

Egocentric Interaction Capture for Mixed Reality



Siwei Zhang¹



Qianli Ma¹



Yan Zhang¹



Zhiyin Qian¹



Taein Kwon¹



Federica Bogo²



Marc Pollefeys^{1,2}



Siyu Tang¹

1 Computer Vision and Learning Group, ETH Zurich

2 Microsoft



Video from Meta AI



Video from Microsoft Mixed Reality

Third-person View

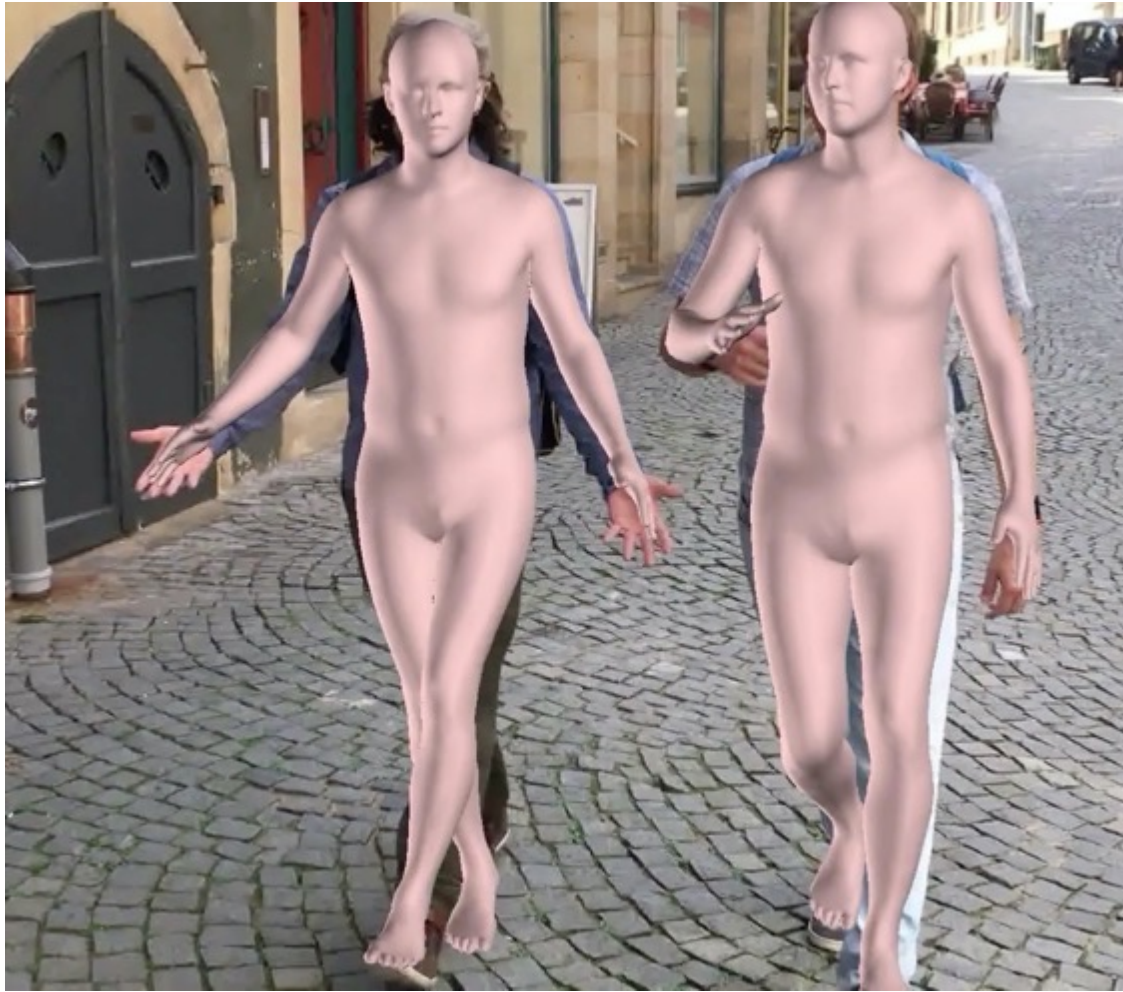


First-person / Egocentric View



Video from Meta AI

Third-person View



First-person / Egocentric View



SPIN (Kolotouros et al.)



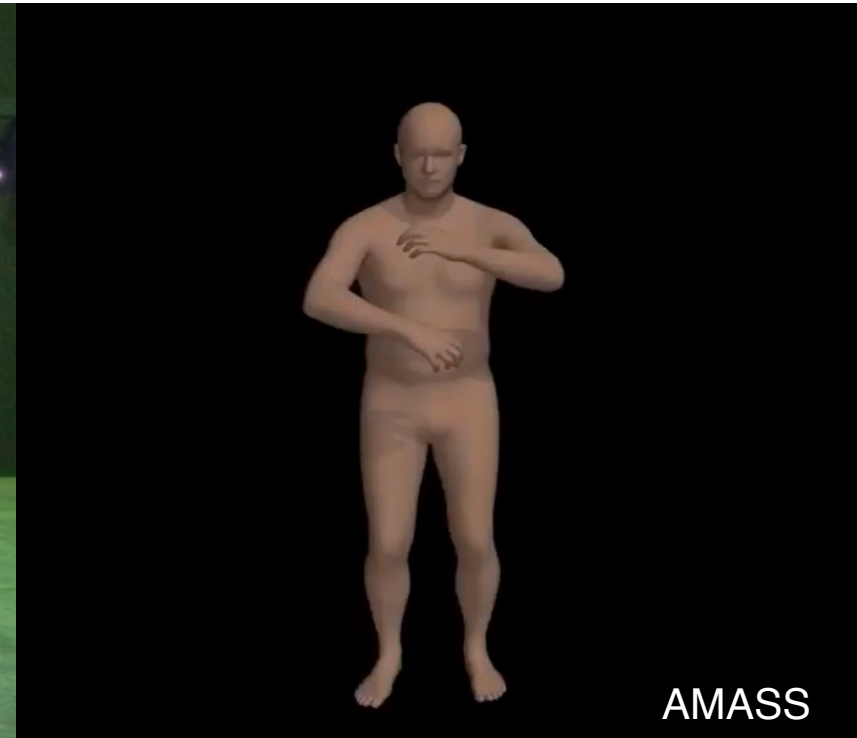
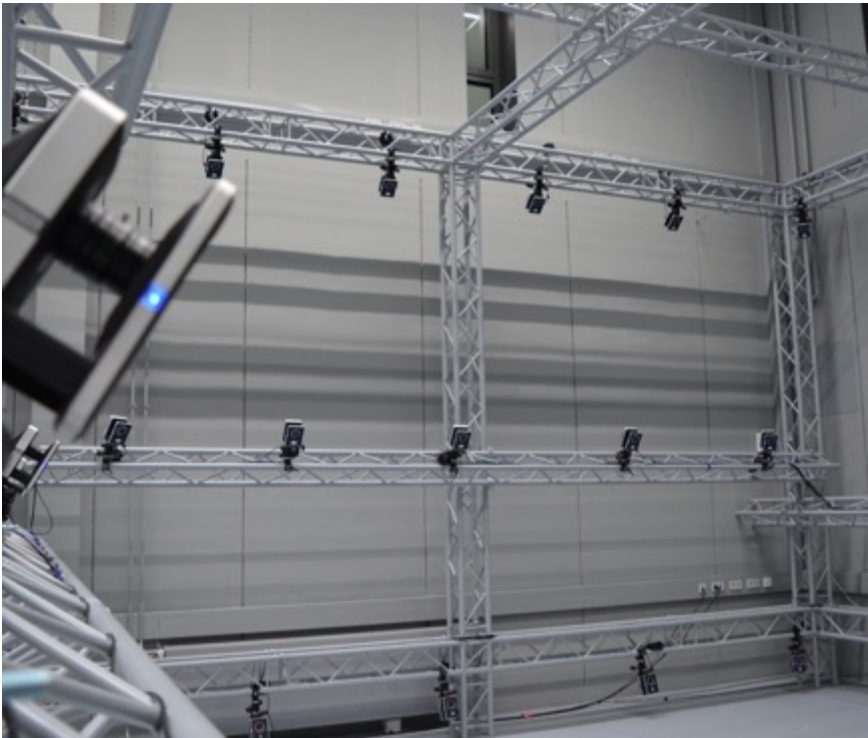


