversarial network (WGAN), RepNet generates 3D skeletons by sending 2D skeletons as input into WGAN. The produced 3D skeletons are reprojected to 2D domain by camera loss. A substitutive fusion approach [6] tracks the status of joints in two views of 2D skeletons and labels the joints as "well tracked", "inferred", or "not tracked". For a joint in different views, the one diagnosed as "well tracked" has the highest priority, and then the inferred one is better than the one of "not tracked". Furthermore, dynamic time warping with K-nearest neighbor is adopted to classify learned skeleton representation. In regard to multiple persons, a multi-way matching algorithm [7] clusters detected 2D skeletons by sticking the keypoints of a person in various views.

Two networks (named as VA-RNN and VA-CNN) work together to discover skeleton representation of actions [8], [9]. An automatic selection scheme is involved in both the nets to choose prior viewpoints of 2D skeletons rather than a fixed criterion. In VA-RNN, 2D skeletons are rotated to obtain a complex feature by training multiple LSTM. For VA-CNN, a intrinsic feature is abstracted by a convolutional network. These two features are fused to determine action classes.

Additionally, loss functions (e.g., regression loss and consistency loss.) in deep neural networks (e.g., CNN) are designed to improve pose estimation within multi-view skeletons, and distance biases of skeleton information in multiple views are minimized to improve fusion accuracy [10].

**Challenges:** Relation evaluation (e.g., similarity) between skeletons in two views is still an important issue.

*2) Object Segmentation:* Object segmentation based on skeletons aims at cutting targeted items (including biological and non-biological things) from others. In this type of tasks, object positions should be ascertained first and then object margins are drawn out.

Pose2Seg [11] focuses on human segmentation in images. In the model, detected skeletons are compared with established standard skeletons of each pose, and a affine transformation matrix calculates similarity. SegModule is presented to understand skeleton features and corresponding visualization. In addition, semantics is introduced as a supplementary perspective. By training a Part FCNto obtain semantic part score map, the body is further divided into head, torso, left arm, left leg, right arm, and right leg [12]. Another human segmentation approach [13] discovers edge information through adding all widths of physical parts (e.g., head and shoulders, neck, chest, hip) into body skeletons. In consideration of cascade connection of curve skeletons like a tree and its leaves, body segmentation are realized by trimming branches and isolating intersectional leaves [14]. In a human counting case, blurry margins are
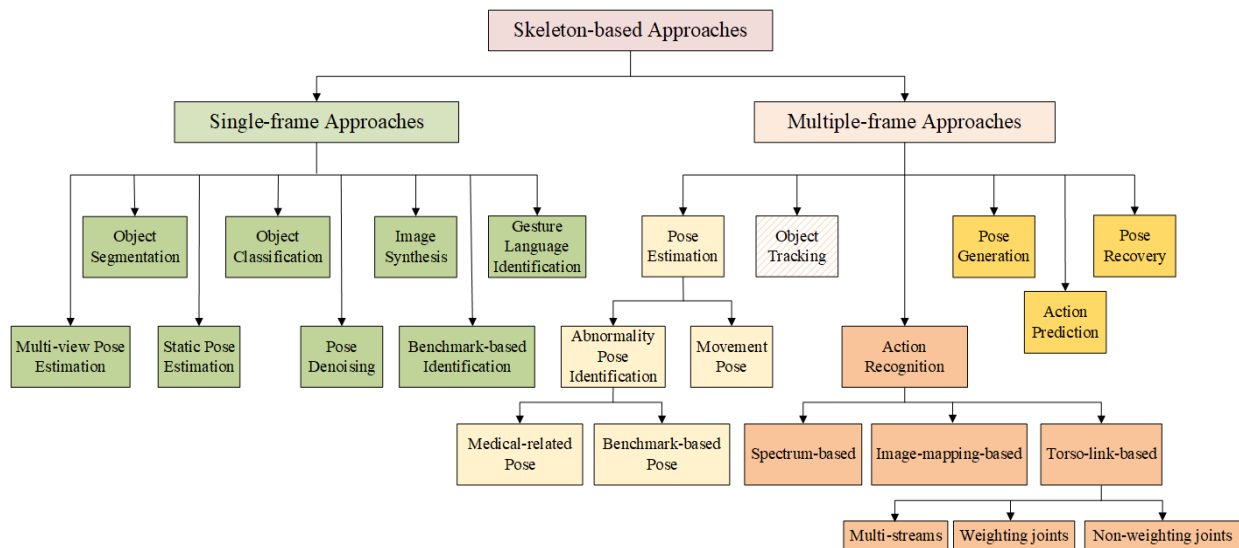
Fig. 1: Types of skeleton-based approaches.

allowed, let alone body parts. After removing background, all heads are identified within skeleton graph segmentation and utilized to gain human amount [15].

Besides of human body segmentation, animal and non-biological items are also being studied to be segmented. Critical points are selected on the basis of surface skeletons and expand to component sets separately [16]. With skeleton matching techniques, skeletal branches over a particular position are reconstructed by estimating the distances between two views and used to reconstruct the non-biological things [17].

**Challenges:** Pinpoint object margins of diverse items are confirmed not only depending on skeletons but also item structures.

*3) Static Pose Estimation:* Without comparison with any benchmark, it is hard to evaluate poses by only one image.

**(1) Single person**.

A path (i.e., 2D-3D-2D) learns the features of skeletons severely through projecting into 3D and being reprojected into 2D with lifting networks. The original 2D skeleton is as input [18]. A parsing induced learner exploits parsing information to enhance skeleton information through a pose encoder. A pose encoder abstracts pose features while another pose encoder fuses residual information into pose representation [19]. HR-Net has parallel high-to-low resolution subnetworks to gain both high and low resolutions of skeletons and sums them up [20]. Based on ConvNets, 2D images are translated into 3D skeleton models. For instance, heatmaps and silhouettes are extracted from a 2D body, and then pose and shape parameters are abstracted separately. These parameters are meshed together into a 3D body with 2D annotations [21]. In another case, depth knowledge is integrated with ConvNets between paired joints in skeletons to mix into a 3D pose [22]. PoseRefiner implements skeletons in binary channels and pose classes to train ConvNets learning likelihood heatmaps to refine the skeletons [23]. A dual-source approach learns the representations both from 2D and 3D skeletons [24]. 3D

poses are projected into multiple 2D skeletons, and used to find the highest likelihood along with a test image. To address skeletons seriously, a upper-body visualization uses different colors and polylines to distinguish between the left and right body [25]. With these representations, 16 poses are clustered with high accuracy. Furthermore, ConvNets parameters are analyzed to find prior settings to identify poses with skeletons [26].

**(2) Multiple persons**.

Cascaded Pyramid Network (CPN) containing two networks (i.e., GlobalNet and RefineNet) [1] adopts a top-down pipeline which means locating skeleton keypoints first with ResNet backbone, extracting features of these keypoints as HyperNet, and then assembling them. A multi-person pose estimation (RMPE) involves spatial transformer networks to rectify various ground truth bounding boxes for people and a final box is obtained for each person [27]. DeepCut [28] uses adapted fast R-CNN (AFR-CNN) to detect body parts with integer linear programming and dense CNN (Dense-CNN) to obtain the intensity with geometric and appearant constraints. Additionally, based on DeepCut, DeeperCut [29] is proposed to improve body part detection with a bottom-up pipeline and an image-conditioned pairwise assembling strategy is designed. Angles among body parts are observed meticulously to assist in searching pairwise joints.

**Challenges:** Without consecutive images of an action, features learned from static poses can resolve the estimation.

*4) Object Classification:* The purpose of object classification is to identify the items in images through observing skeletons without any limitations of object types.

In order to offset skeleton noise in classification, skeleton contours are trimmed as a trade-off between shape reconstruction error and skeleton simplicity based on Bayesian theory [30]. Tree representation of skeletons as a graph is turned into strings of skeleton edges, and a deformable contour method is used to compare these strings of diverse objects [31]. Node pairs in curve skeletons are obtained by cascade of symmetry

filters and a symmetry correspondence matrix is designed to gain symmetry cloud which is classified by spectral analysis [32].

