

Supplementary Material for A Dual-Masked Auto-Encoder for Robust Motion Capture with Spatial-Temporal Skeletal Token Completion

Junkun Jiang¹, Jie Chen^{1,*}, Yike Guo^{1,2}

¹Department of Computer Science, Hong Kong Baptist University, Hong Kong SAR, China

²Data Science Institute, Imperial college London, London, UK

{csjkjiang, chenjie, yikeguo}@comp.hkbu.edu.hk, yg@doc.ic.ac.uk

1 NETWORK STRUCTURE

As introduced in the main paper, the proposed D-MAE consists of one Encoder for transforming encoded tokens into the latent embedding space and one Decoder for predicting the masked token and completing the 3D skeletal motion. Both the Encoder and the Decoder are vanilla Transformers [11] consisting of Multi-head Self Attention *Attn* blocks and Multilayer Perceptron *MLP* blocks. D-MAE takes a masked motion sequence as input. We treat each joint’s 3D coordinates $s \in \mathbb{R}^3$ as an independent token encoded by the proposed dual-position encoders (the signal encoder γ_s for 3D spatial coordinates s and the context encoder γ_v for each joint’s structural and temporal context $v \in \{J, T\}$, please refer to Eq. 5,6,7 in the main paper). Following [5], a linear layer *Enc-to-Dec* is used to bridge the Encoder and the Decoder. To make the final prediction, the masked tokens are then projected by a linear layer *Prediction* to the signal dimension.

Table 1 shows the Transformer encoder-decoder structure. Also we will release the whole code of our multi-view multi-person motion capture system in the project page: https://github.com/HKBU-VSComputing/2022_MM_DMAE-Mocap.

2 TRAINING DETAILS

D-MAE completes the motion sequences via encoding motion signals in both skeletal-structural and temporal domains. We further provide the training configurations for its easy reproduction.

Table 2 summarizes the training configurations. Most of the configurations are shared by *training* and *fine-tune* processes, without specific tuning. Due to the flexibility of the D-MAE’s encoding mechanism, we can simply transfer one skeleton labeling definition to another definition without any changes on our motion capture system.

3 ADDITIONAL VISUAL EXAMPLES

We provide additional visual comparisons with Dong et al [3] and Zhang et al [12] on the Shelf dataset and one of the clip from the BU-Mocap dataset. We also provide the synthetic video for continuous evaluation in <https://youtu.be/zCOIGwWISoI>. Note that, we only train our model on the Shelf dataset, and directly evaluate the whole system on the BU-Mocap dataset. Without any additional tuning, our system outperforms than [3, 12] in the qualitative comparisons, which also indicates that with the help of the dual directional encoding, the D-MAE learns how to model the human motion.

Table 1: The details of the proposed D-MAE network structure. ‘num’ is short for ‘number’. ‘dim’ is short for ‘dimension’.

<i>Encoder</i>	num
Depth	6
Attn dim	256
Attn head	8
Attn dim per head	64
MLP dim	512
<i>Enc-to-Dec</i>	num
In dim	256
Out dim	128
<i>Decoder</i>	num
Depth	6
Attn dim	128
Attn head	8
Attn dim per head	64
MLP dim	512
<i>Prediction</i>	num
In dim	128
Out dim	3

Table 2: Configurations for training and fine-tune. [†] We use the linear *lr* scaling rule [4]: i.e., $lr = base_lr \times batch_size/256$.

Configuration	training	fine-tune
optimizer		AdamW [8]
optimizer momentum		$\beta_1, \beta_2 = 0.9, 0.999$
weight decay		0.05
learning rate schedule		cosine decay [7]
warmup epochs [4]	30	5
augmentation	Normalization, RandomRotation	
gradient clipping	0.02	0.005
drop out		\times
base learning rate [†]	2e-4	8e-4
batch size	256	128
epoch size	1000	500

As shown in Figure 4, the 3D skeletons reconstructed by [3] have different colors. The blue skeleton at the first frame is painted into pink at the second frame. This phenomenon indicates [3] has weak identification consistency, i.e. during the continuous frames, [3] reconstructs the same candidate with different identities. [12]

* Corresponding author: Jie Chen.

Table 3: Current records for BU-Mocap dataset. Each of them captures multi-view multi-person motion with massive interactions and occlusions.