Egocentric View



**How to capture 3D human pose, shape and motions
in 3D scenes?**

Marker-based Motion Capture System



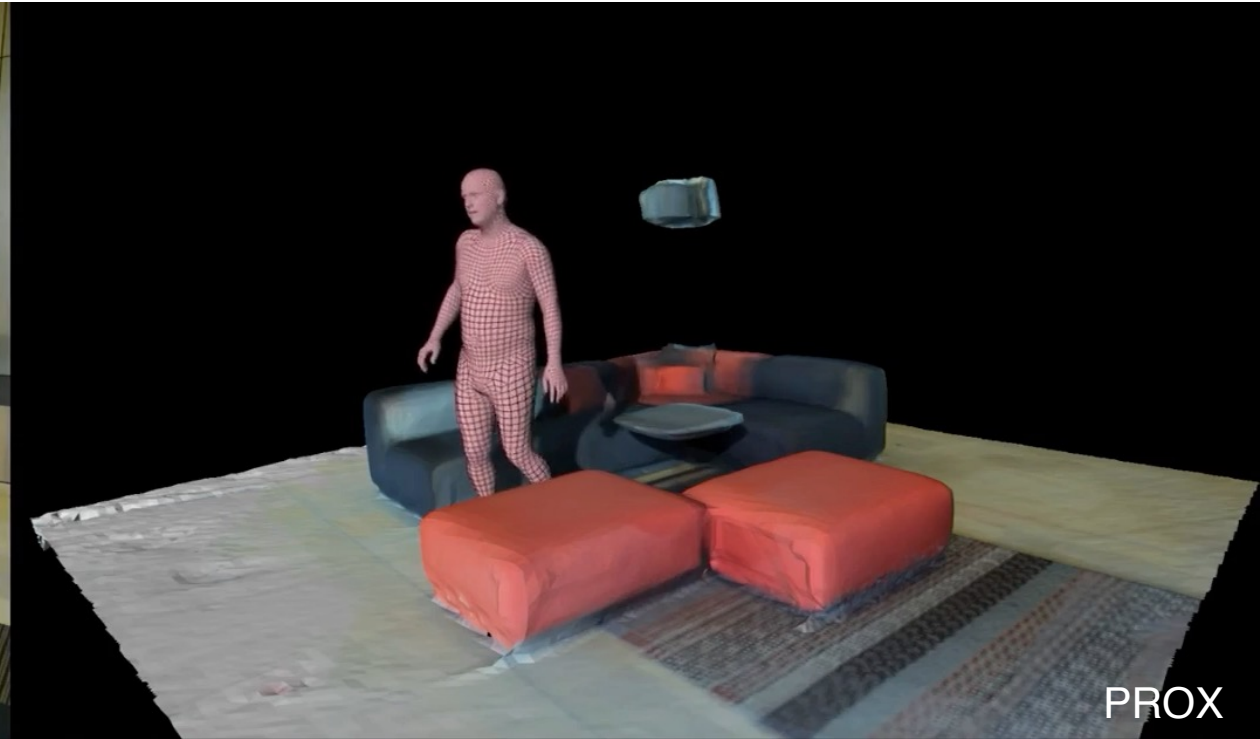
Expensive setup / expert knowledge

[1] <https://ps.is.mpg.de/pages/motion-capture>

[2] <https://sentimentalflow.wordpress.com/2017/01/30/first-blog-post/>

[3] AMASS: Archive of Motion Capture as Surface Shapes, Mahmood et al, ICCV 2019

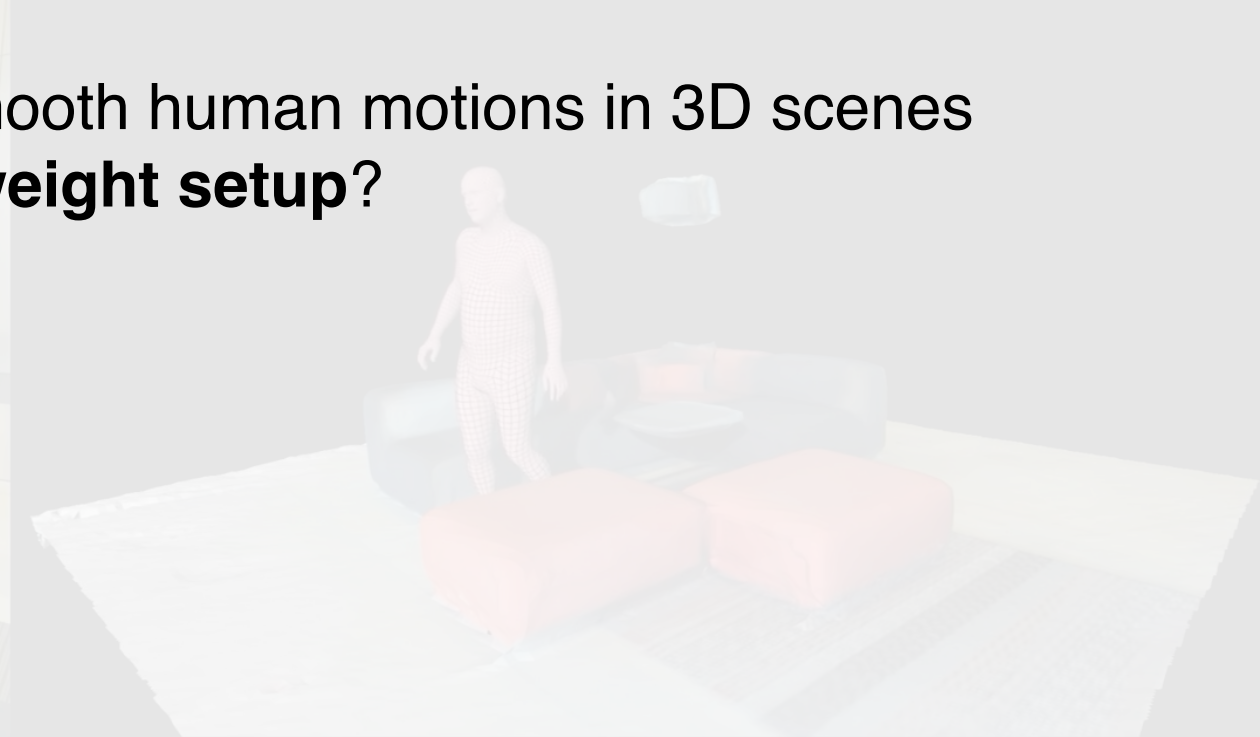
Monocular RGB-D Kinect Setting



Lightweight, noisy motion reconstructions

Monocular RGB-D Kinect Setting

How to reconstruct natural and smooth human motions in 3D scenes with a **lightweight setup**?



