

TaiChi Action Capture and Performance Analysis with Multi-view RGB Cameras

Jianwei Li, Siyu Mo and Yanfei Shen

Beijing Sport University
{jianwei,msy,syf}@bsu.edu.cn,

Abstract

Recent advances in computer vision and deep learning have influenced the field of sports performance analysis for researchers to track and reconstruct freely moving humans without any marker attachment. However, there are few works for vision-based motion capture and intelligent analysis for professional TaiChi movement. In this paper, we propose a framework for TaiChi performance capture and analysis with multi-view geometry and artificial intelligence technology. The main innovative work is as follows: 1) A multi-camera system suitable for TaiChi motion capture is built and the multi-view TaiChi data is collected and processed; 2) A combination of traditional visual method and implicit neural radiance field is proposed to achieve sparse 3D skeleton fusion and dense 3D surface reconstruction. 3) The normalization modeling of movement sequences is carried out based on motion transfer, so as to realize TaiChi performance analysis for different groups. We have carried out evaluation experiments, and the experimental results have shown the efficiency of our method.

1 Introduction

Human motion capture (Mocap) is usually used to obtain 3D human movement information in proactive health care and intelligent sports. Sports performance analysis and evaluation can improve athletes' competitive ability or promote public scientific fitness. The widely used inertial and optical motion capture systems can track and record human movement well, but need to bind sensors or paste marks on human body, which may affect human movement. Moreover, most current optical and inertial motion capture systems are expensive, and how to stick the markers also requires certain professional knowledge. Visual motion capture methods use cameras to non-invasively capture human motion images and then obtain the motion data through human pose estimation (HPE) and 3D reconstruction. Vision-based human action analysis is an important research topic in computer vision, and in recent years it has been widely applied in intelligent sports. Accurate 3D human motion modeling is the prerequisite for reliable human motion analysis.

Human action recognition (HAR) and action quality assessment (AQA) are two tasks of performance analysis in intelligent sports. The former aims to identify the action classification, while the latter aims to automatically quantify the performance of the action or to score its performance. Traditional methods for action analysis are mainly based on artificial design features, and compare the action sequences by estimating the distance error or dynamic time warping. Deep-learning methods use the deep network to directly learn action features and have shown more powerful performance. Generally speaking, deep-learning methods consist of video-based methods and skeleton-based methods. Algorithms based on video generally extract features directly from images, such as C3D [Parmar and Morris, 2020], I3D [Carreira and Zisserman, 2017], and TSN [Xiang *et al.*, 2018] and Pseudo3D [Qiu *et al.*, 2017], and then extract time domain features by LSTM, pooling, and so on. The finally score prediction is performed by a fully connected neural network. Skeleton-based methods first detect the human skeleton in the images or video, and then model the correlation information between human joints, so as to realize human motion modeling and motion quality evaluation.

At present vision-based performance analysis is mature in motion recognition and has made remarkable progress, but the performance of action quality assessment in sports motion scoring and intelligent sports training is still lower than the current application needs. Many studies on human movement evaluation with computer vision have been proposed, however most of them only select a few relatively simple fitness or clinical rehabilitation movements to identify and evaluate. As the current mainstream method, deep learning needs large-scale human motion dataset to train a better model, which limits their effects on sports action analysis. Although some human professional sports datasets have been presented in recent years, however most of them are RGB images or videos collected from the Internet, such as AQA-7 [Parmar and Morris, 2019a] and Yoga-82 [Verma *et al.*, 2020]. The performance of sports scoring or quality evaluation is still below the current application requirements.

According to above analyses, existing studies focus more on the recognitions of regular actions or assessments of competitive sports, and lack of 3D action dataset. As a complement to above work, we focus on how to capture and intelligently assess the quality of TaiChi actions. In summary the

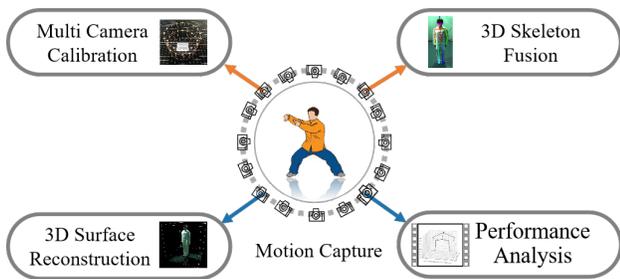


Figure 1: The system of TaiChi performance capture and analysis.

main **contributions** of this paper are the followings:

1. A professional TaiChi dataset consists of 23,232 action samples captured through multi-view cameras;
2. An effective 3D human modelling framework with multi-camera calibration, 3D skeleton fusion and 3D surface reconstruction;
3. A normalized modeling method for skeleton sequences based on motion transfer to analyse TaiChi performance from different groups.

