

4DFA: Four-Dimensional Full-Anatomy Reconstruction of Individualized Digital Human Models Based on Motion Videos

Rui Zhao^{1, a}, Jiachen Mi^{2,3, b}, Yanxin Jiang^{4, c}, Zhefu Chen^{*5, d}, Hongkai Wang^{*6,7, e}

¹ School of Biomedical Engineering, Dalian University of Technology, Dalian 116024, China;

² School of Biomedical Engineering, Dalian University of Technology, Dalian 116024, China;

³ INTESIM(Dalian)co, LTD, Dalian, 116023, China;

⁴ School of Biomedical Engineering, Dalian University of Technology, Dalian 116024, China;

⁵ School of Kinesiology and Health Promotion Dalian University of Technology, Dalian 116024, China;

⁶ School of Biomedical Engineering, Dalian University of Technology, Dalian 116024, China;

⁷ Liaoning Key Laboratory of Integrated Circuit and Biomedical Electronic System, Dalian University of Technology, Dalian 116024, China.

^a Zrui@mail.dlut.edu.cn, ^b 1196026755@qq.com, ^c JYX@mail.dlut.edu.cn,

^d chenzhefu@dlut.edu.cn, ^e wang.hongkai@dlut.edu.cn

Abstract. Reconstructing personalized models of the human body is a key step for digital twin modeling of sports motion. There have been many studies reconstructing the body surface models based on video sequences, but these algorithms lack the modeling of internal structures (bones, muscles, blood vessels, etc.). Some researchers map internal organ structures into the 3D point cloud scanning of the body surface, but such mapping does not guarantee realistic tissue deformation and is not developed for video processing. This paper focuses on the video-based reconstruction of a personalized full-anatomy digital model in sports motion. Our method creates a four-dimensional (4D) model of the moving process composed of the 3D models of each time moment. We first use a deep network to regress the body surface of each time point and then estimate the internal structures by registering a previously constructed deformable human anatomy atlas to the body surfaces. To mimic realistic internal structure deformation, we applied articulated linear transform to the bony structures to avoid bone shape distortion and applied smooth nonlinear transform to the soft tissue to mimic the motion deformation. We also developed an intersection removal method to prevent the possible intersection between bones and soft tissues during the deformation. Our method is tested with athletes' sports videos and automatically generates full-anatomy motion models which cannot be achieved by previous methods. Our method also surpasses the existing video-based body surface modeling method by providing additional internal structures. The 4D human motion model constructed by our method is useful in sports modeling, biomechanics simulation, wearable device design, etc.

Keywords: Pose deformation; Personalized human modeling; Human anatomy atlas; Internal organ adaptive modeling.

1. Introduction

In the field of kinesiology and health promotion, digital human models are widely used for biomechanical simulation of body motion, with the purpose of sports injury prediction, muscle load assessment, etc. Although simulation based on a generic human model has been studied for decades, personalized simulation using a digital twin model is just an emerging field. The simulation of individual motion not only requires accurate modeling of personal anatomy but also demands precise capturing of the body movement. Thanks to the fast development of deep learning techniques in recent years, individualized motion modeling based on personal sports videos becomes feasible.

So far, a series of deep learning algorithms have been proposed to address the above challenges. Most of these methods register deformable human surface models (such as SPACE[11] and SMPL[12]) to the sports photo or video sequence. Early methods in this field first infer bone joint positions from the 2D images and then use the joints to predict the 3D body surface model. Later, Kanazawa et al.[2] proposed an end-to-end method that realizes direct image-to-model mapping. There were also methods focusing on body part pose estimation such as the trunk and limbs[3], [4], hand posture[5], [6], and facial expression[7], [8]. Rong et al.[9] proposed a modular system that integrates the regression outputs of three parts to unify the output of whole-body pose estimation. Xiu et al.[10] introduced the normal graph and made geometric constraints on the body to improve the pose-level modeling accuracy.