Lightweight, noisy motion reconstructions

Outline

- **LEMO: Learning Motion Priors for 4D Human Body Capture in 3D Scenes**

Siwei Zhang, Yan Zhang, Federica Bogo, Marc Pollefeys, Siyu Tang

[ICCV 2021, Oral presentation](#)

- **EgoBody: Human Body Shape and Motion of Interacting People from Head-Mounted Devices**

Siwei Zhang, Qianli Ma, Yan Zhang, Zhiyin Qian, Taein Kwon, Marc Pollefeys, Federica Bogo, Siyu Tang

Outline

- **LEMO: Learning Motion Priors for 4D Human Body Capture in 3D Scenes**

Siwei Zhang, Yan Zhang, Federica Bogo, Marc Pollefeys, Siyu Tang

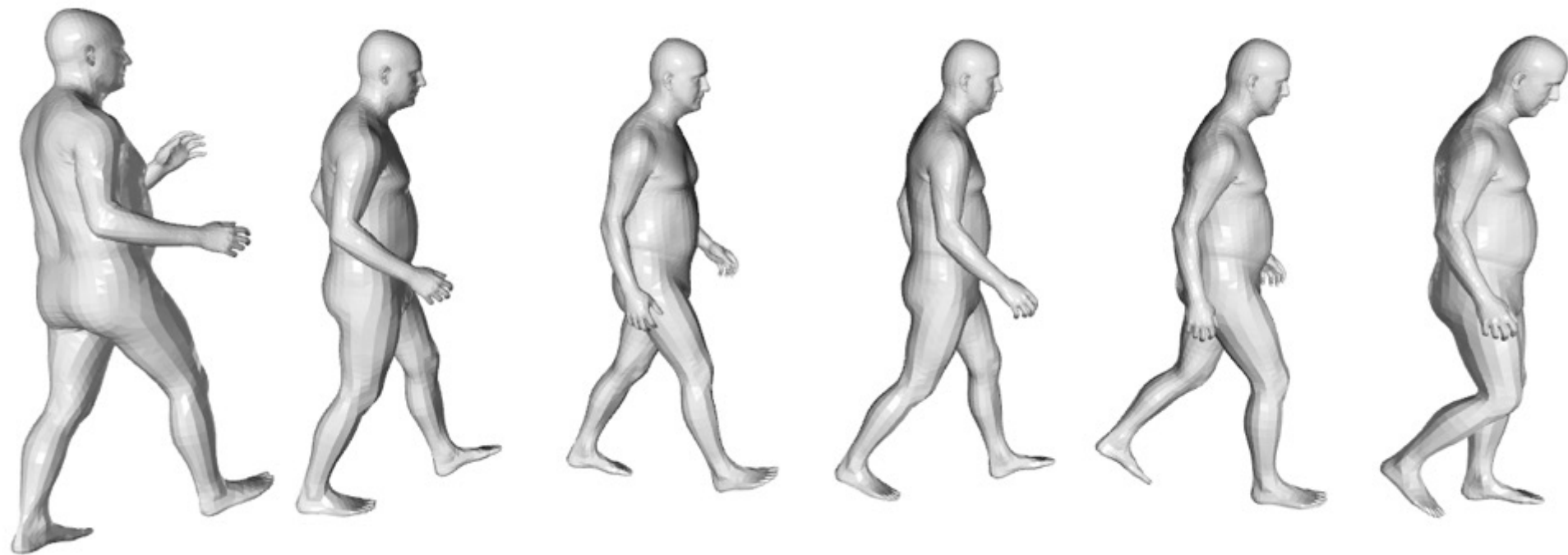
[ICCV 2021, Oral presentation](#)

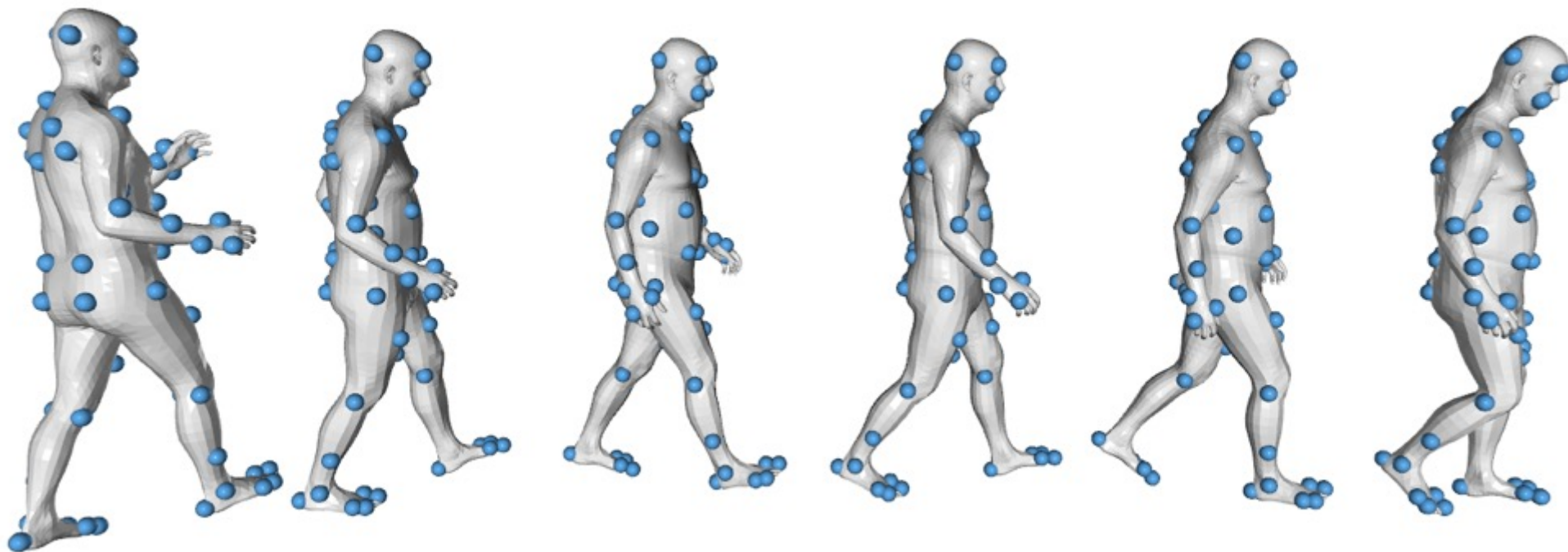
- **EgoBody: Human Body Shape and Motion of Interacting People from Head-Mounted Devices**