2 Related Work

Vision-based motion capture. Human Mocap system based on vision technology can obtain 3D movement information non-invasively, and is gradually applied in the field of sports performance analysis. The type of human motion description includes human skeleton model (e.g., MPII [Andriluka *et al.*, 2014]), human parametric model (e.g., SMPL [Loper *et al.*, 2015]) and dense shape model (e.g., HumanNeRF [Weng *et al.*, 2022]). Among them, the skeleton model of human body describes the non-rigid motion of 3D surface with high degrees of freedom as surface motion driven by the dynamic chain. As a structured representation of human pose, skeleton model can conveniently and effectively represent quantitative information of human movements, and is widely used in action analysis. According to the number of viewpoints, visual Mocap systems can be divided into single-view system and multi-view system. The single-view system generally uses a single camera to capture human motion from a fixed perspective, while the multi-view system obtains human motion images from multiple perspectives based on the multi-camera system. The main challenge of single-view method is the problems of occlusion and depth uncertainty. Compared with the single-view method, the multi-view method can provide multi-view information, which can alleviate the occlusion problems and better restore 3D human posture. Imocap [Dong *et al.*, 2020] can capture human motion from multiple Internet videos, and open a new direction for 3D HPE. DeepMultiCap [Zheng *et al.*, 2021] uses the pixel alignment implicit function based on the parameterized model to reconstruct the invisible region response to the severe occlusion problem in the close-range interaction scene, and captures geometric details of human surface over time based on the attention module. Multi-view video data can contain more temporal and spatial information, but human posture

may vary dramatically in continuous frames of video data, so how to integrate the data effectively remains to be solved. Currently, mainstream deep learning methods often require a large number of labeled data, which increases the difficulty of model training for sports action.

Vision-based human motion datasets. NTU RGB+D [Liu *et al.*, 2019] is so far the largest Kinect-based action dataset collected from 106 distinct subjects and contains more than 114 thousand video samples and 8 million frames. The dataset contains 120 different action classes including daily, mutual, and health-related activities. Human 3.6M [Ionescu *et al.*, 2014] is another large dataset with 3.6 million human poses and corresponding images. There are 11 subjects and 17 action scenes, and the data is made up of four digital cameras, one time sensor and ten motion cameras. UCF-sport [Soomro and Zamir, 2014] is the first sports action dataset, and contains close to 200 action video sequences collected from various sports which are typically featured on broadcast television channels such as BBC and ESPN. Since then a number of sports motion datasets [Li *et al.*, 2018a; Shao *et al.*, 2020; Verma *et al.*, 2020] used for action recognition have emerged. FineGym [Shao *et al.*, 2020] provides coarse-to-fine annotations both temporally and semantically for gymnastics videos. There are three levels of categorical labels, and the temporal dimension is also divided into two levels, i.e., actions and sub-actions. UMONS-TAICHI [Tits *et al.*, 2018] includes 2,200 sequences of 13 classes (relative to different Taijiquan techniques) performed by 12 participants of different levels of expertise. Fitness-AQA [Parmar *et al.*, 2022] is a new exercise dataset comprising of three exercises (*BackSquat*, *BarbellRow* and *OverheadPress*), has been annotated by expert trainers for multiple crucial and typically occurring exercise errors. At present, most sports action datasets used for performance analysis are competition data based on publicly available RGB images or videos, such as public FSD-10 [Liu *et al.*, 2020] and Finediving [Xu *et al.*, 2022], but few of them have multi-view 3D skeleton poses.

Vision-based action analysis. With the development of intelligent sports and computer vision, many action analysis methods for sports have been gradually proposed in recent years. Deep learning is currently the mainstream method for vision-based action analysis, where the most widely used models are RNNs, CNNs, GCNs and Transform-based. According to the type of input data, there are mainly image-based [Duan *et al.*, 2022; Parmar and Morris, 2019b] and skeleton-based [Yan *et al.*, 2018; Pan *et al.*, 2019] methods. ScoringNet [Li *et al.*, 2018b] and SwingNet [McNally *et al.*, 2019] are based on images, and support fine-grained action classification and action scoring. These methods focus on the visual activity information of the whole scene including the performer’s body and background, but may tend to ignore the motion relationship within the human skeleton joints. Skeleton-based methods generally begin by extracting human skeleton, and then conduct spatio-temporal modeling on the association information between skeleton joints. For example, the joint relational graph method proposed by Pan *et al.* [Pan *et al.*, 2019] models the conventional motion and motion difference between different parts of human body according to the joint common module and joint difference

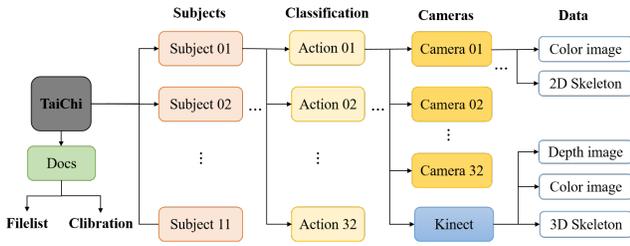


Figure 2: The organization of TaiChi data.