Despite the fast progress in deep learning-based motion modeling, the current methods only generate body surfaces. However, biomechanical simulation requires not only the external surface but also the internal structures like bones and muscles. To tackle this problem, some studies register a human atlas with internal anatomy to the body surface scan data or volumetric tomographic images. For example, Keller et al.[13] inferred the skeleton from a 3D surface model in an arbitrary pose, Ali-Hamadi et al.[14] combined the fat distribution and other information inferred from MRI data, and Sugar et al.[15] calculated the deformation of anatomical structures based on the displacement transformation of the predefined virtual skeleton in different poses[16]–[18]. However, these methods require complex data acquisition (surface scanning or tomographic imaging) which is not suitable for fast motion capturing. Moreover, these approaches use nonlinear spatial transform to map the internal structures into the individual body, resulting in unrealistic tissue deformation for large posture changes.

In summary, state-of-the-art (SOTA) deep learning methods only model the body surface without internal structures, while the current internal anatomy modeling approaches are not developed for video data and are prone to unrealistic tissue deformation. Our study aims to address these problems with the following key features:

- i) Our method combines the deep learning-based body surface model with the internal structure mapping of a whole-body anatomy atlas. We for the first time realizes full anatomy modeling based on simple sports video.
- ii) We use a previously constructed deformable human atlas trained with thousands of computed tomographic images for internal structure modeling. The prior knowledge of internal anatomy from a large population dataset ensures reasonable estimation of unseen organs from the external surface.
- iii) To simulate the body motion, we apply articulated rigid transform and smooth nonrigid deformation to the bones and soft tissues, respectively. Unlike the existing methods which apply nonrigid deformation to all the structures, our method avoids anatomically implausible distortion of the bones and realizes anatomically realistic motion deformation of the whole body anatomy.
- iv) We also develop a method to remove the possible intersections between the bones and the soft tissue during large pose changes. The generation of an intersection-free model is crucial for subsequent finite element (FEM) simulation of the motion kinetics.

2. Materials and Methods

A. 3D digital human anatomical model

In this study, we used a deformable digital human atlas [19]constructed by our group based on 1063 whole-body CT images. Anatomical structures including the skin, bones, muscles, torso organs, and great vessels were segmented from the CT images and the inter-subject anatomical variations are learned via the statistical shape model (SSM) method. With *prior* anatomical knowledge learned from large population data, this atlas can model the shape correspondence

between the external skin and internal structures, facilitating the estimation of internal anatomy in the subsequent steps.

B. Motion armature

To endow our human atlas with posture-changing ability, this study defines an articulated motion armature with the armature joints located at the skeleton joint centers. For the balance between the modeling reality and the computation complexity, the motion armature is defined with 26 individual segments. The control range of each armature segment is defined via the closest point searching between the model vertices and the armature segments. Fig. 1 illustrates the deformable human atlas with internal organs and the defined articulated armature with the individual segment control ranges.

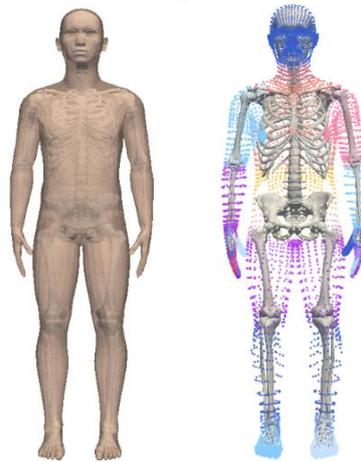


Fig. 1. Deformable human anatomy atlas with articulated motion armature. (a) The atlas in standing pose with transparent skin rendering to see the internal structures. (b) Control range of each armature segment, rendered with different pseudo colors for different segments.

C. Personalized pose deformation

Fig. 2 shows the workflow of video-based full-anatomy motion modeling. First, the video sequence is input to a deep neural network for body surface pose estimation of every single moment. Then, for every single moment, the articulated bone rotation of every armature segment is obtained by computing the rigid spatial transform (including translation and rotation) from the standing pose to the motion pose of this segment. Based on the obtained transform, articulated rigid and nonrigid deformations are applied to the bones and soft tissues respectively. Since the bones and the soft tissues are deformed using different types of spatial transforms, the possible intersection between them may occur, thus the intersection removal step is finally performed to correct possible intersections.