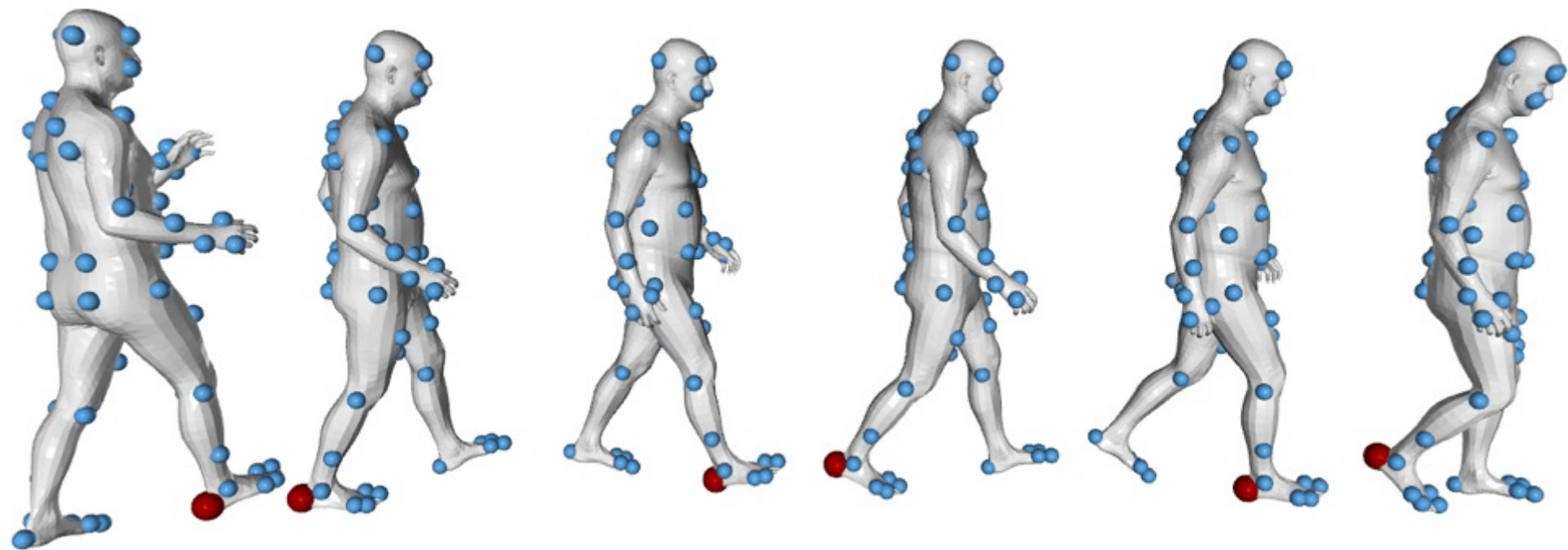
Siwei Zhang, Qianli Ma, Yan Zhang, Zhiyin Qian, Taein Kwon, Marc Pollefeys, Federica Bogo, Siyu Tang

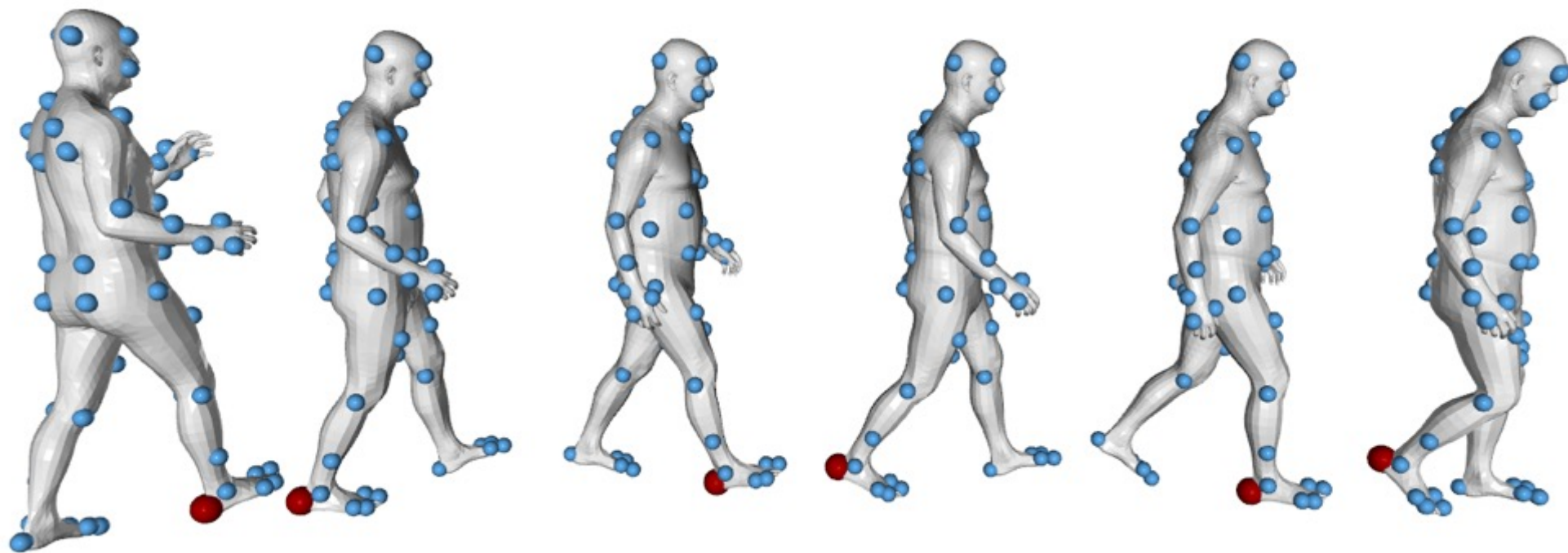
How to reconstruct natural and smooth human motions in 3D scenes
with a **lightweight setup**?

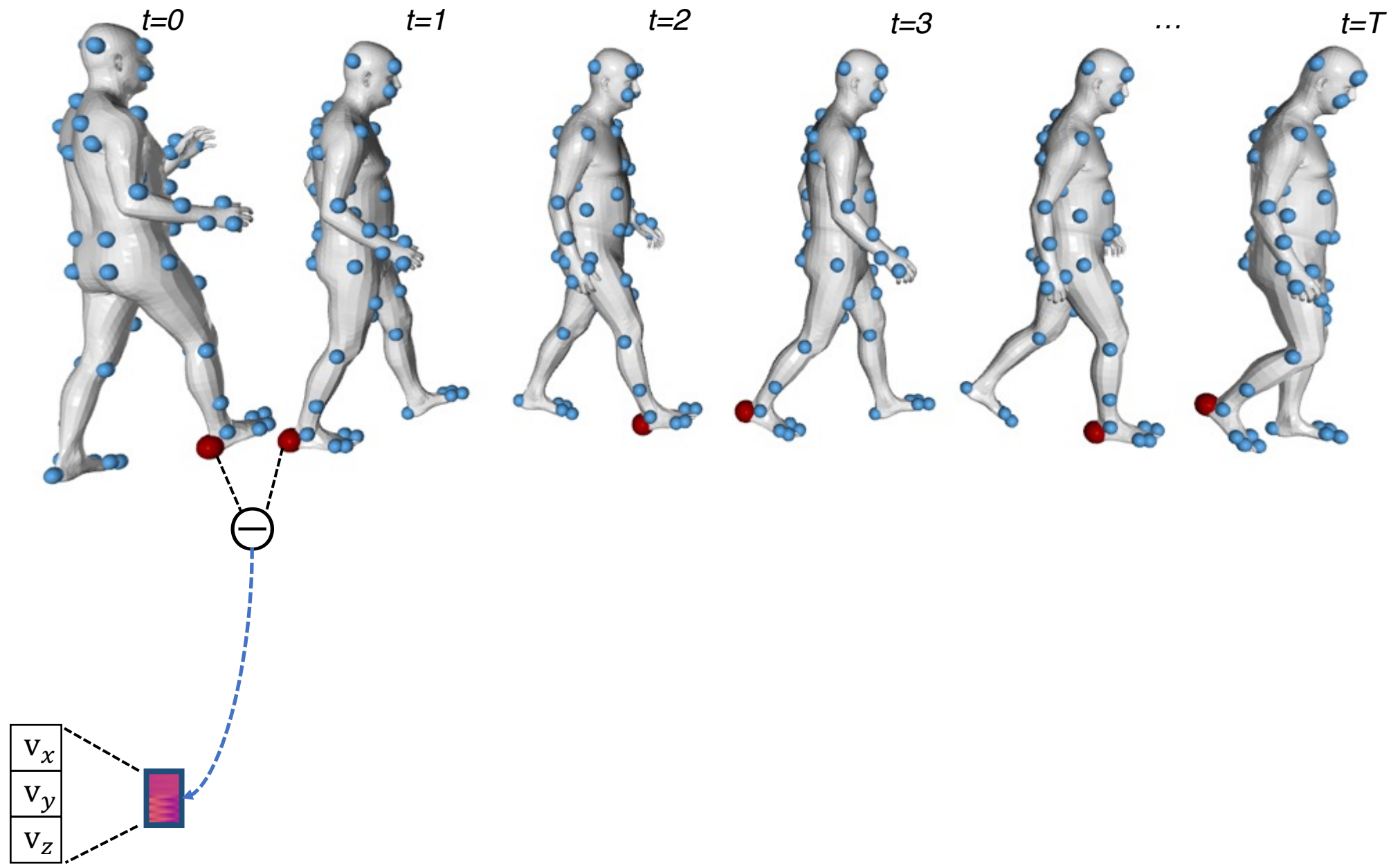
Key insight: learning motion priors from the high quality mocap dataset