**Challenges:** Curve skeletons are main tools for object classification with obvious transformation of each object.

*5) Pose Denoising:* Typically, skeletons are dominant for special structures, especially for human body. When odd poses occur and parts of bodies are overlapped, it is hard to peel skeletons of each body. Redundant joints and limbs are deemed as noise, which can be eliminated by filtering linear transformations [33] and comparing joint positions and limb angles with standard ones [34].

**Challenges:** Advanced denoising algorithms may be used to further improve the capability of pose denoising.

*6) Image Synthesis:* Commonly, image synthesis tries to produce other views of skeletons by learning the representation of skeletons in a view. GAN is a popular technique to solve this problem, but traditional GAN has weak ability of obtaining the relationship among joints and limbs. Bidirection GAN is adopted to search the mapping from an initial pose skeleton to another pose skeleton [35]. Deformable GANs [36] use heatmaps from skeletons as conditional information and human poses as originial images feeding into the generator to generate other pose images. The generated poses are shifted into the discriminator with corresponding heatmaps. Furthermore, CRF-RNN is established to predict conditional human body with pose transforamtion from a given skeleton to a target skeleton [37]. Mask R-CNN model is built to transform a pose skeleton into another pose skeleton, and paralleling models operate simultaneously with diverse keypoints sampled from initial pose skeletons. These mappings are combined together to creat a new pose skeleton [38]. Some researches conduct physical structure characteristics to reconstruct images other than studying the inherence. For example, body skeletons are divided into multiple components and the background is drawed out [39]. These components are rotated and integrated into a new body which is blended into a synthesized new background.

**Challenges:** GAN is a useful tool in generation although its performance is still unstable. While applying GAN on skeleton synthesis, how to improve the quality is still a challenge.

*7) Benchmark-based Identification:* Comparison with benchmark is a straight way to judge the class of items. For skeletons, key points in standard and objective skeletons are contrasted in turn.

Skeleton information (joints and limbs) is transformed into a tree with key points which are compared with given trees. The best match of two shape trees are treated as identical thing [40]. A skeleton graph of any object is also used to contrast with standard units in both high geometric and topological similarities [41]. To highlight the joints in comparisons, lines among joints address the relations of two joints instead of limbs. Unlike a fixed number of compared joints in previous two techniques, comparison steps in skeleton graphs here are random [42]. Similarity metrics along with 3D skeleton features are designed to obtain the similarity between the human skeleton in a single frame image and the templates [43]. Body dot clouds are cooperated with main curve skeletons to judge

a matching ratio of human motions, which works better than the case only involved curve skeletons [44]. Skeletons using Kinect from upper and lower body are analyzed separately to identify human gaits with ANN as a classifier [45].

**Challenges:** The evaluation criterion of similarity between the object and template skeletons is controlled by experience. Complex metrics may play well on similarity assessment.

*8) Gesture Language Identification:* Body gestures have rich information as well as language, also called as body language. Deep learning algorithm is a novel tool to understand gesture meanings, such as RNN [46] and LSTM [47]. Deep networks extract features of gesture skeletons, which is easier than those of intact bodies.

**Challenges:** Unlike full body skeletons, most body gestures only contain partial joints and limbs. Thus, application scans should be settled to ensure the involved groups of joints and limbs.

### B. Multi-frame Approaches

*1) Dynamic Pose Estimation:* **(1) Abnormality pose identification.**

**1) Medical-related pose.** A view adaptive LSTM (VA-LSTM) aims at detecting medical condition (i.e., sneeze/cough, headache, neck pain, staggering, chest pain, vomiting, falling, and back pain), containing a classification and regression subnetwork [48]. Original skeletons are rotated and translated into new architectures and then sent into the subnetworks to learn the corresponding medical classes. Neuromusculoskeletal disorders are also detected by pose estimation [49]. Asymmetry features are extracted with splitted results of body joints according to the left and right body, which is used to catch normal motion patterns by a probabilistic normalcy model. The likelihood between a test action and a normal motion is computed to determine the abnormality.

**2) Benchmark-based pose.** Deep learning structures (e.g. CNN [50] and LSTM [51]) are applied to abstract the features which are compared with a series of skeleton benchmarks (i.e., joints and limbs).

**(2) Movement pose.**

Graph convolutional network (GCN) is a core solution for movement pose estimation on account of its strong ability of capturing spatio and temporal features [52], [53], [54], [55], [56], [54]. Except GCN, deep learning structures (e.g., CNN [57], [58], RNN [59], [57], and LSTM [60]) are also key techniques for movement pose estimation through learning representations under given conditions.

Moreover, physical analysis of skeletons and probability estimation are also utilized in movement pose estimation. A exemplar-based method is explored to adjust initial estimated poses with inhomogeneous systematic bias while skeletons are defined as a simple directed graph and limbs are directed arrows. A regression function is proposed to predict the pose with rooted-mean-squared differences between templates and objective skeletons [61]. 2D keypoints extracted by pose estimators with skeletons are fused with SMPL regressors to create 3D models with accurate camera parameters [62].

Semantic representations of volume occupancy and ground plane support are helpful for distinguishing multiple persons after evaluating each single person with 3D skeletons [63]. On the strength of spatial positions and joint changes, a decision tree can quickly recognize basic action events and monitor action types under graph constraints of state transition [64].

**Challenges:** Unlike static pose estimation, dynamic pose estimation usually discusses both spatio and temporal features with a series of images. The mixture process is a key point to control the quality of final representations.

*2) Object Tracking:* PoseTrack [65], [66] follows the particular persons in videos, including multi-person pose estimation in a image, multi-person pose estimation in videos, and multi-person articulated tracking. ArtTrack [67] draws body part graphs in temporal aspect and abandons joints with loose relation by a feed-forward convolutional architecture. For robot, visual distances and lighting intensity are considered while following human using Kinect [68]. Also based on Kinect and Kalman filter, foot points are detected and followed with depth information and pairwise curve matching in 3D space from given views to a virtual bird's eye view [69]. With capturing keypoints in motion and fitting skeleton, human pose is tracked with transformed 3D models [70]. Unlike deep learning structures, images in 3D models are filtered and the likelihood is calculated between shape models deformed by skeleton poses and image data with regard to probability theory [71]. Fast movements of animals and persons are traced with non-rigid temporal deformation of 3D surface [72]. Other than tracking human, people handling objects are traced with GCNs after detecting hand joints in body skeletons [73].

**Challenges:** The relation of a given object between sequences is a crucial issue in object tracking.

*3) Action Recognition:* Action recognition is a crucial part of skeleton approaches, which has been successfully applied to a lot of real projects. Among action recognition techniques about skeletons, there are three categories: (1)Spectrum-based, (2)Image-mapping-based, and (3)Torso-link-based approaches. In torso-link-based techniques, three subtypes are involved: 1) Multiple Streams, 2) Weighting joints, and 3) Non-weighting joints.

**(1) Spectrum-based**.

Skeleton sequences are turned into color spectrum dots with ConvNets, involving a procedure with 1) joint distribution mapping, 2) spectrum coding of joint trajectories and body parts, and 3) joint velocity weighted saturation and brightness. The spectrum actions are transferred from different filmed angles and then sent to the corresponding ConvNets which separately give scores with regard to actions [74].

**(2) Image-mapping-based**.

In this area, action skeletons are tranformed into feature maps which are learned by convolutional neural networks as the main solution backbone. With a single CNN, both temporal and spatial information obtained by physical information of joints and limbs are transformed into representation images which are sent to a CNN-family model (e.g., VGG [75] and CNN-LSTM [76]) for learning features [77], [78]. Furthermore, the representation images are stretched with various sizes and fed into multiple CNNs to combine a highly representative explanation, in which these CNNs have the same structure [79] or diverse structures [80]. Feature maps in traditional three coordinate dimensions (i.e., X, Y, and Z axis) of 3D skeletons are dissembled to corresponding CNNs, resulting in multi-column representations of actions [81]. Besides of coordinate dimensions and time dimension in 3D action skeletons, color (i.e., RGB) is also deemed as a dimension to learn deep features in action recognition [82], [83]. Features of different factors of action skeletons (i.e., joint-joint distances, joint-joint orientations, joint-joint vectors, joint-line distances, line-line angles) obtained by physical computations are encoded into images and loaded into multiple CNNs to further extract features [84], [85], [86]. For a joint in skeleton sequences, potential relations are explored by position chains of a physical body and mixed with features learned from each frame by CNN for compact representations of actions [87]. Additionally, temporal consistency is deeply discovered by establishing more networks to analyze extra information in actions [88]. Thereinto, DD-net [89] is a popular framework.