For video-based body surface estimation, we adopt the deep neural network model used for 3D human pose modeling [20]. This network uses the Resnet50 [21] structure as the image encoder to compute the spatial features of each time frame and then uses a two-layer gated recurrent unit with a hidden size of 1024 to compute the temporal features, followed by two fully-connected layers with 1024 neurons to regress the pose and body shape parameters of the SMPL model. The reconstructed SMPL model is used to get the 3D joint points and 2D projection coordinate points. The loss function used for training is,

$$L_g = L_{3D} + L_{2D} + L_{SMPL} + L_{adv} \#(1)$$

where L_{3D} , L_{2D} , L_{SMPL} , and L_{adv} are the losses of 3D joint points, 2D projection coordinate points, and the regressed SMPL's pose and shape parameters losses, respectively.

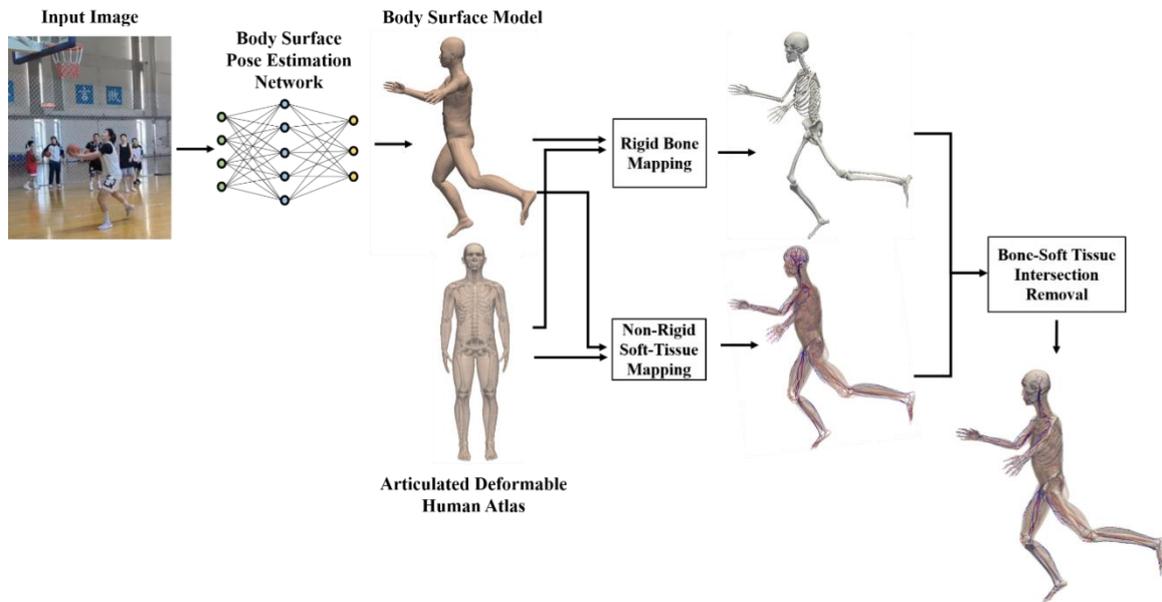


Fig.2. Workflow of video-based full-anatomy motion modeling.

D. Articulated rigid bone

The articulated rigid bone transform is applied in a bottom-to-top manner following the parent-child joint inheritance relationship in the moving armature. The child joints with the lowest generation are transformed first, and then its relative parent joints are transformed, with all the children joints of this parent joints transformed simultaneously. This process is recurrently applied to all the armature segments and the articulated posing of the entire body is realized.

E. Smooth nonrigid deformation of the soft tissues

To mimic the smooth nonrigid motion of the soft tissues, we use the Centers of Rotation (CoR) method [17] instead of the commonly used Linear Blending Skinning (LBS). The CoR method was originally proposed for solving the deformation artifacts of the LBS method for the body surface. In this study, we use CoR to produce smooth nonrigid deformation of both the skin and the internal soft tissues. Due to the page limitation of this paper, we refer the readers to [17] for detailed principles of the CoR method.

F. Intersection removing algorithm

In the process of digital human pose deformation, the intersection between the soft tissue and bones may be caused because the bones and soft tissues are deformed using different types of transforms. It is necessary to take intersection removal operations after the above deformation process to eliminate the possible surface crossing and generate a model available for FEM simulation.