module respectively. HDVR [Hu and Ahuja, 2021] proposes a hierarchical dance video recognition framework by estimating 3D human pose from the corresponding 2D human pose sequences. For competitive gymnastics, SportsCap [Chen *et al.*, 2021] uses ST-GCN [Yan *et al.*, 2018] method to predict a fine-grained semantic action attributes, and adopts a semantic attribute mapping block to assemble various correlated action attributes into a high-level action label for the overall detailed understanding of the whole movement. In recent years, vision-based HAR methods are mature and have made remarkable progress, and AQA technologies have also been developed gradually. However, human body is often in self-occlusion with large folding or bending in sports, the performance of existing methods in rehabilitation training and sports scoring is still lower than current application requirements. The recognition accuracy of uncommon or highly similar human motion is still limited, and how to effectively model and analyze human motion with challenging situation, such as complex movements and view change, needs to be further studied.

3 Methods

For TaiChi performance capture and analysis, we design a non-invasive system with multi-view geometry and artificial intelligence technology.

3.1 Experimental Setup

Since there has a lot of body rotation in TaiChi movements, we set up a ring array multi-camera system to realize a better motion capture. The installation bracket of the system is a positive 16-sided shape with a diameter of 450 *cm* and a height of 250 *cm*, with a total of 16 columns. Each two RGB cameras (2448×2048p) from the FLIR company are installed on each column. There are 32 cameras in total, 16 of them are 100 *cm* from the ground and the rest of them are 200 *cm* from the ground. Cameras on the top are tilted down about 20 degrees, and cameras on the bottom are tilted down about 10 degrees. Each four cameras are connected to a server, and the 8 servers are networked through a 10 GB router. All cameras are synchronously controlled through a special trigger device. Data process and experiments are conducted on the PyTorch deep-learning framework on a standard desktop PC with 11 GB 1080Ti GPUs.

3.2 System Composition

As shown in Figure 1, the proposed system mainly contains five modules:

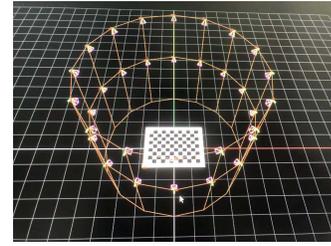


Figure 3: Visual simulation of our multi-camera calibration.

- Multi-camera calibration. Before motion capture, the system is calibrated to obtain the internal reference matrix of each camera and the pose relationship between multi cameras.
- Motion capture. The 32 high frame rate RGB cameras through 8 servers are synchronously controlled to capture and store the high-definition TaiChi motion data.
- 3D skeleton fusion. The 2D skeletons are obtained by the HPE method from each RGB image and then fused into 3D human skeleton by multi-visual geometry matching.
- 3D surface reconstruction. The 3D human surface model is reconstructed with multi-view images through camera poses estimation and neural radiance field rendering.
- Performance analysis. The skeleton sequences of different subjects are transferred to a standard model to eliminate individual appearance differences, and then TaiChi action quality is compared and evaluated by comparing the trajectory and angle changes of the re-targeted skeletons.

4 TaiChi Performance Capture

4.1 Multi-view Data Organization

Figure 2 shows the organization of multi-view TaiChi action data, which contains 23,232 action samples, including each sample’s RGB image, depth image, 2D skeleton and 3D skeleton data. Each TaiChi action sample is captured by 32 RGB cameras from 32 different views and a RGB-D camera (Kinect Azure) from the front view simultaneously. During TaiChi data acquisition, 11 subjects (3 female and 8 male) have performed the 24-form TaiChi actions same as TaiChi-24 [Li *et al.*, 2022]. Each action sample is manually segmented and labeled with category and action quality. The 2D skeletons are computed through Openpose [Cao *et al.*, 2017] algorithm, while the 3D skeletons are obtained from Kinect Azure SDK.

4.2 Multi-camera Calibration

To get the pose relationship of the 32 RGB cameras, we design a multi-camera calibration tool by the 2D planar checkerboard calibration method [Zhang, 2000]. Figure 3 shows a visual simulation of our multi-camera calibration process. The 2D checkerboard is located in the center of the multi-camera system, and its orientation and pitch angle are changed uniformly. More than 100 checkerboard images are selected in

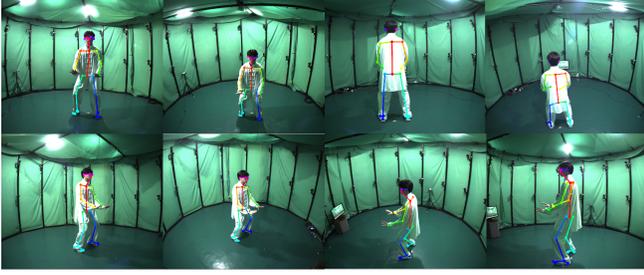


Figure 4: Example of 3D skeleton fusion results from 8 different views.