Unlike learning abstract features in action images, traditionally physical and mathematical techniques are still useful in action recognition under certain conditions. Physical information of joint changes is also used to detect actions through calculating precise relations (e.g., distances and angles) among joints for each action [90], [91]. Geometrical relationship of limbs and joints draws a tree in which actions as a father node link key poses. The actions are viewed as tree nodes and features derived from those actions are as child nodes [92].

**(3) Torso-link-based**.

Torso-link (also called stick-link) is a most commonly used technique in skeleton-based approaches.

**1) Multiple Streams**. More than two aspects of representations gained from torso-link skeletons are used to learn features of actions by dependent and parallel networks, e.g., spatial and temporal [93], [94], dots and lines [95], [96], joints and time [97], [98], position and feature [99], spatial, temporal, structural, and actional [100]. The types of networks depend on action characteristics, e.g., RNN [93], RRN [93], [97], CNN [96], [98], [100], GCN [94], and LSTM [99].

**2) Weighting joints**. In other categories, all items in joints or/and limbs are given equal importance. However, in particular cases, actions have typical changes of partial joints and limbs which can represent the actions severely. Significant joints and limbs can be chosen by covariance matrix [101], filtering function on the basis of skeleton graphs [102], CNN [103], LSTM [104], multi-head attention model with itereative attention on diverse parts of a body [105], projecting skeleton angles onto a unit sphere [106], information gain with regard to position and velocity histogram from skeletons [107]. Significance sorting techniques of all joints and limbs are conducive to reducing the hardness and complexity of skeleton feature extraction in action recognition, including spatial pyramid model (SPM) [56]. Weighting the features of different parts of human body is also valid for identifying actions, e.g., bidirectional RNN [108], [109].

**3) Non-weighting joints**. This part is the main component of action recognition, where overviews are sufficient over last decades [110], [111], [112], [113], [114], [115], [116], [117],

[118], [119]. Therefore, we here introduce the sketch briefly. In the field, inherent representation of skeletons of actions is gained by two ways: physical computation and feature extraction.

For the former, based on empirical knowledge of human body, the relations among joints and limbs under a particular action is calculated with explicit equations [120], [121], [122], [123], [124], [125], [126].

For the latter, deep neural networks and other techniques devote to learning skeleton features for each action. Typically, deep neural networks have gained great success on action recognition, e.g., DNN [127], [128], [129], CNN [130], [131], [132], [133], RNN [134], [135], [136], LSTM [137], [138], [139], [140], [141], [142], [143], [144], [145], GCN [146], [147], [148], [149], [150], [151], [152].

Moreover, traditional machine learning algorithms are designed to identify actions, e.g., kNN [153], RBF [154]. HMM discovers the semantic information of actions which assists in action identification [155], [104]. Reinforcement learning is also a technique to obtain effective representation of an action [156]. Bayesian varies across different sequences of an action [157].

Additionally, probability and mathematical theories are useful in action recognition, e.g., analogical generalization and retrieval [158], screw matrices [159], gradient vector flow comparison [160], discriminative metric [161].

**Challenges:** In action recognition, the greatest challenge is to determine the start and end moment of an action. Usually, fixed time interval is adopted and leads to high deviation while actions have widely different time intervals.

*4) Action Prediction:* Unlike HMM, CRF, RNN, LSTM, CNN, a latent global network with latent long-term global information is designed to predict an action [162]. Based on the competition in GAN, two nets (i.e., I-Net and D-Net) are trained iteratively. Full and partial sequences are sent into I-Net to learn representations, separately. Afterwards, the representations are distinguished by D-Net.

**Challenges:** The evaluation of likelihood between the existed partial images and the intact sequences of an action is a key point of action prediction.

*5) Pose Generation:* FAAST [163] provides a toolkit to create animated virtual characters using natural interaction from OpenNI-compliant depth sensors. Mesh body is also produced on basis of rigid limb motions and skinning weights both for humans and animals [164], [165].

**Challenges:** Relations both inside skeletons of a pose and among the series of images need be deeply observed to generate precise poses.

*6) Pose Stripping:* Radio frequency (RF) reflections of Wifi back from environment and humans is captured for estimation poses [166]. Heatmaps in both vertical and horizontal directions are parsed with encoders and then fused with keypoint confidence maps from RGB sequences, in which human skeletons can be stripped from background.

**Challenges:** The basic assumption is that the reflections from human body and other items are disparate. This assumption is susceptible to the things with the same reflection with human body.

## II. DATASETS

We conclude top 11 datasets which have high-frequent usage for skeleton approaches.

**1.NTU RGB+D Dataset [167].** This dataset consists of 56,880 action samples containing 4 different modalities of data for each sample: 1) RGB videos 136GB, 2) depth map sequences(including Masked depth maps 83GB and Full depth maps 886GB), 3) 3D skeletal data 5.8GB, 4) Infrared videos 221 GB, Total 1.3TB. In this dataset, the resolution of RGB videos are 1920 by 1080, depth maps and IR videos are all in 512 by 424, and 3D skeletal data contains the three dimensional locations of 25 major body joints at each frame.

**2.UT-Kinect Dataset [168].** This dataset includes two separate datasets. The first dataset(3.42G) is collected using Kinect mounted on top of a humanoid robot. There are 9 action types in the humanoid robot dataset:stand up,wave,hug,point,punch,reach,throw,run,shake hands.The second dataset(3.33G) is collected using a non-humanoid robot.There are 9 action types in the non-humanoid robot dataset:ignore, pass by the robot, point at the robot, reach an object, run away, stand up, stop the robot, throw at the robot, and wave to the robot. Each dataset contains 5 parts: 1) RGB images(.jpg), the resolution is 480x640. 2) Depth images(.png) the resolution is 320x240. 3) Calibrated depth images(.png), the resolution is 320x240. 4) Sketetal joint locations (.txt). Each row contains the data of one frame, the format is: frame number, frame count, skeletonId, (x,y,z) locations of joint 1-20. 5) Labels of action sequence (.txt).

**3.Florence 3D Actions Dataset [169].** The dataset collected at the University of Florence during 2012,has been captured using a Kinect camera. It includes 9 activities: wave, drink from a bottle, answer phone, clap, tight lace, sit down, stand up, read watch, bow.During acquisition,10 subjects were asked to perform the above actions for 2 or 3 times,which resulted in a total of 215 activity samples.

**4. Kinetics dataset [170].** This dataset stems from the competition of ActivityNet Large Scale Activity Recognition Challenge 2018, which started from 2016 CVPR. The dataset provided by Deepmind Team of Google currently includes a total of 600 categories and 500 thousand video clips all from Youtube. In the collected 600 categories, each one has at least 600 videos. Each video lasts about 10 seconds. The categories are classified into three main types: 1) Interaction between humans and objects such as playing musical instruments. 2) Human interaction such as handshake, hug. 3) Sports, etc. These three main types can also be described as Person, Person-Person, Person-Object.

**5. N-UCLA Dataset [120], [171].** Northwestern-UCLA dataset(N-UCLA) was collected by three Kinect cameras, which contains 1494 sequences covering 10 action classes from 10 performers. And these 10 actions are: pick up with one hand, pick up with two hands, drop trash, walk around, sit down, stand up, donning, doffing, throw and carry. The subjects perform an action only one time in an action sequence, which contains an average of 39 frames.

**6.SBU Interaction Dataset [137].** SBU Interaction dataset was collected with Kinect. It contains 8 classes of two-person

interactions, and includes 282 skeleton sequences with 6822 frames. Each body skeleton consists of 15 joints.