Let the two surfaces that need to do intersection removal be S_1 and S_2 , and both must ensure that they are closed mesh surfaces in 3D space. The intersection-removing process we developed is as follows,

- i) Iterate through all vertices on S_1 and detect whether each vertex is inside S_2 . To do this, let any vertex $v \in S_1$, make n rays from v , and check the number of intersection points of each ray with S_2 . If the number of intersection points between all n rays and S_2 is odd, it means that v is inside S_2 , otherwise, it is outside. In this project, $n=10$ can satisfy all test cases.
- ii) Find the point to be adjusted. If the object to be detected is between soft tissue and bone or between bone and bone, the point to be adjusted is the point in S_1 that is inside S_2 ; if the object to be detected is between soft tissue and skin, the point to be adjusted is the point in S_1 that is outside S_2 .

iii) For each point v^* to be adjusted, through v make vertical lines for all triangular faces in S_2 and record the distance from v^* to the vertical foot if the foot falls inside the triangle. Take the shortest one among all distances as the shortest distance from v^* to S_2 , and then take the vertical foot o corresponding to the shortest distance as the target point of v^* , then the coordinates of v^* should be updated as,

$$v^* \leftarrow 1.5 \cdot (o - v^*) \quad (2)$$

- iv) Repeat the process described in (i) - (iii) above until no point to be adjusted is detected.
- v) Based on the result of (iv), do the faces cross detection? Iterate through all the triangular faces on S_1 and detect their intersection relationship with all the faces on S_2 . Specific method: set t_1 as any face of S_1 , and detect the cross relationship between its three edges and t_2 of any face of S_2 . As long as any one of the three edges intersects with t_2 , t_1 is considered to intersect with t_2 .
- vi) For the vertices of all faces on S_1 that intersect S_2 , adjust their positions as in (ii), (iii), and (iv).
- vii) Repeat the process of (v)-(vi) until no intersecting faces are detected.

3. Result and Analysis

Table I compares our method with several existing methods, including modeling internal anatomical structures, a 3D human reconstruction based on videos or medical images, and complete automation of the reconstruction algorithm. In the table Framkmocap[9] and vibe[20] are methods to reconstruct body surface models based on video sequences, where Framkmocap effectively integrates modules for hand pose and face expression estimation and outputs full body pose estimation results; Vibe proposes a grid structure based on time series, using an adversarial learning framework to generate realistic and reasonable human actions without real 3D labels. We used a deep learning method with high accuracy as the first step to reconstruct the human skin model. Anatomy transfer[14] proposed a semi-automatic method to estimate the internal anatomical structure by combining MRI images, while osso[13] combined the DXA detection results to generate a skin model that fits the parametric body model to achieve a 3D body skin to internal skeletal mapping.

Our proposed method is based on the digital atlas of personalized human whole-body anatomy, combined with the topological structure of human motion armature, to automatically rotate the joints and deform the posture of the anatomical structure of the whole body. At the same time, it solves the unrealistic distortion during large pose deformation and is more suitable for the reconstruction of 3D models with a large number of motion moments during the complete motion.

Table 1. Comparison between our method and existing methods

Method	Skin	Bone	Internal Anatomy	Input Data	Automation Level
Framkmocap[9]	√	×	×	Video	Automatic
vibe[20]	√	×	×	Video	Automatic
Anatomy transfer [14]	√	√	√	MRI	Semi-automatic
OSSO[13]	√	√	×	DXA	Semi-automatic
Ours	√	√	√	Video	Automatic

A. Realistic anatomical structure deformation

In the process of digital human overall posture deformation, we use the Center of Rotation algorithm to simulate the real flexible soft tissue deformation, which solves the volume loss at the joints. At the same time, due to the inability to strictly simulate the complex rotation and sliding processes of human joints, human posture deformation may cause an intersection between adjacent anatomical structures. Among the currently available studies, Keller et al. et al.[13] and Ali-Hamadi et al.[14] obtained visually plausible results for the postural transformation of human anatomical structures. However, the limitation is that the results still have some slight skin interpenetration or overlapping of internal tissues. For this reason, we developed a fully automated de-crossing algorithm to prevent possible crossings of bones and soft tissues during the deformation process.