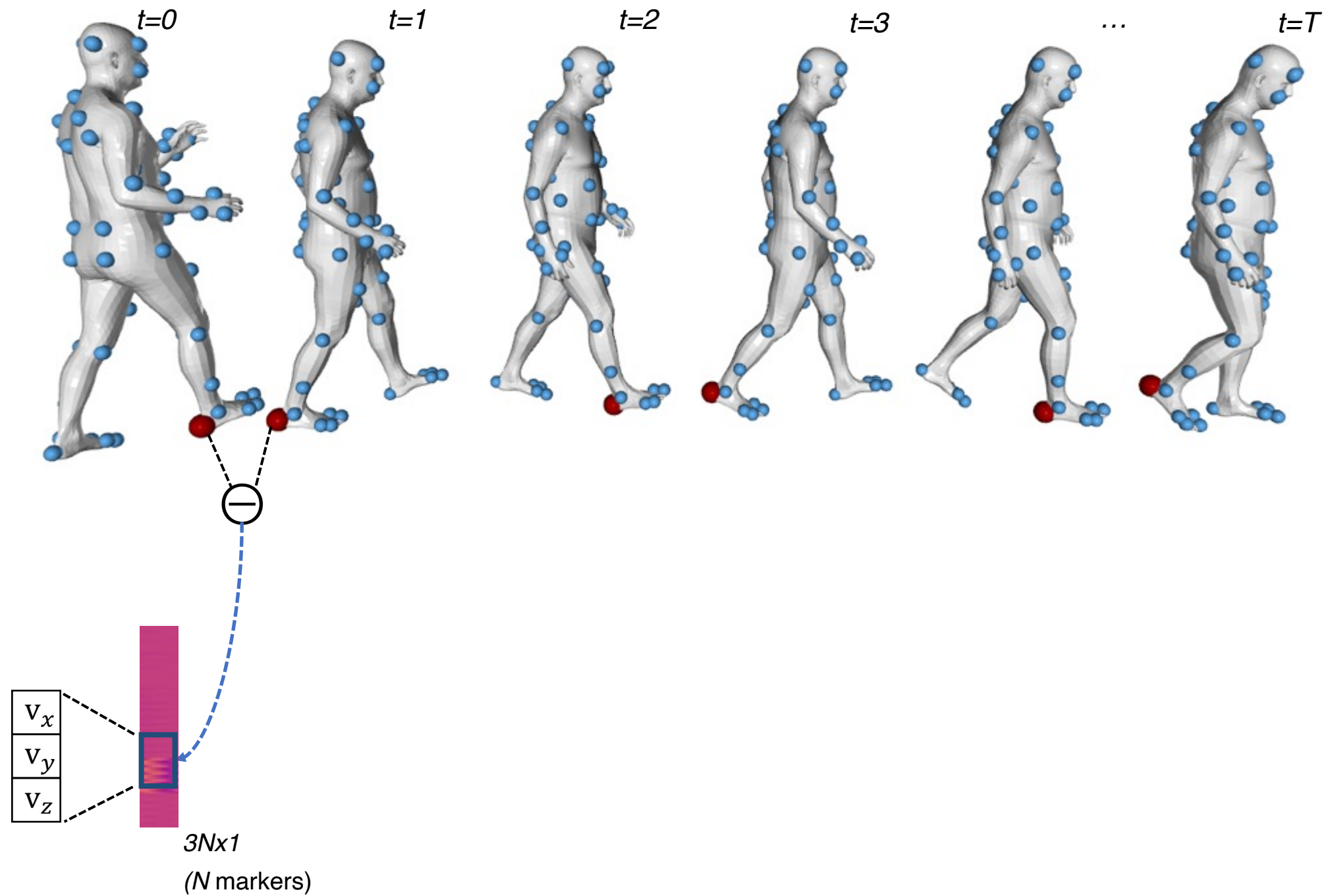


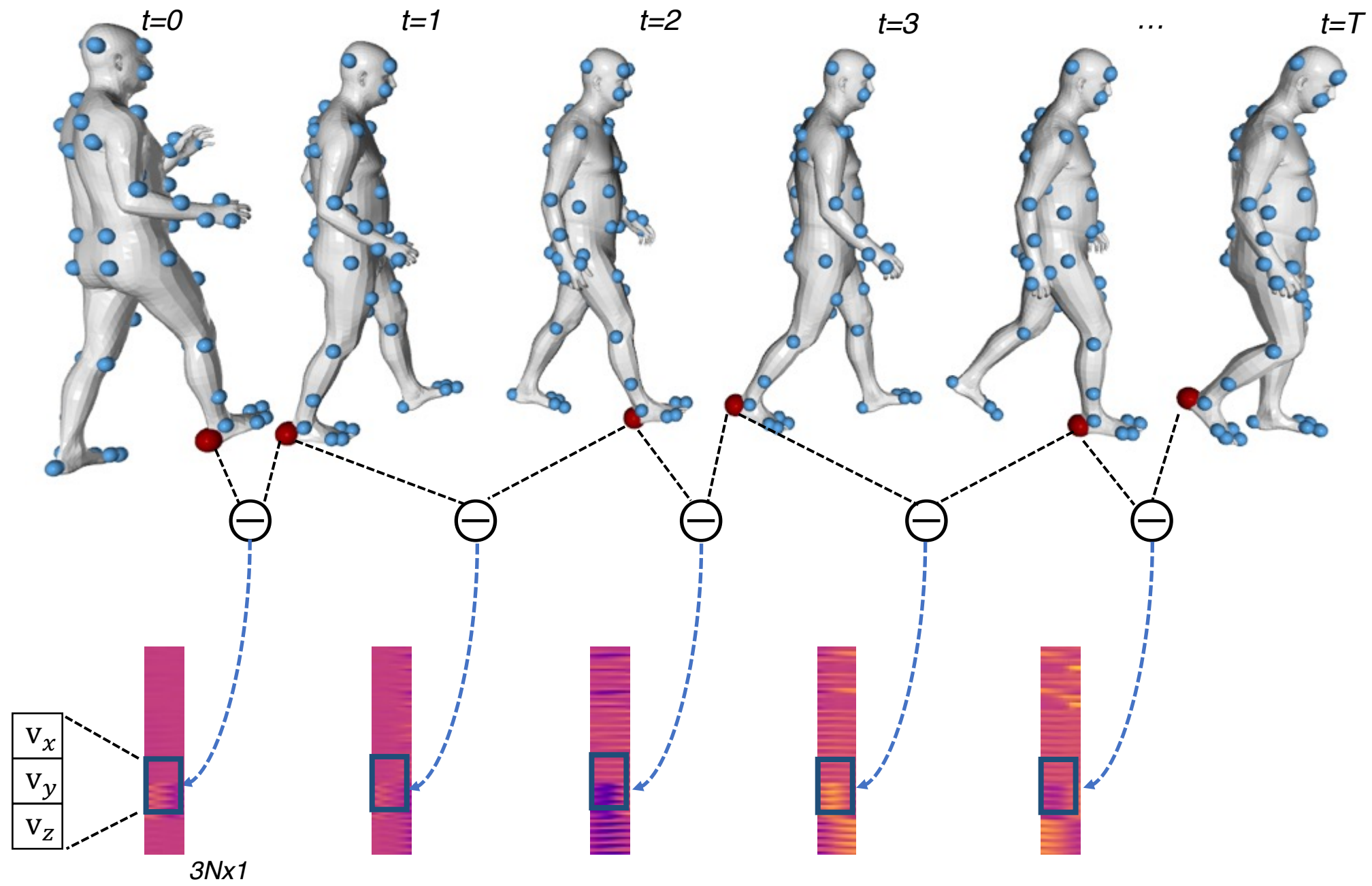


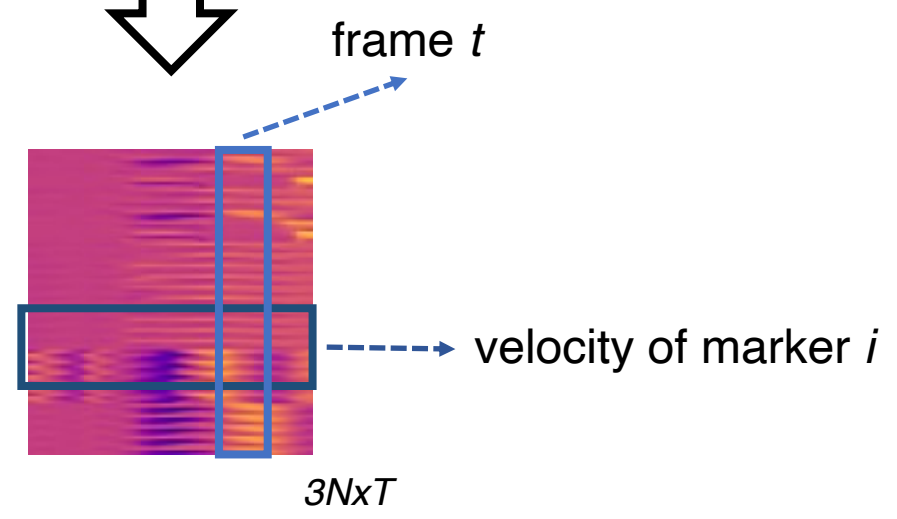
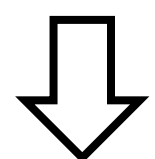
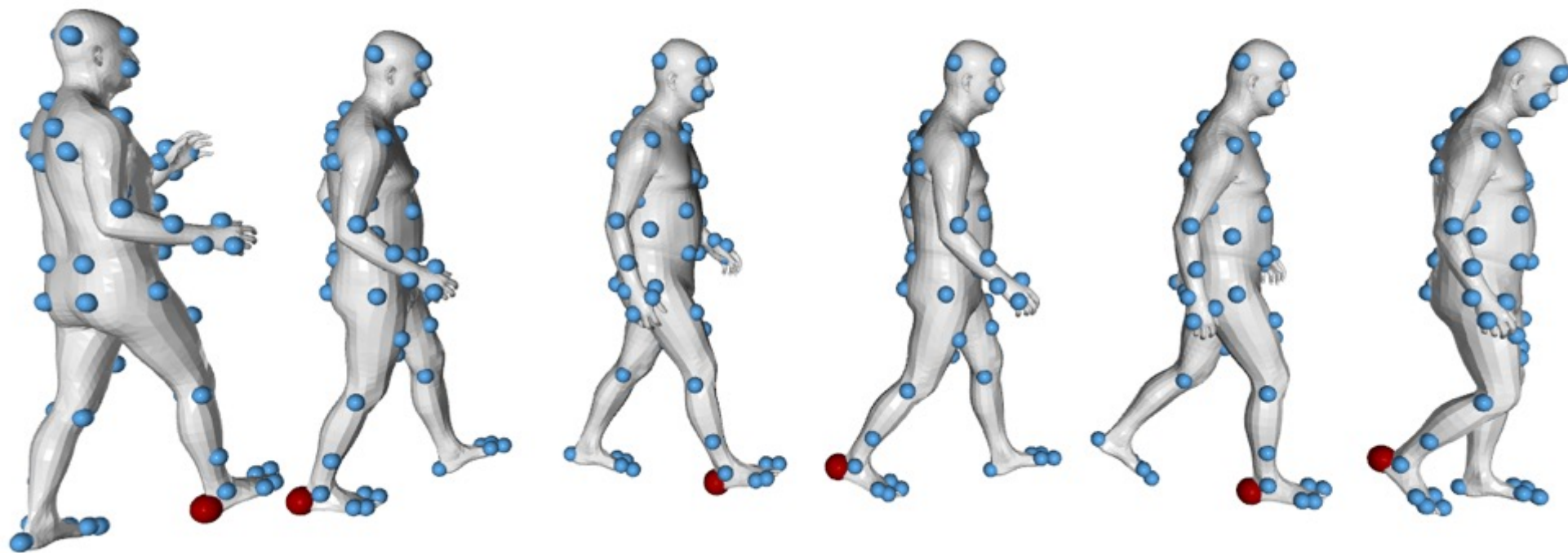