**7.SYSU Dataset [129], [172].** SYSU 3D Human-Object Interaction (SYSU) dataset is collected by Kinect camera. It contains 480 skeleton clips of 12 action categories performed by 40 subjects and each clip has 20 joints.

**8.MSR Action3D Dataset (publiced by Microsoft Research Redmond).** MSR-Action3D dataset is an action dataset of depth sequences captured by a depth camera. This dataset contains twenty actions: high arm wave, horizontal arm wave, hammer, hand catch, forward punch, high throw, draw x, draw tick, draw circle, hand clap, two hand wave, side-boxing, bend, forward kick, side kick, jogging, tennis swing, tennis serve, golf swing, pick up & throw. It is created by Wanqing Li during his time at Microsoft Research Redmond.

**9.Berkeley MHAD Dataset [161].** The Berkeley Multimodal Human Action Database (MHAD) contains 11 actions performed by 7 male and 5 female subjects in the range 23-30 years of age except for one elderly subject. All the subjects performed 5 repetitions of each action, yielding about 660 action sequences which correspond to about 82 minutes of total recording time.

**10.UTD-MHAD Dataset [173].** This dataset was collected as part of our research on human action recognition using fusion of depth and inertial sensor data. For our multimodal human action dataset reported here, only one Kinect camera and one wearable inertial sensor were used. This was intentional due to the practicality or relatively non-intrusiveness aspect of using these two differing modality sensors. Both of these sensors are low cost, easy to operate, and do not require much computational power for the real-time manipulation of data generated by them. A picture of the Kinect camera can capture a color image with a resolution of 640 by 480 pixels and a 16-bit depth image with a resolution of 320 by 240 pixels. The frame rate is approximately 30 frames per second.

**11.HDM05 Dataset [174].** HDM05 dataset is a motion capture database which contains more than three hours of systematically recorded and well-documented motion capture data in the C3D as well as in the ASF/AMC data format. Furthermore, HDM05 contains for more than 70 motion classes in 10 to 50 realizations executed by various actors. The HDM05 database has been designed and set up under the direction of Meinard Müller Tido Röder, Michael Clausen, Bernhard Eberhardt, Björn Krüger, and Andreas Weber. The motion capturing has been conducted in the year 2005 at the Hochschule der Medien (HDM), Stuttgart, Germany, supervised by Bernhard Eberhardt.

## III. CONCLUSION

Skeleton-based approach as a significant part was evolving along with the blooming development of artificial intelligent applications (such as object detection, action identification, pose estimation, and so on) which had attracted high attentions. This paper observed skeleton-based approaches and categorized these techniques in accordance with target tasks rather than theoretical frameworks, which is useful for introducing this scope.

## REFERENCES

[1] Y. Chen, Z. Wang, Y. Peng, Z. Zhang, G. Yu, and J. Sun, "Cascaded pyramid network for multi-person pose estimation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7103–7112.

[2] G. Moon, J. Y. Chang, and K. M. Lee, "Posefix: Model-agnostic general human pose refinement network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7773–7781.

[3] Z. Cao, G. Hidalgo, T. Simon, S.-E. Wei, and Y. Sheikh, "Openpose: realtime multi-person 2D pose estimation using part affinity fields," *arXiv preprint arXiv:1812.08008*, 2018.

[4] M. Kocabas, S. Karagoz, and E. Akbas, "Self-supervised learning of 3D human pose using multi-view geometry," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1077–1086.

[5] B. Wandt and B. Rosenhahn, "Repnet: Weakly supervised training of an adversarial reprojection network for 3D human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7782–7791.

[6] N. A. Azis, H.-J. Choi, and Y. Iraqi, "Substitutive skeleton fusion for human action recognition," in *2015 International Conference on Big Data and Smart Computing (BigComp)*. IEEE, 2015, pp. 170–177.

[7] J. Dong, W. Jiang, Q. Huang, H. Bao, and X. Zhou, "Fast and robust multi-person 3D pose estimation from multiple views," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7792–7801.

[8] P. Zhang, C. Lan, J. Xing, W. Zeng, J. Xue, and N. Zheng, "View adaptive recurrent neural networks for high performance human action recognition from skeleton data," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2117–2126.

[9] ——, "View adaptive neural networks for high performance skeleton-based human action recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 8, pp. 1963–1978, 2019.

[10] H. Rhodin, J. Spörri, I. Katircioglu, V. Constantin, F. Meyer, E. Müller, M. Salzmann, and P. Fua, "Learning monocular 3D human pose estimation from multi-view images," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8437–8446.

[11] S.-H. Zhang, R. Li, X. Dong, P. Rosin, Z. Cai, X. Han, D. Yang, H. Huang, and S.-M. Hu, "Pose2seg: detection free human instance segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019, pp. 889–898.

[12] F. Xia, P. Wang, X. Chen, and A. L. Yuille, "Joint multi-person pose estimation and semantic part segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6769–6778.

[13] J. C. J. Junior, C. R. Jung, and S. R. Musse, "Skeleton-based human segmentation in still images," in *2012 19th IEEE International Conference on Image Processing*. IEEE, 2012, pp. 141–144.

[14] C. Lovato, U. Castellani, and A. Giachetti, "Automatic segmentation of scanned human body using curve skeleton analysis," in *International Conference on Computer Vision/Computer Graphics Collaboration Techniques and Applications*. Springer, 2009, pp. 34–45.

[15] D. Merad, K.-E. Aziz, and N. Thome, "Fast people counting using head detection from skeleton graph," in *2010 7th IEEE International Conference on Advanced Video and Signal Based Surveillance*. IEEE, 2010, pp. 233–240.

[16] D. Reniers and A. Telea, "Skeleton-based hierarchical shape segmentation," in *IEEE International Conference on Shape Modeling and Applications 2007 (SMI'07)*. IEEE, 2007, pp. 179–188.

[17] B. Durix, G. Morin, S. Chambon, C. Roudet, and L. Garnier, "Skeleton-based multiview reconstruction," in *2016 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2016, pp. 4047–4051.

[18] C.-H. Chen, A. Tyagi, A. Agrawal, D. Drover, S. Stojanov, and J. M. Rehg, "Unsupervised 3D pose estimation with geometric self-supervision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5714–5724.

[19] X. Nie, J. Feng, Y. Zuo, and S. Yan, "Human pose estimation with parsing induced learner," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2100–2108.

[20] K. Sun, B. Xiao, D. Liu, and J. Wang, "Deep high-resolution representation learning for human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5693–5703.

[21] G. Pavlakos, L. Zhu, X. Zhou, and K. Daniilidis, "Learning to estimate 3D human pose and shape from a single color image," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 459–468.

[22] G. Pavlakos, X. Zhou, and K. Daniilidis, "Ordinal depth supervision for 3D human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7307–7316.

[23] M. Fieraru, A. Khoreva, L. Pishchulin, and B. Schiele, "Learning to refine human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 205–214.

[24] H. Yasin, U. Iqbal, B. Kruger, A. Weber, and J. Gall, "A dual-source approach for 3D pose estimation from a single image," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4948–4956.

[25] S. Johnson and M. Everingham, "Learning effective human pose estimation from inaccurate annotation," in *CVPR 2011*. IEEE, 2011, pp. 1465–1472.

[26] U. Rafi, B. Leibe, J. Gall, and I. Kostrikov, "An efficient convolutional network for human pose estimation." in *BMVC*, vol. 1, 2016, p. 2.

[27] H.-S. Fang, S. Xie, Y.-W. Tai, and C. Lu, "RMPE: Regional multi-person pose estimation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2334–2343.

[28] L. Pishchulin, E. Insafutdinov, S. Tang, B. Andres, M. Andriluka, P. V. Gehler, and B. Schiele, "Deepcut: Joint subset partition and labeling for multi person pose estimation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4929–4937.

[29] E. Insafutdinov, L. Pishchulin, B. Andres, M. Andriluka, and B. Schiele, "Deepercut: A deeper, stronger, and faster multi-person pose estimation model," in *European Conference on Computer Vision*. Springer, 2016, pp. 34–50.