Fig. 3 compares the results after skin optimization, using the skin at the arm joint as an example. The left figure shows the results of the commonly used linear blending skinning(LBS), and the right figure shows the results of the COR skinning method we applied. The comparison shows that COR improves the volume loss problem when skinning the LBS, which has obvious artifact depressions in both wrist and shoulder, and COR solves most of the artifacts introduced in the LBS and improves the mesh quality. Our proposed method is based on the digital atlas of personalized human whole-body anatomy, combined with the topological structure of human motion armature, to automatically rotate the joints and deform the posture of the anatomical structure of the whole body. At the same time, it solves the unrealistic distortion during large pose deformation and is more suitable for the reconstruction of 3D models with a large number of motion moments during the complete motion.

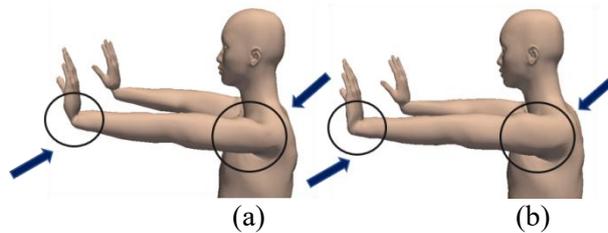


Fig. 3 The results of smooth deformation. (a) Linear blending skinning result with obvious artifact at the joints; (b) Our realistic deformation result.

Fig. 4 compares the results of our method with those of existing studies in which 3D point cloud scans based on the body surface achieve mapping of internal structural structures. Our method uses both rigid and flexible spatial deformation for bones and soft tissues respectively, which avoids unrealistic distortion deformation during the simulation of postural deformation. The results show that our method produces more reasonable internal anatomical structure postural deformation results so that no unreasonable shape distortion occurs in bones and blood vessels.

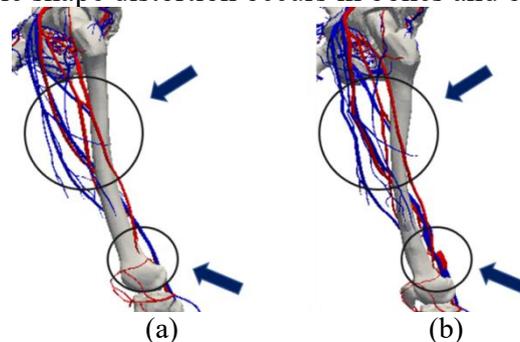


Fig. 4 The results of internal anatomical structure mapping.(a) Our method produces smooth soft tissue deformation and rigid bone transform (b) The existing nonlinear transformation method may produce distortions of soft tissue and bones

We developed a fully automated de-crossing algorithm to prevent possible intersections between bones and soft tissues during deformation. Fig. 5 shows the results of our fully automated

de-crossing algorithm for an intersection between bones and soft tissues during deformation. Here the intersection between blood vessels and bones near the knee joint is used as an example, while the whole de-crossing process takes about 8 minutes to run on a laptop with an i7 CPU.

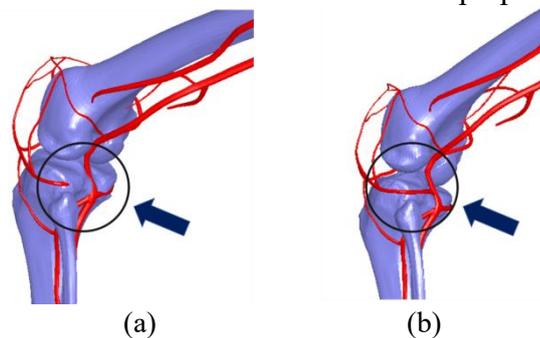


Fig. 5 The results of the intersection removal between the blood vessels and bones. (a) Before intersection removal; (b) After intersection removal.

B. Adaptive modeling results of whole-body anatomical structures

Our method is based on the digital atlas of personalized human body anatomy structure, designing a whole set of human motion digital modeling methods. We reconstruct the 3D human skin model from the video, combine human motion armature topology for joint rotation and posture deformation of whole-body anatomical structures, where the skeleton and soft tissue are rigid and flexible deformation respectively, and introduce the de-crossing mechanism. Finally, reasonable internal anatomical structure and posture deformation structure are generated.