Clip name	length (s)	actor num.	view num.	depth	IMU
coop-1	120	2	5	✓	✗
coop-2	120	2	5	✓	✗
coop-3	120	4	5	✓	✗
coop-4	120	4	4	✓	✗
coop-5	120	5	5	✓	✓

reconstructs failed on the feet due to limited 2D clues under severe occlusion. Ours can reconstruct well because of the D-MAE’s completion. The *Adaptive Triangulation* module first filters the inaccurate 2D detection. The D-MAE completes the missing 3D parts via dual directional encoding.

As shown in Figure 5, Dong et al [3] reconstructs 3D poses in the good skeletal structure while one of the reprojection is misaligned to its human body. Zhang et al [12] fails to reconstruct 3D poses from the third frame to the fifth frame. The joint *left shoulder* of the blue skeleton is assigned to another skeleton erroneously making the twisted failure reconstruction. The twisted skeletons usually appear in the rotation case. When the actor spins, the 3D reconstruction becomes twisted. Ours performs better than [3, 12] in this scene.

4 FINE-TUNE ON BU-MOCAP

As shown in Table. 3 in the main paper, we further evaluate the generalization of the proposed D-MAE and compare it with a learning-based method [6]. Row 3 and 4 demonstrate the result for both of the two models trained on the Shelf [1] dataset respectively. Row 5 and 6 demonstrate the fine-tune results. We fine-tune models on the BU-Mocap (70% for training, 30% for testing). Specifically, we follow [6]’s statement that to use HRNet [10] as the 2D pose detector. Due to BU-Mocap’s full annotation (all actors in every frame are labelled), the data is much more sufficient than the Shelf dataset leading to earlier fitting. Both of ours and [6]’s training epoch is set to 600. Other configurations are consistent with each statement.

5 BU-MOCAP DATASET

As the one of our main contributions, the proposed BU-Mocap includes but is not limited to the following modalities: 1) RGB color image, 2) depth map, 3) motion trajectory recorded by inertial sensors (inertial measurement units, IMUs), 3) 3D motion coordinate, 4) reconstructed point cloud. Table 3 summarizes the existing motion clips. As shown in Table 3, 5 clips with total 600 seconds are provided. The capturing rate is set to 30 frame per second (fps). The inertial sensor’s recording rate is set to 120 Hz.

Figure 1 demonstrates the shooting set for the BU-Mocap dataset. We setup 5 RGB cameras [9] to record motion. Before every shooting, we use a chessboard to get the extrinsic parameter of each camera. As for the filming script, we invite up to 5 volunteers to imitate the dancing choreography on the screen. Each clip has its own dancing style, including *Jazz dance*, *Cha Cha*, *Free dance* and so on. We design each actor’s appearance, i.e. 5 actors wear different

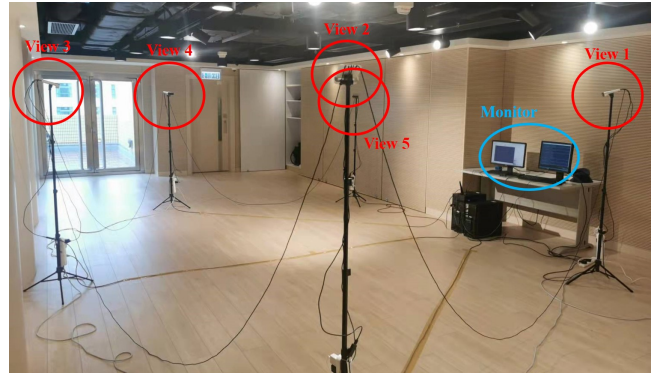


Figure 1: Demonstration of the shooting set for the BU-Mocap dataset. Five Azure Kinect cameras [9] are placed to surround the capture area. Two monitors are used for motion display.