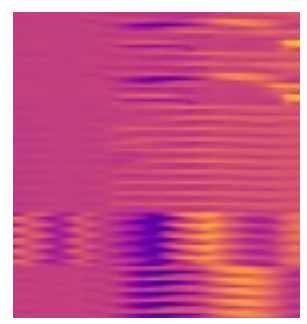
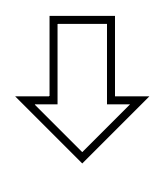
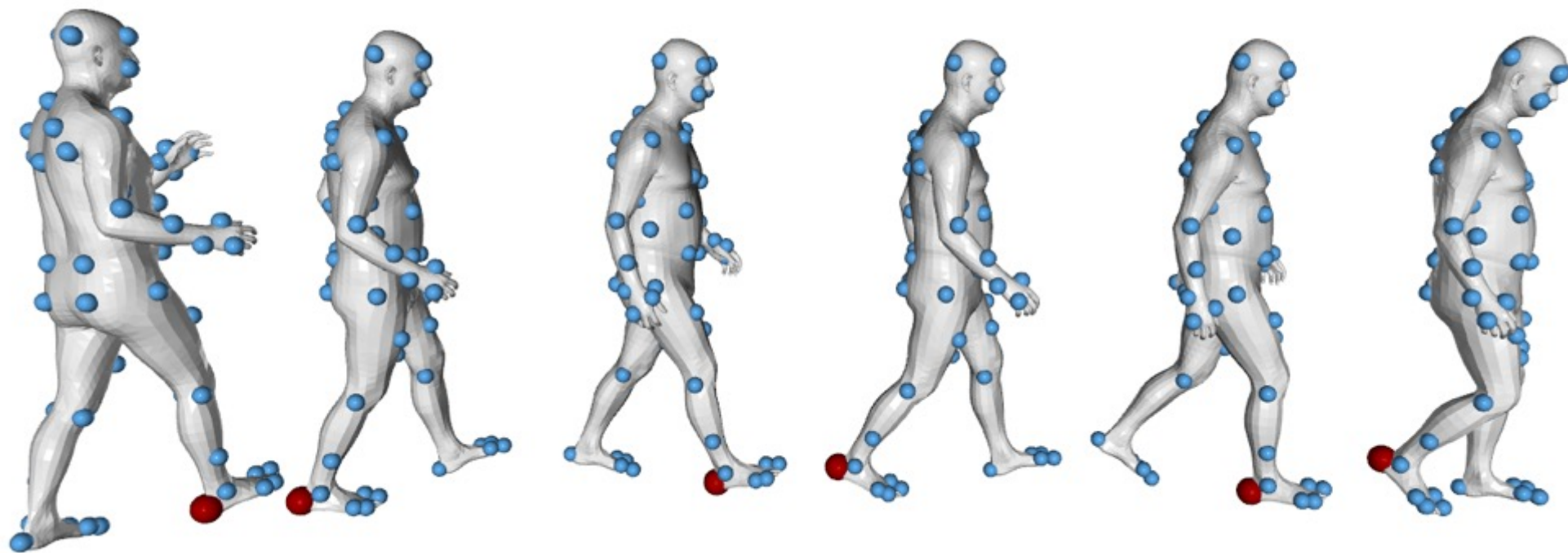