Fig. 6 (a) shows a series of action images acquired from a video sequence. (b) The 3D human model reconstructed from the motion images (c) The human anatomy model corresponding to each motion instant. Arteries and veins, soft tissues, and bones are shown in different colors, respectively. (d) Local zoomed-in view of the model. The human anatomy after postural transformation is accurate and without unreasonable distortion.

The limitation of the current study is that it is not possible to quantitatively assess the accuracy of the postural deformation results of the whole-body anatomical structure, but the obtained postural deformation results are visually plausible. These deformation results have no irrational shape distortion among bones, muscles, and internal organs, and no geometric intersection between adjacent anatomical structures, which can provide digital models for subsequent human finite element simulation or biomechanical simulation for simulation testing, thus effectively meeting the needs of different applications.

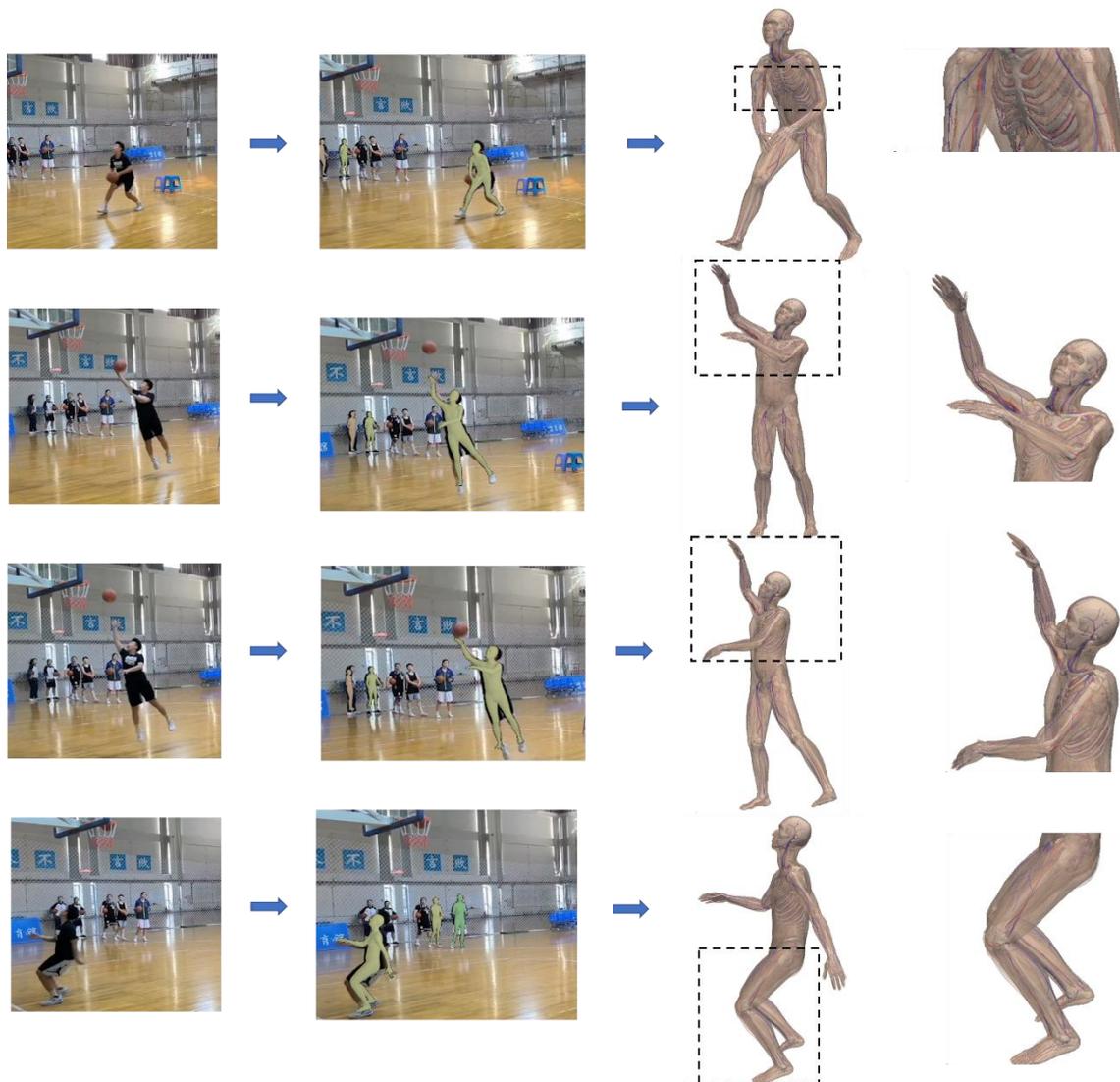


Fig. 6 Video sequence-based modeling of full body anatomy for each moment. (a) A series of action images acquired from the video sequence. (b) The body surface model reconstructed for each moment (c) The human resultant anatomical model corresponding to each moment. Arteries and veins, soft tissues and bones are shown in different colors (d) Local zoomed-in view of the dashed boxes in (c).

4. Conclusion

In this paper, we present an algorithm for video-based full-anatomy reconstruction of personalized digital twin human models in motion. We use a deep neural network to produce the body surface models of each time moment, and develop an internal anatomy method by registering a deformable digital human anatomy atlas to the body surfaces. The integrated articulated deformation of bones and soft tissues produces anatomically realistic models of the moving anatomy, and the intersection removal method generates high-quality mesh for subsequent FEM analysis. The method in this paper provides a novel technical tool for digitally evaluating motion results and performing personalized motion mechanics simulations.

Acknowledgment

This work was supported in part by the National Key Research and Development Program No. 2020YFB1711500, 2020YFB1711501, and 2020YFB1711503, the general program of the National

Natural Science Fund of China (No. 81971693, 61971445), the funding of Dalian Engineering Research Center for Artificial Intelligence in Medical Imaging, Hainan Province Key Research and Development Plan ZDYF2021SHFZ244, the Fundamental Research Funds for the Central Universities (No. DUT22YG229), the funding of Liaoning Key Lab of IC & BME System and Dalian Engineering Research Center for Artificial Intelligence in Medical Imaging.

References

- [1] C. Lassner, J. Romero, M. Kiefel, F. Bogo, M. J. Black, and P. V. Gehler, "Unite the People: Closing the Loop Between 3D and 2D Human Representations," in 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Honolulu, HI, Jul. 2017, pp. 4704–4713. DOI: 10.1109/CVPR.2017.500.