each calibration. The grids in the checkerboard are 10×15 , and the actual side length of each grid is 5 cm . In the process of data acquisition and processing, we have carried out 10 times of calibration. Based on the camera projection model, we construct a minimization objective function with re-projection error to solve the camera parameters:

$$\min \sum \| \mathbf{P}\mathbf{X}_i - \mathbf{x}_i \|^2, \quad (1)$$

where \mathbf{x} is the feature in RGB image, \mathbf{X} is the of in the checkerboard, and $\mathbf{P} = \mathbf{K}[\mathbf{R}|\mathbf{t}]$ is the camera projection matrix. \mathbf{K} is the i -th camera internal reference matrix, and $[\mathbf{R}_i|\mathbf{t}_i]$ is the i -th camera external reference matrix. In order to further improve the accuracy of calibration, bundle-adjustment (BA) [Triggs *et al.*, 2000] is used for optimization.

4.3 3D Skeleton Fusion

2D skeletons estimated from a single view RGB image often have the problem of occlusion, which will affect the accuracy of performance analysis. Therefore, we make 3D skeleton fusion by direct linear transformation (DLT) [Adbel-Aziz, 1971] algorithm. The main calculation process is shown in the following formula:

$$\mathbf{s}_i = \mathbf{K}_i[\mathbf{R}_i|\mathbf{t}_i]\mathbf{S}, i \in (1, m) \quad (2)$$

where $\mathbf{s}_i = \{s_1, s_2, \dots, s_N\}$ is the 2D skeleton with N joints in the i -th camera view, $\mathbf{S} = \{S_1, S_2, \dots, S_N\}$ is the corresponding 3D skeleton, and m is the number of fused camera views. Figure 4 shows an example of 3D skeleton fusion results (rendered on 2D images) for a subject from 8 different views.

4.4 3D Surface Reconstruction

Considering the excellent modeling and rendering capabilities of neural radiation fields (NeRFs), we realize 3D human surface reconstruction by joint using the traditional Colmap [Schonberger and Frahm, 2016] and deep-learning based Instant NeRF [Müller *et al.*, 2022]: 1) Firstly, we detect and extract SIFT features from each input image; 2) And then the positions of the multi cameras are estimated by feature matching; 3) Finally, data conversion is performed for NeRF rendering. Given a 3D point and a viewing direction $d \in \mathbb{R}^3$, NeRF estimates RGB color values and density (c, σ) that are

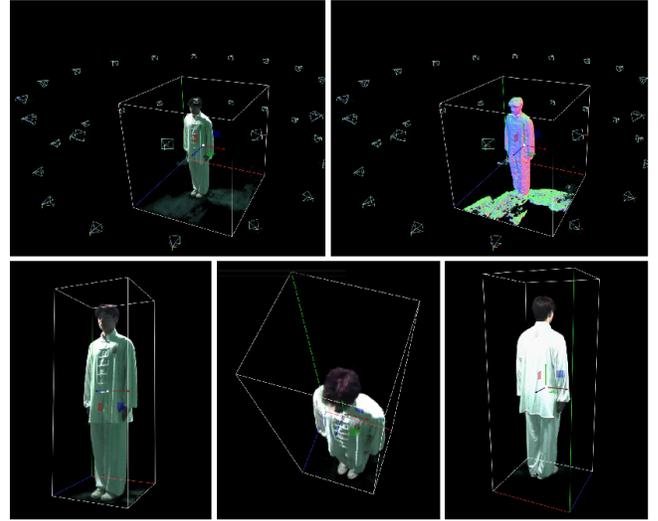


Figure 5: 3D human surface reconstruction with NeRFs.

then accumulated via quadrature to calculate the expected color of each camera ray:

$$C(r) = \int_{t_n}^{t_f} \exp\left(-\int_{t_n}^t \sigma(s)ds\right)\sigma(t)c(t,d)dt, \quad (3)$$

where t_f and t_n define the near and far bounds, and the camera ray is indicated as $r(t) = o + td$.

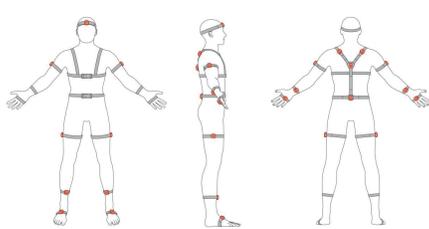
In order to speed up the signed distance function (SDF) training, 3D training positions are uniformly sampled and computed to the triangle mesh. Figure 5 shows the 3D surface reconstruction results of a subject with Instant NeRF. The top are the overall relationship between the camera poses and human models, while the bottom are the reconstruction details from the front, top and back views respectively.