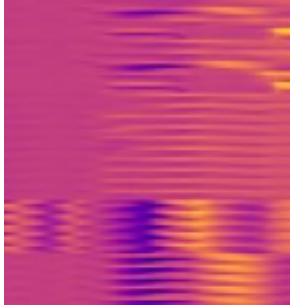


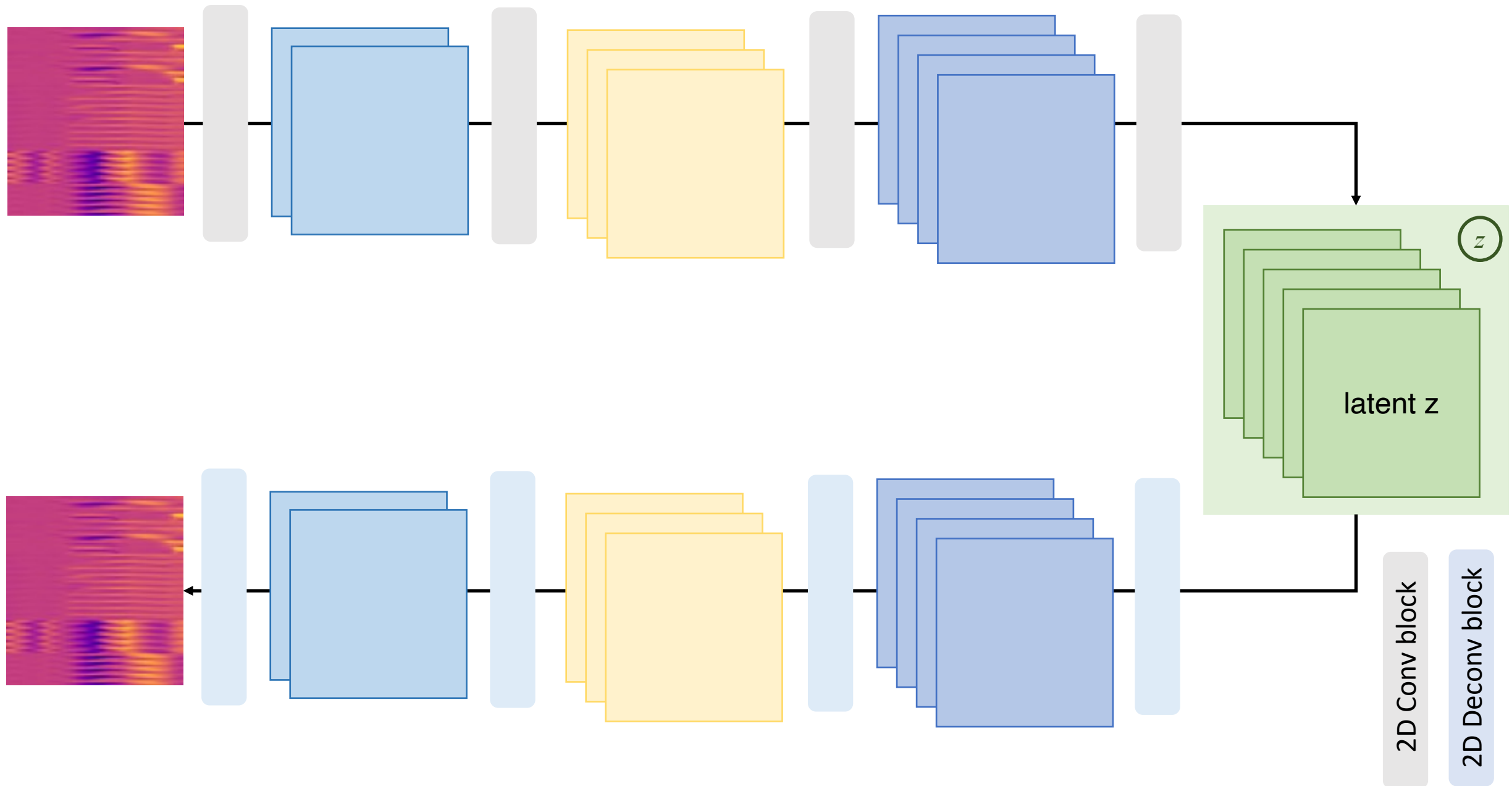


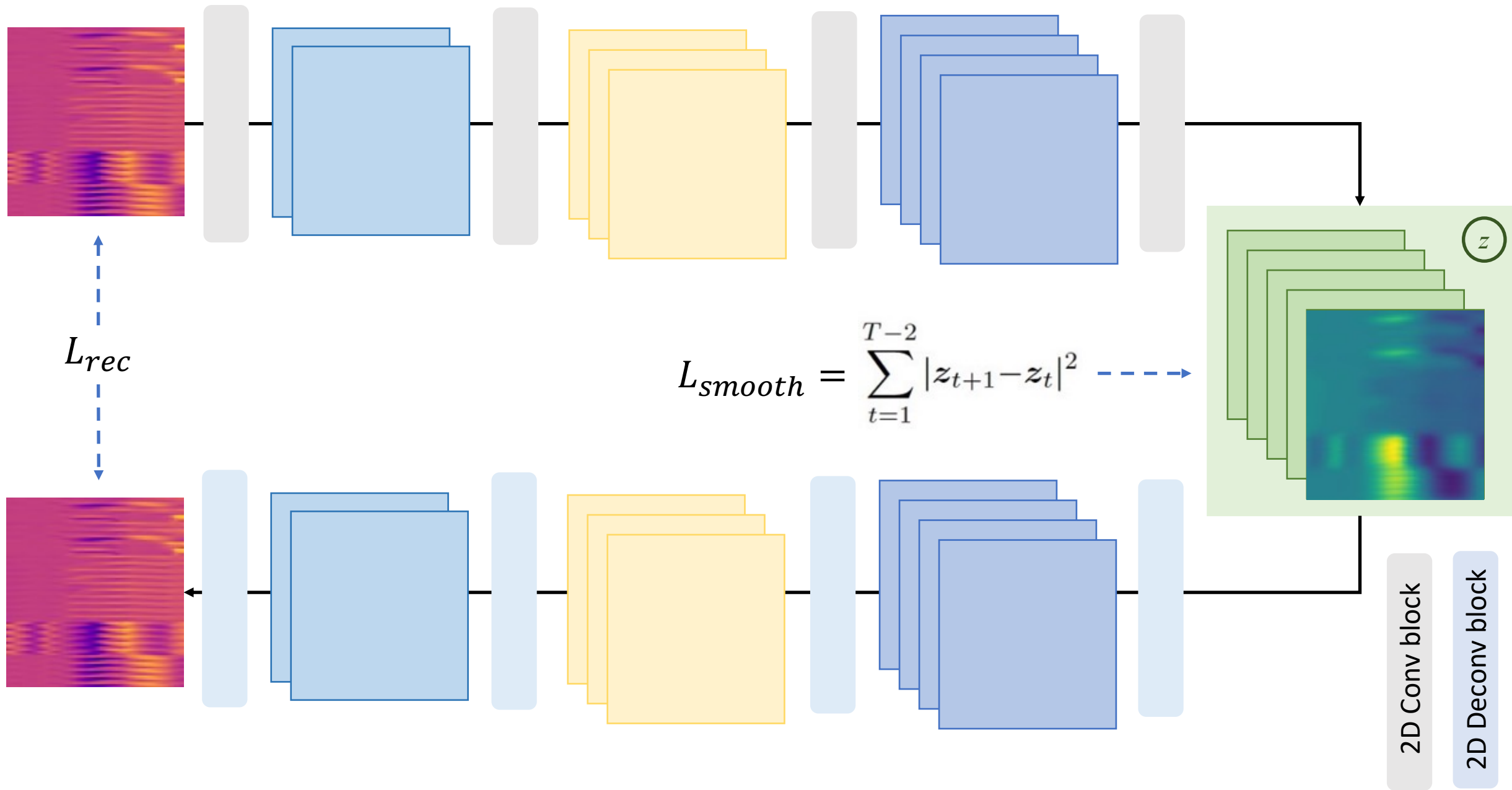