- [2] A. Kanazawa, M. J. Black, D. W. Jacobs, and J. Malik, "End-to-End Recovery of Human Shape and Pose," in 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, Jun. 2018, pp. 7122–7131. DOI: 10.1109/CVPR.2018.00744.
- [3] F. Bogo, A. Kanazawa, C. Lassner, P. Gehler, J. Romero, and M. J. Black, "Keep It SMPL: Automatic Estimation of 3D Human Pose and Shape from a Single Image," in Computer Vision – ECCV 2016, vol. 9909, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds. Cham: Springer International Publishing, 2016, pp. 561–578. DOI: 10.1007/978-3-319-46454-1_34.
- [4] N. Kolotouros, G. Pavlakos, M. Black, and K. Daniilidis, "Learning to Reconstruct 3D Human Pose and Shape via Model-Fitting in the Loop," in 2019 IEEE/CVF International Conference on Computer Vision (ICCV), Seoul, Korea (South), Oct. 2019, pp. 2252–2261. doi: 10.1109/ICCV.2019.00234.
- [5] Y. Cai, L. Ge, J. Cai, and J. Yuan, "Weakly-Supervised 3D Hand Pose Estimation from Monocular RGB Images," in Computer Vision – ECCV 2018, vol. 11210, V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, Eds. Cham: Springer International Publishing, 2018, pp. 678–694. DOI: 10.1007/978-3-030-01231-1_41.
- [6] L. Ge et al., "3D Hand Shape and Pose Estimation From a Single RGB Image," in 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Long Beach, CA, USA, Jun. 2019, pp. 10825–10834. DOI: 10.1109/CVPR.2019.01109.
- [7] Y. Feng, F. Wu, X. Shao, Y. Wang, and X. Zhou, "Joint 3D Face Reconstruction and Dense Alignment with Position Map Regression Network," in Computer Vision – ECCV 2018, vol. 11218, V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, Eds. Cham: Springer International Publishing, 2018, pp. 557–574. DOI: 10.1007/978-3-030-01264-9_33.
- [8] S. Sanyal, T. Bolkart, H. Feng, and M. J. Black, "Learning to Regress 3D Face Shape and Expression From an Image Without 3D Supervision," in 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Long Beach, CA, USA, Jun. 2019, pp. 7755–7764. DOI: 10.1109/CVPR.2019.00795.
- [9] Y. Rong, T. Shiratori, and H. Joo, "FrankMocap: A Monocular 3D Whole-Body Pose Estimation System via Regression and Integration," in 2021 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW), Montreal, BC, Canada, Oct. 2021, pp. 1749–1759. doi: 10.1109/ICCVW54120.2021.00201.
- [10] Y. Xiu, J. Yang, D. Tzionas, and M. J. Black, "ICON: Implicit Clothed humans Obtained from Normals," in 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), New Orleans, LA, USA, Jun. 2022, pp. 13286–13296. DOI: 10.1109/CVPR52688.2022.01294.
- [11] D. Anguelov, P. Srinivasan, D. Koller, S. Thrun, J. Rodgers, and J. Davis, "SCAPE: shape completion and animation of people," *ACM Trans. Graph.*, vol. 24, no. 3, Art. no. 3, Jul. 2005, DOI: 10.1145/1073204.1073207.
- [12] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black, "SMPL: a skinned multi-person linear model," *ACM Trans. Graph.*, vol. 34, no. 6, pp. 1–16, Nov. 2015, DOI: 10.1145/2816795.2818013.
- [13] M. Keller, S. Zuffi, M. J. Black, and S. Pujades, "OSSO: Obtaining Skeletal Shape from Outside," in 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), New Orleans, LA, USA, Jun. 2022, pp. 20460–20469. doi: 10.1109/CVPR52688.2022.01984.