5 TaiChi Performance Analysis

5.1 Data Precision Analysis

Accurate recovery of 3D human movement information is the premise of reliable analysis of human movement. In order to verify the accuracy of our multi-camera system with IMU-based motion capture system, we conduct experiments using these two systems in time and space synchronization. A subject has wore the IMU equipments (Perception Neuron Studio, PNS) and performed three upper limb exercises (*lateral hand lift, left punch and right punch*) and three lower limb exercises (*left knee lift, right knee lift and lunging squat*) when the frame rates of the multi-camera system are 30 fps and 60 fps respectively. Figure 6 shows the wearing position of the PNS sensors and human skeleton nodes (Body-25) extracted by Openpose. It can be seen that there has a position deviation for the joints obtained by these two Mocap systems.

We have compared and analyzed the Mocap data by the coordinate information of *shoulder, elbow, hip and knee* joints. Figure 7 shows the MSE error between our multi-camera system and IMU-based Mocap system. The average angle error



(a) The wearing positions of the PNS sensor



(b) The joints extracted with Openpose algorithm

Figure 6: Comparison of the wearing positions of IMU and human skeleton joints extracted from our visual Mocap system.

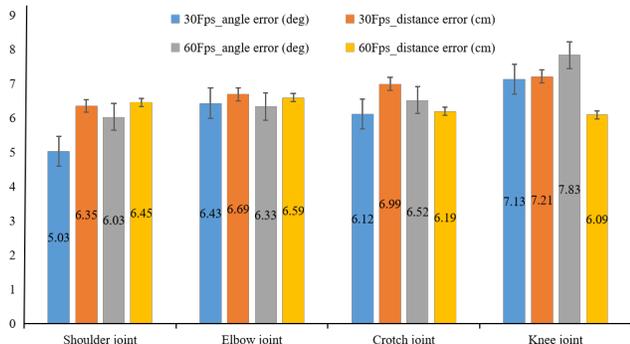


Figure 7: Comparison of multi-camera and IMU-based Mocap system by angle error (*deg*) and distance error (*cm*) of key joints.

and distance error are 6.18 *deg* and 6.81 *cm* respectively in 30 fps. The average angle error and distance error are 6.68 *deg* and 6.41 *cm* respectively in 60 fps. The main reason for the error lies in the deviation between the paste positions of the PNS sensor and skeleton nodes obtained by visual estimation (as shown in Figure 6).

5.2 Performance Analysis

For TaiChi performance analysis, we combine 32 camera views from 16 orientations into 16 stereo pairs to obtain 3D skeletons. 3D skeleton of each camera pair also is computed with 3D reconstruction module of Openpose. The quantization accuracies of TaiChi action recognition and assessment have been discussed in [Li *et al.*, 2022]. Therefore, we mainly conduct performance analysis by comparing the movements between the coach and the students. To avoid the influence of individual differences, the action skeleton sequences from the coach and students are uniformly transferred to a standard virtual human model for comparison.

The flowchart of performance analysis with the motion transfer network is shown in Figure 8. The model decomposes the skeleton sequences and recombines the elements to generate a new skeleton sequence, which can be viewed at any desired view-angle. We implement the action transfer based on Transmomo [Yang *et al.*, 2020] without using any paired data for supervision. The transfer network is trained in an unsupervised manner by exploiting invariance properties of three orthogonal factors of variation including motion,

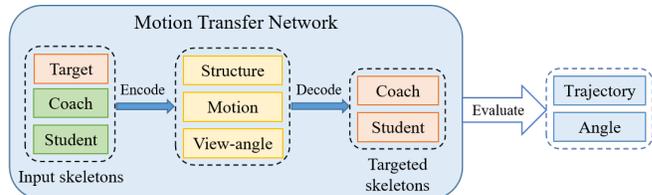


Figure 8: The flowchart of performance analysis with motion transfer network.

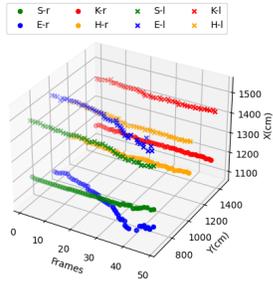
structure, and view-angle. The motion is invariant despite structural and view-angle perturbations, the structure is consistent through time and invariant despite view-angle perturbations, and the view-angle is consistent through time and invariant despite structural perturbations.