[30] W. Shen, X. Bai, X. Yang, and L. J. Latecki, "Skeleton pruning as trade-off between skeleton simplicity and reconstruction error," *Science China Information Sciences*, vol. 56, no. 4, pp. 1–14, 2013.

[31] L. He, C. Y. Han, B. Everding, and W. G. Wee, "Graph matching for object recognition and recovery," *Pattern Recognition*, vol. 37, no. 7, pp. 1557–1560, 2004.

[32] W. Jiang, K. Xu, Z.-Q. Cheng, and H. Zhang, "Skeleton-based intrinsic symmetry detection on point clouds," *Graphical Models*, vol. 75, no. 4, pp. 177–188, 2013.

[33] G. G. Demisse, K. Papadopoulos, D. Aouada, and B. Ottersten, "Pose encoding for robust skeleton-based action recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 188–194.

[34] L.-C. Liao, Y.-H. Yang, and L.-C. Fu, "Joint-oriented features for skeleton-based action recognition," in *2019 IEEE International Conference on Systems, Man and Cybernetics (SMC)*. IEEE, 2019, pp. 1154–1159.

[35] A. Pumarola, A. Agudo, A. Sanfeliu, and F. Moreno-Noguer, "Unsupervised person image synthesis in arbitrary poses," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8620–8628.

[36] A. Siarohin, E. Sangineto, S. Lathuilière, and N. Sebe, "Deformable GANs for pose-based human image generation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3408–3416.

[37] C. Si, W. Wang, L. Wang, and T. Tan, "Multistage adversarial losses for pose-based human image synthesis," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 118–126.

[38] I. Radosavovic, P. Dollár, R. Girshick, G. Gkioxari, and K. He, "Data distillation: Towards omni-supervised learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4119–4128.

[39] G. Balakrishnan, A. Zhao, A. V. Dalca, F. Durand, and J. Guttag, "Synthesizing images of humans in unseen poses," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8340–8348.

[40] B. Jiang, J. Tang, B. Luo, Z. Chen, and Z. Chen, "Skeleton graph matching based on a novel shape tree," in *2009 ISECS International Colloquium on Computing, Communication, Control, and Management*, vol. 4. IEEE, 2009, pp. 636–639.

[41] H. Sundar, D. Silver, N. Gagvani, and S. Dickinson, "Skeleton based shape matching and retrieval," in *2003 Shape Modeling International*. IEEE, 2003, pp. 130–139.

[42] N. T. Giang, N. Q. Tao, N. D. Dung, and N. T. The, "Skeleton based shape matching using reweighted random walks," in *2013 9th*

[43] International Conference on Information, Communications & Signal Processing. IEEE, 2013, pp. 1–5.

[43] M. Munaro, S. Ghidoni, D. T. Dizmen, and E. Menegatti, "A feature-based approach to people re-identification using skeleton keypoints," in *2014 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2014, pp. 5644–5651.

[44] A. Hajdu, C. Giamas, and J. Pitas, "Object simplification using a skeleton-based weight function," in *2007 International Symposium on Signals, Circuits and Systems*, vol. 2. IEEE, 2007, pp. 1–4.

[45] A. Sinha, K. Chakravarty, B. Bhowmick *et al.*, "Person identification using skeleton information from kinect," in *Proc. Intl. Conf. on Advances in Computer-Human Interactions*, 2013, pp. 101–108.

[46] Z. Zhang, Y. Song, and Y. Zhang, "Motion-pose recurrent neural network with instantaneous kinematic descriptor for skeleton based gesture detection and recognition," in *2017 4th IAPR Asian Conference on Pattern Recognition (ACPR)*. IEEE, 2017, pp. 764–769.

[47] X. Liu, H. Shi, X. Hong, H. Chen, D. Tao, and G. Zhao, "Hidden states exploration for 3D skeleton-based gesture recognition," in *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2019, pp. 1846–1855.

[48] J. Yin, J. Han, C. Wang, B. Zhang, and X. Zeng, "A skeleton-based action recognition system for medical condition detection," in *2019 IEEE Biomedical Circuits and Systems Conference (BioCAS)*. IEEE, 2019, pp. 1–4.

[49] A. Elkholy, M. Hussein, W. Gomaa, D. Damen, and E. Saba, "Efficient and robust skeleton-based quality assessment and abnormality detection in human action performance," *IEEE journal of biomedical and health informatics*, 2019.

[50] J. Wu, K. Wang, B. Cheng, R. Li, C. Chen, and T. Zhou, "Skeleton based fall detection with convolutional neural network," in *2019 Chinese Control And Decision Conference (CCDC)*. IEEE, 2019, pp. 5266–5271.

[51] S. Jeong, S. Kang, and I. Chun, "Human-skeleton based fall-detection method using LSTM for manufacturing industries," in *2019 34th International Technical Conference on Circuits/Systems, Computers and Communications (ITC-CSCC)*. IEEE, 2019, pp. 1–4.

[52] L. Zhao, X. Peng, Y. Tian, M. Kapadia, and D. N. Metaxas, "Semantic graph convolutional networks for 3D human pose regression," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3425–3435.

[53] H. Zhang, Y. Song, and Y. Zhang, "Graph convolutional LSTM model for skeleton-based action recognition," in *2019 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2019, pp. 412–417.

[54] X. Gao, K. Li, Y. Zhang, Q. Miao, L. Sheng, J. Xie, and J. Xu, "3D skeleton-based video action recognition by graph convolution network," in *2019 IEEE International Conference on Smart Internet of Things (SmartIoT)*. IEEE, 2019, pp. 500–501.

[55] H. Ryu, S.-h. Kim, and Y. Hwang, "Skeleton-based human action recognition using spatio-temporal geometry," in *2019 19th International Conference on Control, Automation and Systems (ICCAS 2019)*, 2019, pp. 329–332.

[56] P. Li, M. Lu, Z. Zhang, D. Shan, and Y. Yang, "A novel spatial-temporal graph for skeleton-based driver action recognition," in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*. IEEE, 2019, pp. 3243–3248.

[57] D. Pavllo, C. Feichtenhofer, D. Grangier, and M. Auli, "3D human pose estimation in video with temporal convolutions and semi-supervised training," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7753–7762.

[58] C. Dhiman, M. Saxena, and D. K. Vishwakarma, "Skeleton-based view invariant deep features for human activity recognition," in *2019 IEEE Fifth International Conference on Multimedia Big Data (BigMM)*. IEEE, 2019, pp. 225–230.

[59] X. Chen, H. Guo, G. Wang, and L. Zhang, "Motion feature augmented recurrent neural network for skeleton-based dynamic hand gesture recognition," in *2017 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2017, pp. 2881–2885.

[60] Y. Feng, L. Ma, W. Liu, and J. Luo, "Spatio-temporal video re-localization by warp LSTM," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1288–1297.

[61] W. Shen, K. Deng, X. Bai, T. Leyvand, B. Guo, and Z. Tu, "Exemplar-based human action pose correction and tagging," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2012, pp. 1784–1791.

[62] A. Arnab, C. Doersch, and A. Zisserman, "Exploiting temporal context for 3D human pose estimation in the wild," in *Proceedings of the IEEE*

*Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3395–3404.

[63] A. Zanfir, E. Marinoiu, and C. Sminchisescu, "Monocular 3D pose and shape estimation of multiple people in natural scenes-the importance of multiple scene constraints," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2148–2157.

[64] Y. Han, S.-L. Chung, J.-S. Yeh, and Q.-J. Chen, "Real-time skeleton-based indoor activity recognition," in *Proceedings of the 32nd Chinese Control Conference*. IEEE, 2013, pp. 3965–3970.

[65] M. Andriluka, U. Iqbal, E. Insafutdinov, L. Pishchulin, A. Milan, J. Gall, and B. Schiele, "Posetrack: A benchmark for human pose estimation and tracking," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5167–5176.

[66] U. Iqbal, A. Milan, and J. Gall, "Posetrack: Joint multi-person pose estimation and tracking," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2011–2020.