- [14] D. Ali-Hamadi et al., “Anatomy transfer,” *ACM Trans. Graph.*, vol. 32, no. 6, pp. 1–8, Nov. 2013, DOI: 10.1145/2508363.2508415.
- [15] A. Sujar, J. J. Casafranca, A. Serrurier, and M. Garcia, “Real-time animation of human characters’ anatomy,” *Computers & Graphics*, vol. 74, pp. 268–277, Aug. 2018, DOI: 10.1016/j.cag.2018.05.025.
- [16] L. Kavan, S. Collins, J. Žára, and C. O’Sullivan, “Geometric skinning with approximate dual quaternion blending,” *ACM Trans. Graph.*, vol. 27, no. 4, pp. 1–23, Oct. 2008, DOI: 10.1145/1409625.1409627.
- [17] B. H. Le and J. K. Hodgins, “Real-time skeletal skinning with optimized centers of rotation,” *ACM Trans. Graph.*, vol. 35, no. 4, pp. 1–10, Jul. 2016, DOI: 10.1145/2897824.2925959.
- [18] N. Magnenat-Thalmann, R. Laperriere, and D. Thalmann. Joint-Dependent Local Deformations for Hand Animation and Object Grasping. *P~vc. Graphics Interface*, 1988, pp. 26-33.
- [19] H. Wang, X. Sun, L. Huo, X. Tang, and C. Liu, “Construction of Deformable Trunk Atlas of Chinese Human Based on Multiple PET/CT Images: Preliminary Results,” in *Digital Human Modeling. Applications in Health, Safety, Ergonomics, and Risk Management: Ergonomics and Design*, vol. 10286, V. G. Duffy, Ed. Cham: Springer International Publishing, 2017, pp. 69–77. DOI: 10.1007/978-3-319-58463-8_7.
- [20] M. Kocabas, N. Athanasiou, and M. J. Black, “VIBE: Video Inference for Human Body Pose and Shape Estimation,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Seattle, WA, USA, Jun. 2020, pp. 5252–5262. DOI: 10.1109/CVPR42600.2020.00530.
- [21] K. He, X. Zhang, S. Ren, and J. Sun, “Deep Residual Learning for Image Recognition,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA, Jun. 2016, pp. 770–778. DOI: 10.1109/CVPR.2016.90.