For an input sequence $\mathbf{s}^T \in \mathbb{R}^{T \times 2N}$ where T is the length of the skeleton sequence and N is the number of body joints. The motion encoder uses several layers of one dimensional temporal convolution to extract the motion information. The structure encoder has a similar network structure with the difference that the final structure code is obtained after a temporal max pooling. The view code is obtained the same way we obtained the structure code. The decoder takes the motion, body and view codes as input and reconstructs a 3D joint sequence $\mathbf{S}^T \in \mathbb{R}^{T \times 3N}$ through convolution layers, in symmetry with the encoders. The total loss function is derived based on invariance and weighted by the following loss terms:

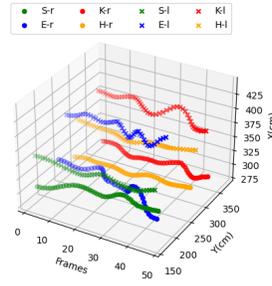
$$\mathbf{L}_{total} = \lambda_{rec} \mathbf{L}_{rec} + \lambda_{crs} \mathbf{L}_{crs} + \lambda_{trip} \mathbf{L}_{trip} + \lambda_{inv} \mathbf{L}_{inv} + \lambda_{adv} \mathbf{L}_{adv}, \quad (4)$$

where \mathbf{L}_{rec} is the reconstruction loss to minimize the difference between real data and 3D reconstructions projected back to 2D, \mathbf{L}_{crs} is the cross reconstruction loss for two sequences, \mathbf{L}_{trip} is the triplet loss to map views, \mathbf{L}_{inv} is the structural invariance loss to ensure that the view code is invariant to structural change estimations from the same sequence to a small neighborhood while alienating estimations from rotated sequences, and \mathbf{L}_{adv} is used to measure the domain discrepancy between the projected 2D sequences and real 2D sequences.

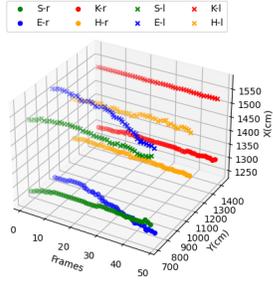
We evaluate the action quality by comparing the trajectory and angle changes of the key joints, such as *shoulder*, *elbow*, *hip* and *knee*. We select the first 15 key points defined in



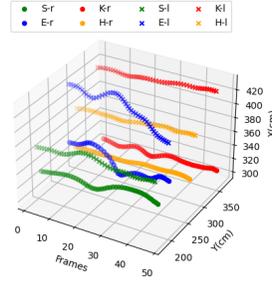
(a) Trajectories before motion transfer from student 1.



(b) Trajectories after motion transfer from student 1.

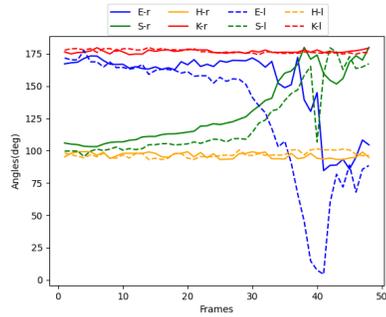


(c) Trajectories before motion transfer from student 2.

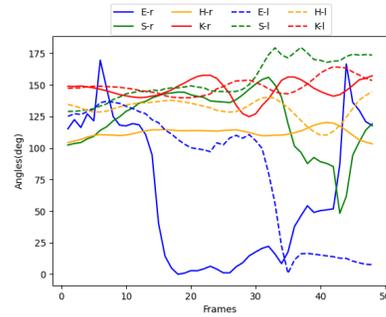


(d) Trajectories after motion transfer from student 2.

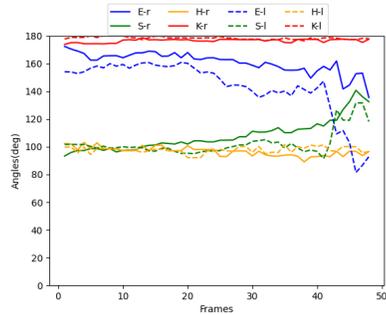
Figure 9: The changes of the key joints' trajectories before and after motion transfer.



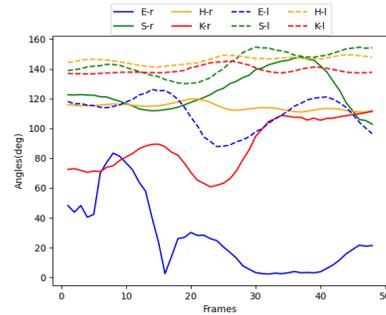
(a) Angles before motion transfer from student 1.



(b) Angles after motion transfer from student 1.



(c) Angles before motion transfer from student 2.



(d) Angles after motion transfer from student 2.

Figure 10: The changes of the key joints' angles before and after motion transfer.



Figure 11: Examples of the re-targeted key-frame skeletons (top: the coach; middle: student 1; bottom: student 2).