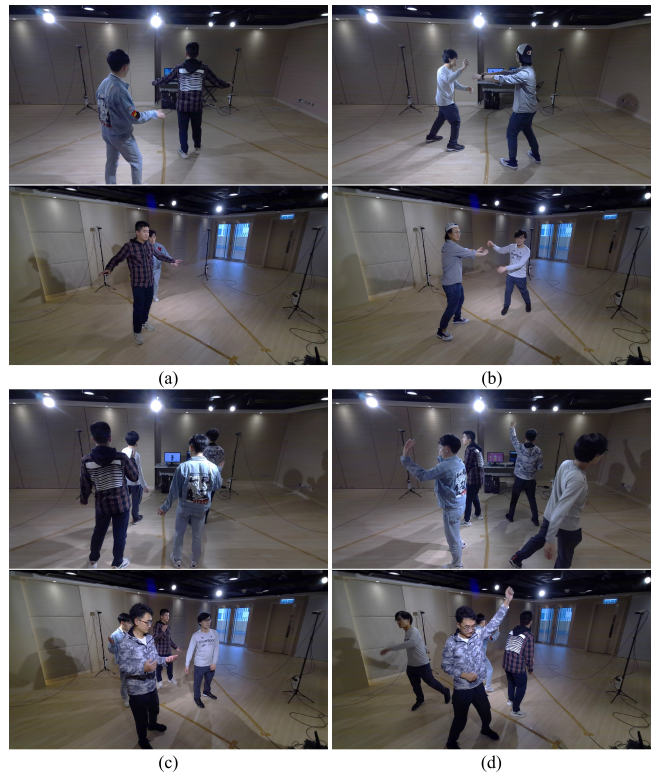


Figure 2: Visual examples of the BU-Mocap dataset. Each sub-figure demonstrates the different actor performance from 2 views.

clothes. For example, actors are asked to wear T-shirts and shorts, loose jacket and trousers respectively. Figure 2 demonstrate the capturing scene. Each sub-figure illustrates one of the dancing clips. The bottom sub-figures are under heavy occlusions and have more interactions. In the future, we will release it to the public and also further extend the BU-Mocap dataset.

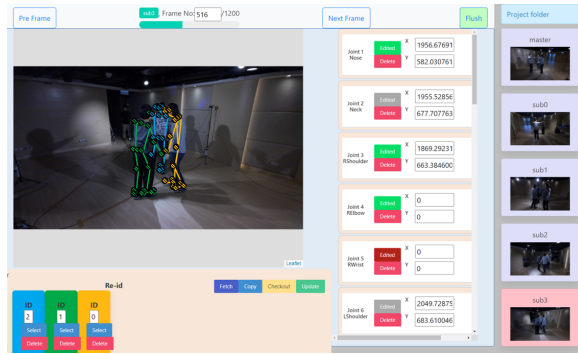


Figure 3: Demonstration of the annotation system. The system is designed for collaborative working. The user-friendly UI allows participants labeling keypoints easily.

6 ANNOTATION SYSTEM

To achieve accurate ground-truth, we develop a multi-view multi-person annotation system (see Figure 3) to manually correct the wrong 2D detection as well as 3D reconstruction. We first adopt off-the-shelf 2D human pose estimator [2] to detect 2D poses from each observation. Then the multi-view multi-person 3D motion capture framework proposed in the paper is used for 3D motion reconstruction. Next we reproject those reconstructed 3D motion back to 2D images. The 2D reprojection is displayed on the annotation system waiting for correction.

As shown in Figure 3, the user interface (UI) of the annotation system consists of 4 parts: 1) the main canvas, demonstrating identity-aware keypoints for each candidate. We can easily drag and delete keypoints on the main canvas; 2) the Re-id panel lying below the main canvas. This panel is used for changing each candidate’s identity; 3) the joint panel lying next to the main canvas. Each column in the joint panel demonstrates the 2D joint location. We can modify it by directly typing XY values; 4) the view-switch panel lying next to the joint panel. We can switch the observing view of the main canvas by right double-clicking other views on this panel.

REFERENCES

- [1] Vasileios Belagiannis, Sikandar Amin, Mykhaylo Andriluka, Bernt Schiele, Nassir Navab, and Slobodan Ilic. 2015. 3D pictorial structures revisited: Multiple human pose estimation. *IEEE transactions on pattern analysis and machine intelligence* 38, 10 (2015), 1929–1942.
- [2] Zhe Cao, Gines Hidalgo, Tomas Simon, Shih-En Wei, and Yaser Sheikh. 2019. OpenPose: realtime multi-person 2D pose estimation using Part Affinity Fields. *IEEE transactions on pattern analysis and machine intelligence* 43, 1 (2019), 172–186.
- [3] Junting Dong, Qi Fang, Wen Jiang, Yurou Yang, Qixing Huang, Hujun Bao, and Xiaowei Zhou. 2021. Fast and Robust Multi-Person 3D Pose Estimation and Tracking from Multiple Views. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).
- [4] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. 2017. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677* (2017).
- [5] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. 2021. Masked autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377* (2021).
- [6] Jiahao Lin and Gim Hee Lee. 2021. Multi-view multi-person 3d pose estimation with plane sweep stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11886–11895.
- [7] Ilya Loshchilov and Frank Hutter. 2016. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983* (2016).
- [8] Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101* (2017).
- [9] Microsoft. 2019. Azure Kinect DK - Build for mixed reality using AI sensors. <https://azure.microsoft.com/en-us/services/kinect-dk/>. Accessed January 29, 2022.
- [10] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. 2019. Deep high-resolution representation learning for human pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 5693–5703.
- [11] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [12] Yuxiang Zhang, Liang An, Tao Yu, Xiu Li, Kun Li, and Yebin Liu. 2020. 4D association graph for realtime multi-person motion capture using multiple video cameras. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1324–1333.