[67] E. Insafutdinov, M. Andriluka, L. Pishchulin, S. Tang, E. Levinkov, B. Andres, and B. Schiele, "Arttrack: Articulated multi-person tracking in the wild," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 6457–6465.

[68] S. A. Shukor, M. A. A. Rahim, and B. Ilias, "Scene parameters analysis of skeleton-based human detection for a mobile robot using kinect," in *2016 6th IEEE International Conference on Control System, Computing and Engineering (ICCSCE)*. IEEE, 2016, pp. 173–178.

[69] C.-H. Kuo, S.-W. Sun, and P.-C. Chang, "A skeleton-based pairwise curve matching scheme for people tracking in a multi-camera environment," in *2013 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference*. IEEE, 2013, pp. 1–5.

[70] L. Herda, P. Fua, R. Plänkers, R. Boulic, and D. Thalmann, "Using skeleton-based tracking to increase the reliability of optical motion capture," *Human movement science*, vol. 20, no. 3, pp. 313–341, 2001.

[71] M. Shaheen, J. Gall, R. Strzodka, L. Van Gool, and H.-P. Seidel, "A comparison of 3D model-based tracking approaches for human motion capture in uncontrolled environments," in *2009 Workshop on Applications of Computer Vision (WACV)*. IEEE, 2009, pp. 1–8.

[72] J. Gall, C. Stoll, E. De Aguiar, C. Theobalt, B. Rosenhahn, and H.-P. Seidel, "Motion capture using joint skeleton tracking and surface estimation," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2009, pp. 1746–1753.

[73] S. Kim, K. Yun, J. Park, and J. Y. Choi, "Skeleton-based action recognition of people handling objects," in *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2019, pp. 61–70.

[74] Y. Hou, Z. Li, P. Wang, and W. Li, "Skeleton optical spectra-based action recognition using convolutional neural networks," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 3, pp. 807–811, 2016.

[75] C. Li, S. Sun, X. Min, W. Lin, B. Nie, and X. Zhang, "End-to-end learning of deep convolutional neural network for 3D human action recognition," in *2017 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*. IEEE, 2017, pp. 609–612.

[76] V.-N. Hoang, T.-L. Le, T.-H. Tran, V.-T. Nguyen *et al.*, "3D skeleton-based action recognition with convolutional neural networks," in *2019 International Conference on Multimedia Analysis and Pattern Recognition (MAPR)*. IEEE, 2019, pp. 1–6.

[77] Y. Du, Y. Fu, and L. Wang, "Skeleton based action recognition with convolutional neural network," in *2015 3rd IAPR Asian Conference on Pattern Recognition (ACPR)*. IEEE, 2015, pp. 579–583.

[78] H.-H. Pham, L. Khoudour, A. Crouzil, P. Zegers, and S. A. Velastin, "Learning to recognise 3D human action from a new skeleton-based representation using deep convolutional neural networks," *IET Computer Vision*, vol. 13, no. 3, pp. 319–328, 2018.

[79] B. Li, Y. Dai, X. Cheng, H. Chen, Y. Lin, and M. He, "Skeleton based action recognition using translation-scale invariant image mapping and multi-scale deep CNN," in *2017 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*. IEEE, 2017, pp. 601–604.

[80] C. Cao, C. Lan, Y. Zhang, W. Zeng, H. Lu, and Y. Zhang, "Skeleton-based action recognition with gated convolutional neural networks," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 11, pp. 3247–3257, 2018.

[81] D.-H. Nguyen, T.-N. Ly, T.-H. Truong, and D.-D. Nguyen, "Multi-column CNNs for skeleton based human gesture recognition," in *2017 9th International Conference on Knowledge and Systems Engineering (KSE)*. IEEE, 2017, pp. 179–184.

[82] M. Liu, C. Chen, and H. Liu, "3D action recognition using data visualization and convolutional neural networks," in *2017 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2017, pp. 925–930.

[83] S. Laraba, M. Brahimi, J. Tilmanne, and T. Dutoit, "3D skeleton-based action recognition by representing motion capture sequences as 2D-RGB images," *Computer Animation and Virtual Worlds*, vol. 28, no. 3-4, p. e1782, 2017.

[84] Z. Ding, P. Wang, P. O. Ogunbona, and W. Li, "Investigation of different skeleton features for CNN-based 3D action recognition," in *2017 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*. IEEE, 2017, pp. 617–622.

[85] Z. Rostami, M. Afrasiabi, and H. Khotanlou, "Skeleton-based action recognition using spatio-temporal features with convolutional neural networks," in *2017 IEEE 4th International Conference on Knowledge-Based Engineering and Innovation (KBEI)*. IEEE, 2017, pp. 0583–0587.

[86] J. Ren, N. Reyes, A. Barczak, C. Scogings, and M. Liu, "An investigation of skeleton-based optical flow-guided features for 3D action recognition using a multi-stream CNN model," in *2018 IEEE 3rd International Conference on Image, Vision and Computing (ICIVC)*. IEEE, 2018, pp. 199–203.

[87] Q. Ke, M. Bennamoun, S. An, F. Sohel, and F. Boussaid, "A new representation of skeleton sequences for 3D action recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3288–3297.

[88] D. C. Luvizon, D. Picard, and H. Tabia, "2D/3D pose estimation and action recognition using multitask deep learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5137–5146.

[89] F. Yang, Y. Wu, S. Sakti, and S. Nakamura, "Make skeleton-based action recognition model smaller, faster and better," in *Proceedings of the ACM Multimedia Asia on ZZZ*, 2019, pp. 1–6.

[90] M. Jiang, J. Kong, G. Bebis, and H. Huo, "Informative joints based human action recognition using skeleton contexts," *Signal Processing: Image Communication*, vol. 33, pp. 29–40, 2015.

[91] V. Bloom, D. Makris, and V. Argyriou, "G3D: A gaming action dataset and real time action recognition evaluation framework," in *2012 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*. IEEE, 2012, pp. 7–12.

[92] J. Wang, X. Nie, Y. Xia, Y. Wu, and S.-C. Zhu, "Cross-view action modeling, learning and recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 2649–2656.

[93] J. Tu, M. Liu, and H. Liu, "Skeleton-based human action recognition using spatial temporal 3D convolutional neural networks," in *2018 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2018, pp. 1–6.

[94] L. Shi, Y. Zhang, J. Cheng, and H. Lu, "Two-stream adaptive graph convolutional networks for skeleton-based action recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 026–12 035.

[95] W. Zheng, L. Li, Z. Zhang, Y. Huang, and L. Wang, "Relational network for skeleton-based action recognition," in *2019 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2019, pp. 826–831.

[96] Y. Li, R. Xia, X. Liu, and Q. Huang, "Learning shape-motion representations from geometric algebra spatio-temporal model for skeleton-based action recognition," in *2019 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2019, pp. 1066–1071.

[97] H. Wang and L. Wang, "Modeling temporal dynamics and spatial configurations of actions using two-stream recurrent neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 499–508.

[98] Y. Xu, J. Cheng, L. Wang, H. Xia, F. Liu, and D. Tao, "Ensemble one-dimensional convolution neural networks for skeleton-based action recognition," *IEEE Signal Processing Letters*, vol. 25, no. 7, pp. 1044–1048, 2018.

[99] C. Si, W. Chen, W. Wang, L. Wang, and T. Tan, "An attention enhanced graph convolutional LSTM network for skeleton-based action recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1227–1236.

[100] Y. Wang, Z. Xu, L. Li, and J. Yao, "Robust multi-feature learning for skeleton-based action recognition," *IEEE Access*, vol. 7, pp. 148 658–148 671, 2019.

[101] T.-N. Nguyen, D.-T. Pham, T.-L. Le, H. Vu, and T.-H. Tran, "Novel skeleton-based action recognition using covariance descriptors on most informative joints," in *2018 10th International Conference on Knowledge and Systems Engineering (KSE)*. IEEE, 2018, pp. 50–55.