Body-25 for motion transfer. For the j -th joint in the skeleton model, the angle sequence (θ_j^T) for an action sample is calculated as follows:

$$\theta_j^T = \arccos\left(\frac{(\mathbf{S}_{j-1} - \mathbf{S}_j) \cdot (\mathbf{S}_{j+1} - \mathbf{S}_j)}{\|\mathbf{S}_{j-1} - \mathbf{S}_j\| \|\mathbf{S}_{j+1} - \mathbf{S}_j\|}\right), j \in (0, 15) \quad (5)$$

In our experiments, we take the RGB image (1920×1080p) sequence of the coach collected by Kinect camera as the target, and two skeleton sequences of the students captured by our multi-camera system as input data. Figure 9 and Figure 10 show the changes of the trajectories and angles of the *shoulder* (S-l and S-r), *elbow* (E-l and E-r), *hip* (H-l and H-r) and *knee* (K-l and K-r) joints respectively. The length of the action sample is 48 frames. By analyzing the changes of trajectories and angles after motion transfer, we can find the major difference joints and emergence moments. In the case of these two students, the differences of them are mainly concentrated in the *shoulder* and *knee* joints during the movements located at the end of the action. Figure 11 shows the visual analysis results of re-targeted skeletons from two students compared with the coach. It can be seen that the TaiChi motion of student1 is more standard. According to the scores of professional coach, student 1 (100) also scored higher than student 2 (86).

6 Conclusions

In this paper, we build a multi-camera system and realize TaiChi performance capture and analysis with multi-view geometry and artificial intelligence technology. Traditional visual method and deep-learning base NeRFs are used in combination to achieve sparse 3D skeleton and dense 3D surface reconstruction. we collect and organize the TaiChi data with multi RGB cameras, and process them to experimental analysis. To realize TaiChi performance analysis for different groups, skeleton sequences are normalized modeled with motion transfer and then further assessed with the changes of the joints' trajectories and angles. We also carry out evaluation experiments and the experimental results have shown the efficiency of our system. In the future, we will research more accurate and robust analysis methods with data from less camera views.

Acknowledgments

The work is assisted by Haiqing Hu, Jinyang Li, Xinyu Wang and Tianhan Zhang from Beijing Sports University, who help us collect and process the data. This work is partially supported by the National Key R&D Program of China (No. 2022YFC3600300, No. 2022YFC3600305), and the Fundamental Research Funds for Central Universities No.2022QN018.

References

- [Adbel-Aziz, 1971] YI Adbel-Aziz. Direct linear transformation from comparator coordinates into object space in close-range photogrammetry. In *ASP Symp. Proc. on Close-Range Photogrammetry, American Society of Photogrammetry, Falls Church, 1971*, pages 1–18, 1971.
- [Andriluka *et al.*, 2014] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014.
- [Cao *et al.*, 2017] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7291–7299, 2017.
- [Carreira and Zisserman, 2017] J. Carreira and A. Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [Chen *et al.*, 2021] Xin Chen, Anqi Pang, Wei Yang, Yuexin Ma, Lan Xu, and Jingyi Yu. Sportscap: Monocular 3d human motion capture and fine-grained understanding in challenging sports videos. *International Journal of Computer Vision*, 129(10):2846–2864, 2021.
- [Dong *et al.*, 2020] Junting Dong, Qing Shuai, Yuanqing Zhang, Xian Liu, Xiaowei Zhou, and Hujun Bao. Motion capture from internet videos. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 210–227. Springer, 2020.
- [Duan *et al.*, 2022] Haodong Duan, Yue Zhao, Kai Chen, Dahua Lin, and Bo Dai. Revisiting skeleton-based action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2969–2978, 2022.
- [Hu and Ahuja, 2021] Xiaodan Hu and Narendra Ahuja. Un-supervised 3d pose estimation for hierarchical dance video recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11015–11024, 2021.
- [Ionescu *et al.*, 2014] C Ionescu, D Papava, V Olaru, and C Sminchisescu. 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):1325–1339, 2014.