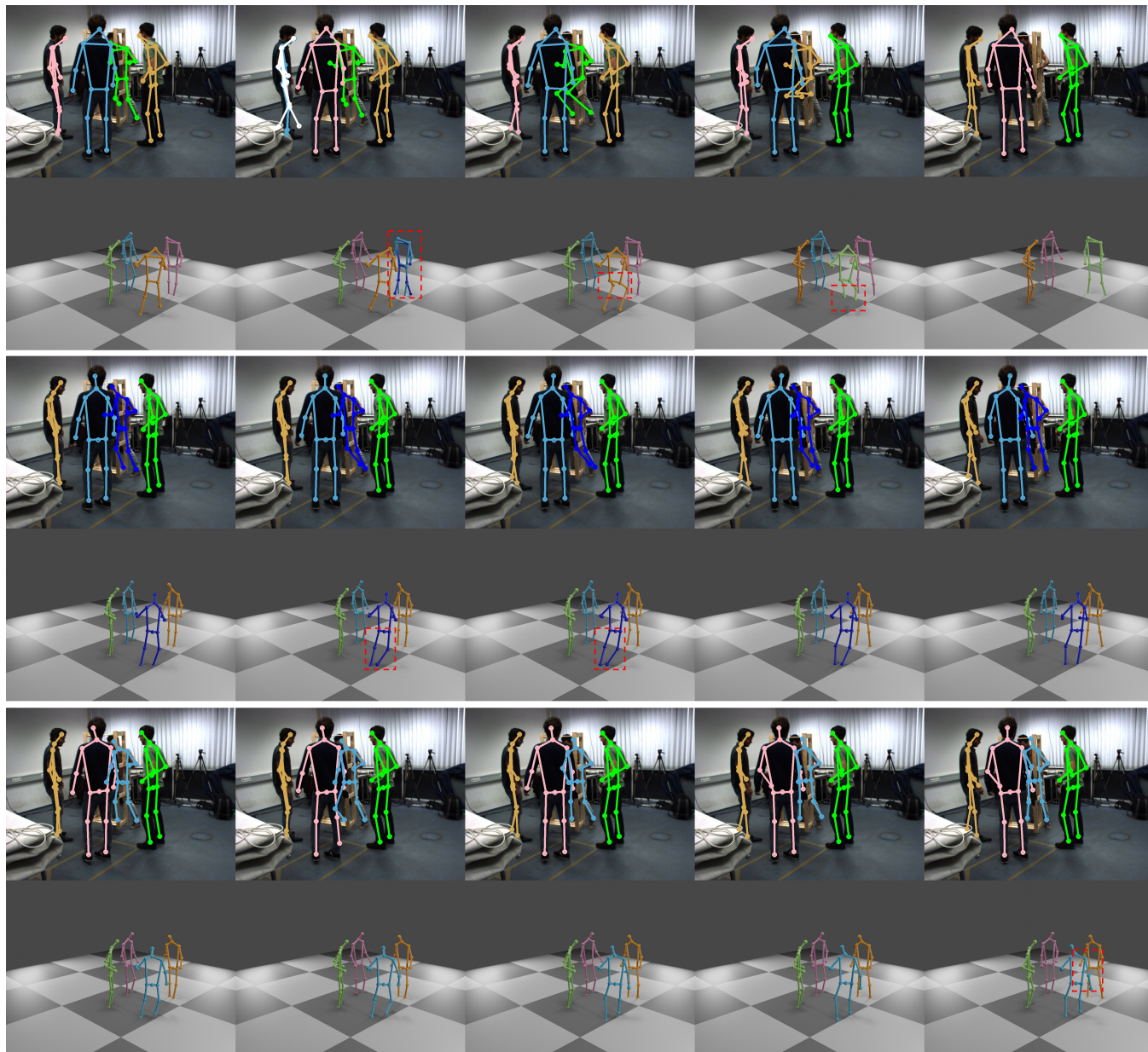


Figure 4: Visual comparisons on the Shelf dataset. We show 3D poses and its reprojection on the 4th view with 5 continuous frames (from top to bottom, Dong et al [3], Zhang et al [12], ours(final)). The red dotted boxes are highlights for attention.