[102] C. Li, Z. Cui, W. Zheng, C. Xu, R. Ji, and J. Yang, "Action-attending graphic neural network," *IEEE Transactions on Image Processing*, vol. 27, no. 7, pp. 3657–3670, 2018.

[103] L. Hu and J. Xu, "Body joints selection convolutional neural networks for skeletal action recognition," in *2017 IEEE 29th International Conference on Tools with Artificial Intelligence (ICTAI)*. IEEE, 2017, pp. 682–686.

[104] W. Ding, K. Liu, F. Cheng, H. Shi, and B. Zhang, "Skeleton-based human action recognition with profile hidden markov models," in *CCF Chinese Conference on Computer Vision*. Springer, 2015, pp. 12–21.

[105] G. Zhang and X. Zhang, "Multi-heads attention graph convolutional networks for skeleton-based action recognition," in *2019 IEEE Visual Communications and Image Processing (VCIP)*. IEEE, 2019, pp. 1–4.

[106] C. Youssef *et al.*, "Spatiotemporal representation of 3D skeleton joints-based action recognition using modified spherical harmonics," *Pattern Recognition Letters*, vol. 83, pp. 32–41, 2016.

[107] Z. Wang, C. Zhang, W. Luo, and W. Lin, "Key joints selection and spatiotemporal mining for skeleton-based action recognition," in *2018 25th IEEE International Conference on Image Processing (ICIP)*. IEEE, 2018, pp. 3458–3462.

[108] Y. Du, W. Wang, and L. Wang, "Hierarchical recurrent neural network for skeleton based action recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1110–1118.

[109] Y. Du, Y. Fu, and L. Wang, "Representation learning of temporal dynamics for skeleton-based action recognition," *IEEE Transactions on Image Processing*, vol. 25, no. 7, pp. 3010–3022, 2016.

[110] R. Poppe, "A survey on vision-based human action recognition," *Image and vision computing*, vol. 28, no. 6, pp. 976–990, 2010.

[111] D. Weinland, R. Ronfard, and E. Boyer, "A survey of vision-based methods for action representation, segmentation and recognition," *Computer vision and image understanding*, vol. 115, no. 2, pp. 224–241, 2011.

[112] M. B. Holte, C. Tran, M. M. Trivedi, and T. B. Moeslund, "Human pose estimation and activity recognition from multi-view videos: Comparative explorations of recent developments," *IEEE Journal of selected topics in signal processing*, vol. 6, no. 5, pp. 538–552, 2012.

[113] J. M. Chaquet, E. J. Carmona, and A. Fernández-Caballero, "A survey of video datasets for human action and activity recognition," *Computer Vision and Image Understanding*, vol. 117, no. 6, pp. 633–659, 2013.

[114] G. Guo and A. Lai, "A survey on still image based human action recognition," *Pattern Recognition*, vol. 47, no. 10, pp. 3343–3361, 2014.

[115] F. Zhu, L. Shao, J. Xie, and Y. Fang, "From handcrafted to learned representations for human action recognition: A survey," *Image and Vision Computing*, vol. 55, pp. 42–52, 2016.

[116] S. Herath, M. Harandi, and F. Porikli, "Going deeper into action recognition: A survey," *Image and vision computing*, vol. 60, pp. 4–21, 2017.

[117] A. G. D'Sa and B. Prasad, "A survey on vision based activity recognition, its applications and challenges," in *2019 Second International Conference on Advanced Computational and Communication Paradigms (ICACCP)*. IEEE, 2019, pp. 1–8.

[118] C. Jobanputra, J. Bavishi, and N. Doshi, "Human activity recognition: A survey," *Procedia Computer Science*, vol. 155, pp. 698–703, 2019.

[119] Y. Chen, Y. Tian, and M. He, "Monocular human pose estimation: A survey of deep learning-based methods," *Computer Vision and Image Understanding*, p. 102897, 2020.

[120] R. Li, H. Fu, W.-L. Lo, Z. Chi, Z. Song, and D. Wen, "Skeleton-based action recognition with key-segment descriptor and temporal step matrix model," *IEEE Access*, vol. 7, pp. 169 782–169 795, 2019.

[121] R. Baptista, E. Ghorbel, A. E. R. Shabayek, D. Aouada, and B. Ottersten, "Key-skeleton based feedback tool for assisting physical activity," in *2018 Zooming Innovation in Consumer Technologies Conference (ZINC)*. IEEE, 2018, pp. 175–176.

[122] F. Ahmed, P. Polash Paul, and M. L. Gavrilova, "Kinect-based gait recognition using sequences of the most relevant joint relative angles," 2015.

[123] A. Taha, H. H. Zayed, M. Khalifa, and E.-S. M. El-Horbaty, "Skeleton-based human activity recognition for video surveillance," *International Journal of Scientific & Engineering Research*, vol. 6, no. 1, pp. 993–1004, 2015.

[124] E. Cippitelli, E. Gambi, S. Spinsante, and F. Flórez-Revuelta, "Evaluation of a skeleton-based method for human activity recognition on a large-scale RGB-D dataset," 2016.

[125] Y. Hbali, S. Hbali, L. Ballihi, and M. Sadgal, "Skeleton-based human activity recognition for elderly monitoring systems," *IET Computer Vision*, vol. 12, no. 1, pp. 16–26, 2017.

[126] W. Ding, K. Liu, F. Cheng, and J. Zhang, "STFC: Spatio-temporal feature chain for skeleton-based human action recognition," *Journal of Visual Communication and Image Representation*, vol. 26, pp. 329–337, 2015.

[127] J. Weng, C. Weng, and J. Yuan, "Spatio-temporal naive-bayes nearest-neighbor (ST-NBNN) for skeleton-based action recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4171–4180.

[128] Z. Huang, C. Wan, T. Probst, and L. Van Gool, "Deep learning on lie groups for skeleton-based action recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 6099–6108.

[129] G. Hu, B. Cui, and S. Yu, "Joint learning in the spatio-temporal and frequency domains for skeleton-based action recognition," *IEEE Transactions on Multimedia*, 2019.

[130] P. Nikolov, O. Boumbarov, A. Manolova, K. Tonchev, and V. Poulkov, "Skeleton-based human activity recognition by spatio-temporal representation and convolutional neural networks with application to cyber physical systems with human in the loop," in *2018 41st International Conference on Telecommunications and Signal Processing (TSP)*. IEEE, 2018, pp. 1–5.

[131] B. Hosseini, R. Montagne, and B. Hammer, "Deep-aligned convolutional neural network for skeleton-based action recognition and segmentation," *arXiv preprint arXiv:1911.04969*, 2019.

[132] K. Zhu, R. Wang, Q. Zhao, J. Cheng, and D. Tao, "A cuboid CNN model with an attention mechanism for skeleton-based action recognition," *IEEE Transactions on Multimedia*, 2019.

[133] Q. Ke, M. Bennamoun, S. An, F. Sohel, and F. Boussaid, "Learning clip representations for skeleton-based 3D action recognition," *IEEE Transactions on Image Processing*, vol. 27, no. 6, pp. 2842–2855, 2018.

[134] H. Wang and L. Wang, "Beyond joints: Learning representations from primitive geometries for skeleton-based action recognition and detection," *IEEE Transactions on Image Processing*, vol. 27, no. 9, pp. 4382–4394, 2018.

[135] S. Wei, Y. Song, and Y. Zhang, "Human skeleton tree recurrent neural network with joint relative motion feature for skeleton based action recognition," in *2017 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2017, pp. 91–95.

[136] J. Ren, R. Napoleon, B. Andre, S. Chris, M. Liu, and J. Ma, "Robust skeleton-based action recognition through hierarchical aggregation of local and global spatio-temporal features," in *2018 15th International Conference on Control, Automation, Robotics and Vision (ICARCV)*. IEEE, 2018, pp. 901–906.

[137] J. Liu, A. Shahroudy, D. Xu, A. C. Kot, and G. Wang, "Skeleton-based action recognition using spatio-temporal LSTM network with trust gates," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 12, pp. 3007–3021, 2017.