- [Li *et al.*, 2018a] Yingwei Li, Yi Li, and Nuno Vasconcelos. Resound: Towards action recognition without representation bias. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 513–528, 2018.
- [Li *et al.*, 2018b] Yongjun Li, Xiujuan Chai, and Xilin Chen. Scoringnet: Learning key fragment for action quality assessment with ranking loss in skilled sports. In *Asian Conference on Computer Vision*, pages 149–164. Springer, 2018.
- [Li *et al.*, 2022] Jianwei Li, Haiqing Hu, Qingjun Xing, Xinyu Wang, Jinyang Li, and Yanfei Shen. Tai chi action quality assessment and visual analysis with a consumer rgb-d camera. In *2022 IEEE 24th International Workshop on Multimedia Signal Processing (MMSP)*, pages 1–6. IEEE, 2022.
- [Liu *et al.*, 2019] Jun Liu, Amir Shahroudy, Mauricio Perez, Gang Wang, Ling-Yu Duan, and Alex C Kot. Ntu rgb+ d 120: A large-scale benchmark for 3d human activity understanding. *IEEE transactions on pattern analysis and machine intelligence*, 42(10):2684–2701, 2019.
- [Liu *et al.*, 2020] Shenglan Liu, Xiang Liu, Gao Huang, Hong Qiao, Lianyu Hu, Dong Jiang, Aibin Zhang, Yang Liu, and Ge Guo. Fsd-10: A fine-grained classification dataset for figure skating. *Neurocomputing*, 413:360–367, 2020.
- [Loper *et al.*, 2015] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *ACM transactions on graphics (TOG)*, 34(6):1–16, 2015.
- [McNally *et al.*, 2019] William McNally, Kanav Vats, Tyler Pinto, Chris Dulhanty, John McPhee, and Alexander Wong. GolfdB: A video database for golf swing sequencing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019.
- [Müller *et al.*, 2022] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022.
- [Pan *et al.*, 2019] Jia-Hui Pan, Jibin Gao, and Wei-Shi Zheng. Action assessment by joint relation graphs. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6331–6340, 2019.
- [Parmar and Morris, 2019a] Paritosh Parmar and Brendan Morris. Action quality assessment across multiple actions. In *2019 IEEE winter conference on applications of computer vision (WACV)*, pages 1468–1476. IEEE, 2019.
- [Parmar and Morris, 2019b] Paritosh Parmar and Brendan Tran Morris. What and how well you performed? a multitask learning approach to action quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 304–313, 2019.
- [Parmar and Morris, 2020] P. Parmar and B. T. Morris. What and how well you performed? a multitask learning approach to action quality assessment. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [Parmar *et al.*, 2022] Paritosh Parmar, Amol Gharat, and Helge Rhodin. Domain knowledge-informed self-supervised representations for workout form assessment. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVIII*, pages 105–123. Springer, 2022.
- [Qiu *et al.*, 2017] Zhaofan Qiu, Ting Yao, and Tao Mei. Learning spatio-temporal representation with pseudo-3d residual networks. In *proceedings of the IEEE International Conference on Computer Vision*, pages 5533–5541, 2017.
- [Schonberger and Frahm, 2016] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016.
- [Shao *et al.*, 2020] Dian Shao, Yue Zhao, Bo Dai, and Dahua Lin. Finegym: A hierarchical video dataset for fine-grained action understanding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2616–2625, 2020.
- [Soomro and Zamir, 2014] Khurram Soomro and Amir R Zamir. Action recognition in realistic sports videos. *Computer vision in sports*, pages 181–208, 2014.
- [Tits *et al.*, 2018] Mickaël Tits, Sohaib Laraba, Eric Caulier, Joëlle Tilmanne, and Thierry Dutoit. Umons-taichi: A multimodal motion capture dataset of expertise in taijiquan gestures. *Data in brief*, 19:1214–1221, 2018.
- [Triggs *et al.*, 2000] Bill Triggs, Philip F McLauchlan, Richard I Hartley, and Andrew W Fitzgibbon. Bundle adjustment—a modern synthesis. In *Vision Algorithms: Theory and Practice: International Workshop on Vision Algorithms Corfu, Greece, September 21–22, 1999 Proceedings*, pages 298–372. Springer, 2000.
- [Verma *et al.*, 2020] Manisha Verma, Sudhakar Kumawat, Yuta Nakashima, and Shanmuganathan Raman. Yoga-82: a new dataset for fine-grained classification of human poses. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 1038–1039, 2020.
- [Weng *et al.*, 2022] Chung-Yi Weng, Brian Curless, Pratul P Srinivasan, Jonathan T Barron, and Ira Kemelmacher-Shlizerman. Humannerf: Free-viewpoint rendering of moving people from monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16210–16220, 2022.
- [Xiang *et al.*, 2018] X. Xiang, Y. Tian, A. Reiter, G. D. Hager, and T. D. Tran. S3d: Stacking segmental p3d for action quality assessment. pages 928–932, 2018.
- [Xu *et al.*, 2022] Jinglin Xu, Yongming Rao, Xumin Yu, Guangyi Chen, Jie Zhou, and Jiwen Lu. Finediving: A fine-grained dataset for procedure-aware action quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2949–2958, 2022.

- [Yan *et al.*, 2018] Sijie Yan, Yuanjun Xiong, and Dahua Lin. Spatial temporal graph convolutional networks for skeleton-based action recognition. In *Thirty-second AAAI conference on artificial intelligence*, 2018.
- [Yang *et al.*, 2020] Zhuoqian Yang, Wentao Zhu, Wayne Wu, Chen Qian, Qiang Zhou, Bolei Zhou, and Chen Change Loy. Transmomo: Invariance-driven unsupervised video motion retargeting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5306–5315, 2020.
- [Zhang, 2000] Zhengyou Zhang. A flexible new technique for camera calibration. *IEEE Transactions on pattern analysis and machine intelligence*, 22(11):1330–1334, 2000.
- [Zheng *et al.*, 2021] Yang Zheng, Ruizhi Shao, Yuxiang Zhang, Tao Yu, Zerong Zheng, Qionghai Dai, and Yebin Liu. Deepmulticap: Performance capture of multiple characters using sparse multiview cameras. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6239–6249, 2021.