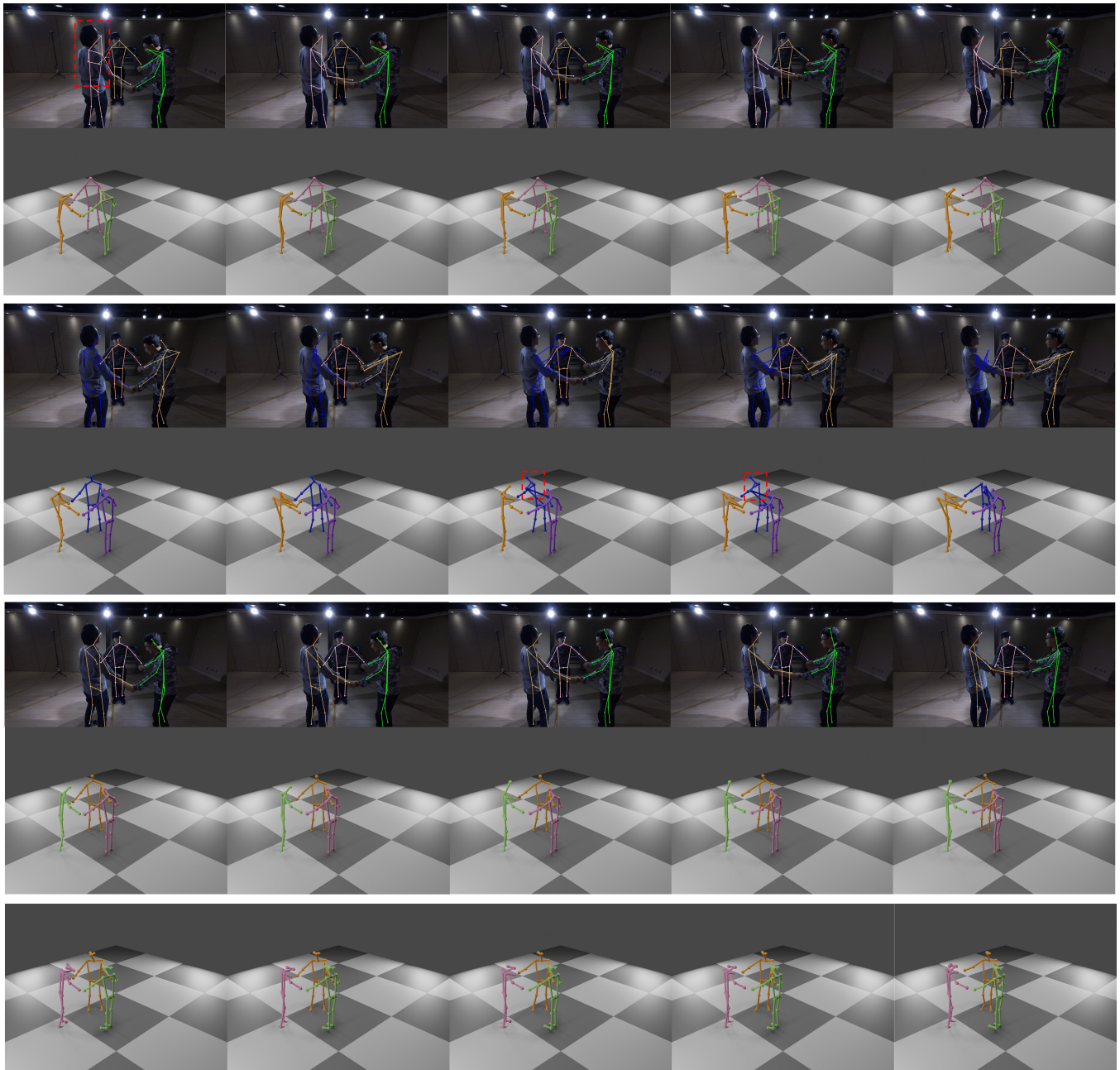


Figure 5: Visual comparisons on the BU-Mocap dataset. We show 3D poses and its reprojection on the 3rd view with 5 continuous frames (from top to bottom, Dong et al [3], Zhang et al [12], ours(final), ground-truth). The red dotted boxes are highlights for attention.