[138] S. Jun and Y. Choe, "Deep batch-normalized LSTM networks with auxiliary classifier for skeleton based action recognition," in *2018 IEEE International Conference on Image Processing, Applications and Systems (IPAS)*. IEEE, 2018, pp. 279–284.

[139] J. Liu, G. Wang, P. Hu, L.-Y. Duan, and A. C. Kot, "Global context-aware attention LSTM networks for 3D action recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1647–1656.

[140] S. Zhang, Y. Yang, J. Xiao, X. Liu, Y. Yang, D. Xie, and Y. Zhuang, "Fusing geometric features for skeleton-based action recognition using multilayer LSTM networks," *IEEE Transactions on Multimedia*, vol. 20, no. 9, pp. 2330–2343, 2018.

[141] R. Cui, A. Zhu, S. Zhang, and G. Hua, "Multi-source learning for skeleton-based action recognition using deep LSTM networks," in *2018 24th International Conference on Pattern Recognition (ICPR)*. IEEE, 2018, pp. 547–552.

[142] F. Meng, H. Liu, Y. Liang, J. Tu, and M. Liu, "Sample fusion network: an end-to-end data augmentation network for skeleton-based human action recognition," *IEEE Transactions on Image Processing*, vol. 28, no. 11, pp. 5281–5295, 2019.

[143] J. Tu, H. Liu, F. Meng, M. Liu, and R. Ding, "Spatial-temporal data augmentation based on LSTM autoencoder network for skeleton-based human action recognition," in *2018 25th IEEE International Conference on Image Processing (ICIP)*. IEEE, 2018, pp. 3478–3482.

[144] H. Liu, J. Tu, M. Liu, and R. Ding, "Learning explicit shape and motion evolution maps for skeleton-based human action recognition,"

in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 1333–1337.

[145] I. Lee, D. Kim, S. Kang, and S. Lee, "Ensemble deep learning for skeleton-based action recognition using temporal sliding LSTM networks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1012–1020.

[146] S. Yan, Y. Xiong, and D. Lin, "Spatial temporal graph convolutional networks for skeleton-based action recognition," in *Thirty-second AAAI conference on artificial intelligence*, 2018.

[147] F. Ye, H. Tang, X. Wang, and X. Liang, "Joints relation inference network for skeleton-based action recognition," in *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 16–20.

[148] H. Yang, Y. Gu, J. Zhu, K. Hu, and X. Zhang, "PGCN-TCA: Pseudo graph convolutional network with temporal and channel-wise attention for skeleton-based action recognition," *IEEE Access*, vol. 8, pp. 10 040–10 047, 2020.

[149] R. Liu, C. Xu, T. Zhang, W. Zhao, Z. Cui, and J. Yang, "Si-GCN: Structure-induced graph convolution network for skeleton-based action recognition," in *2019 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2019, pp. 1–8.

[150] X. Gao, W. Hu, J. Tang, P. Pan, J. Liu, and Z. Guo, "Generalized graph convolutional networks for skeleton-based action recognition," *arXiv preprint arXiv:1811.12013*, 2018.

[151] X. Shi, H. Li, F. Liu, D. Zhang, J. Bi, and Z. Li, "Graph convolutional networks with objects for skeleton-based action recognition," in *2019 IEEE International Conferences on Ubiquitous Computing & Communications (IUCC) and Data Science and Computational Intelligence (DSCI) and Smart Computing, Networking and Services (SmartCNS)*. IEEE, 2019, pp. 280–285.

[152] L. Shi, Y. Zhang, J. Cheng, and H. Lu, "Skeleton-based action recognition with directed graph neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7912–7921.

[153] S. Ubalde, F. Gómez-Fernández, N. A. Goussies, and M. Mejail, "Skeleton-based action recognition using citation-kNN on bags of time-stamped pose descriptors," in *2016 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2016, pp. 3051–3055.

[154] J. Yang, C. Zhu, and J. Yuan, "Spatio-temporal multi-scale soft quantization learning for skeleton-based human action recognition," in *2019 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2019, pp. 1078–1083.

[155] E. Yu and J. K. Aggarwal, "Human action recognition with extremities as semantic posture representation," in *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*. IEEE, 2009, pp. 1–8.

[156] N. Zheng, J. Wen, R. Liu, L. Long, J. Dai, and Z. Gong, "Unsupervised representation learning with long-term dynamics for skeleton based action recognition," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

[157] R. Zhao, W. Xu, H. Su, and Q. Ji, "Bayesian hierarchical dynamic model for human action recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7733–7742.

[158] K. Chen and K. Forbus, "Action recognition from skeleton data via analogical generalization over qualitative representations," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

[159] W. Ding, K. Liu, X. Biao, and F. Cheng, "Skeleton-based human action recognition via screw matrices," *Chinese Journal of Electronics*, vol. 26, no. 4, pp. 790–796, 2017.

[160] S. M. Yoon and A. Kuijper, "Human action recognition based on skeleton splitting," *Expert systems with Applications*, vol. 40, no. 17, pp. 6848–6855, 2013.

[161] R. Chaudhry, F. Ofli, G. Kurillo, R. Bajcsy, and R. Vidal, "Bio-inspired dynamic 3D discriminative skeletal features for human action recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2013, pp. 471–478.

[162] Q. Ke, M. Bennamoun, H. Rahmani, S. An, F. Sohel, and F. Boussaid, "Learning latent global network for skeleton-based action prediction," *IEEE Transactions on Image Processing*, vol. 29, pp. 959–970, 2019.

[163] E. A. Suma, B. Lange, A. S. Rizzo, D. M. Krum, and M. Bolas, "Faast: The flexible action and articulated skeleton toolkit," in *2011 IEEE Virtual Reality Conference*. IEEE, 2011, pp. 247–248.

[164] E. De Aguiar, C. Theobalt, S. Thrun, and H.-P. Seidel, "Automatic conversion of mesh animations into skeleton-based animations," in *Computer Graphics Forum*, vol. 27, no. 2. Wiley Online Library, 2008, pp. 389–397.

[165] X. Chen and J. Feng, "Adaptive skeleton-driven cages for mesh sequences," *Computer Animation and Virtual Worlds*, vol. 25, no. 3-4, pp. 445–453, 2014.

[166] M. Zhao, T. Li, M. Abu Alsheikh, Y. Tian, H. Zhao, A. Torralba, and D. Katabi, "Through-wall human pose estimation using radio signals," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7356–7365.

[167] A. Shahroudy, J. Liu, T.-T. Ng, and G. Wang, "NTU RGB+D: A large scale dataset for 3d human activity analysis," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1010–1019.

[168] L. Xia, I. Gori, J. K. Aggarwal, and M. S. Ryoo, "Robot-centric activity recognition from first-person RGB-Dvideos," in *2015 IEEE Winter Conference on Applications of Computer Vision*. IEEE, 2015, pp. 357–364.

[169] L. Seidenari, V. Varano, S. Berretti, A. Bimbo, and P. Pala, "Recognizing actions from depth cameras as weakly aligned multi-part bag-of-poses," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2013, pp. 479–485.

[170] B. Ghanem, J. C. Niebles, C. Snoek, F. C. Heilbron, H. Alwassel, V. Escorcia, R. Krishna, S. Buch, and C. D. Dao, "The activitynet large-scale activity recognition challenge 2018 summary," *arXiv preprint arXiv:1808.03766*, 2018.

[171] B. Ghanem and N. Ahuja, "Maximum margin distance learning for dynamic texture recognition," in *European Conference on Computer Vision*. Springer, 2010, pp. 223–236.

[172] J.-F. Hu, W.-S. Zheng, J. Lai, and J. Zhang, "Jointly learning heterogeneous features for RGB-D activity recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 5344–5352.

[173] C. Chen, R. Jafari, and N. Kehtarnavaz, "UTD-MHAD: A multimodal dataset for human action recognition utilizing a depth camera and a wearable inertial sensor," in *2015 IEEE International conference on image processing (ICIP)*. IEEE, 2015, pp. 168–172.

[174] M. Müller, T. Röder, M. Clausen, B. Eberhardt, B. Krüger, and A. Weber, "Documentation mocap database HDM05," 2007.