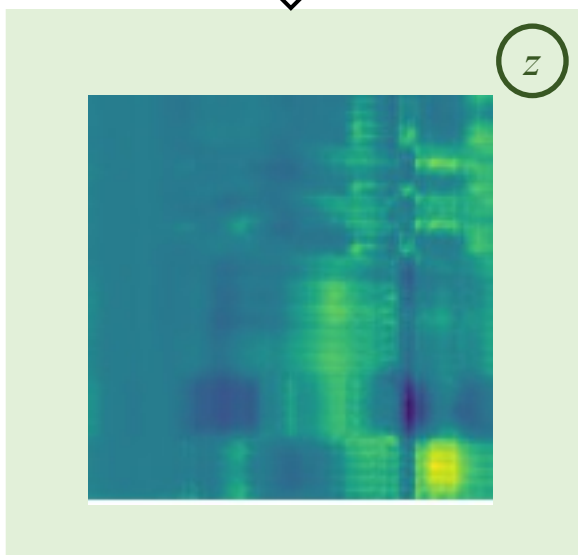
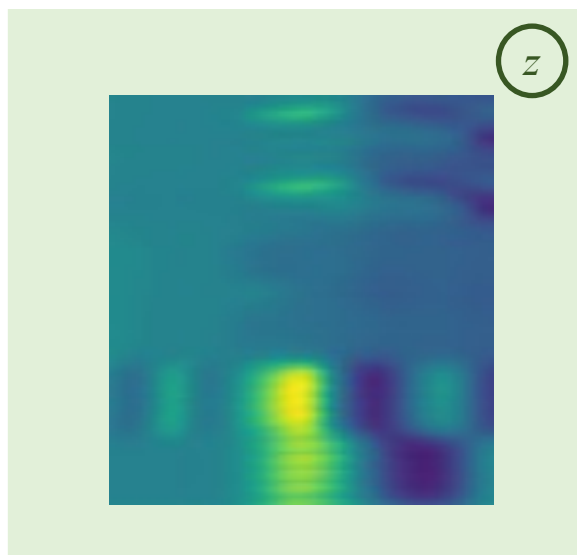
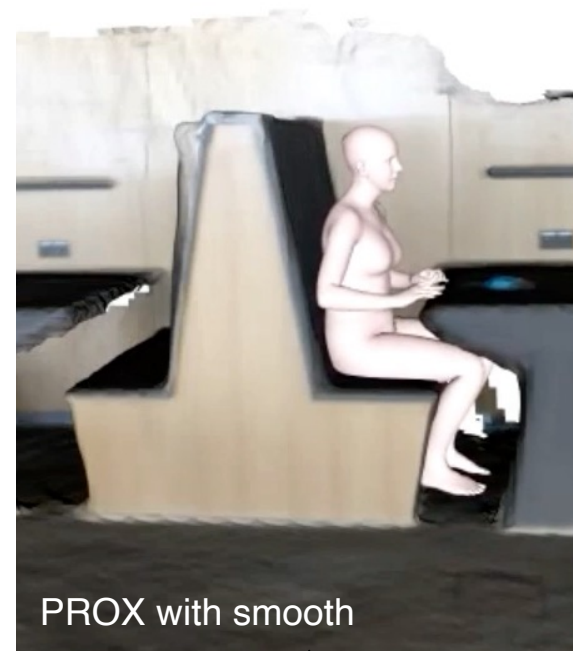
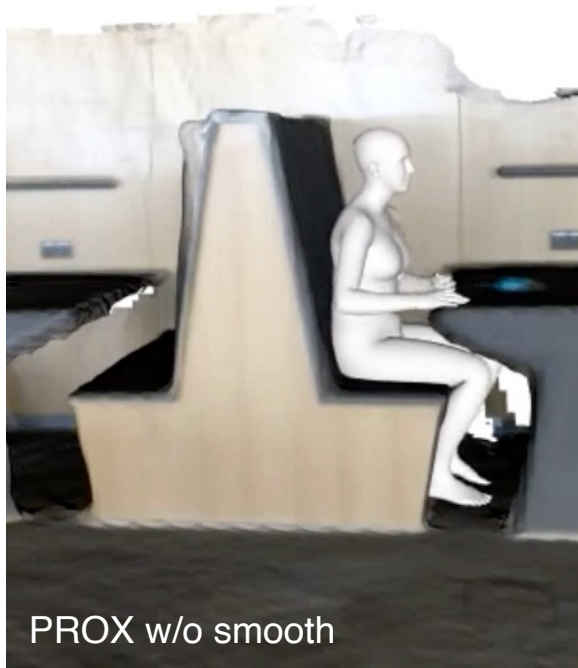




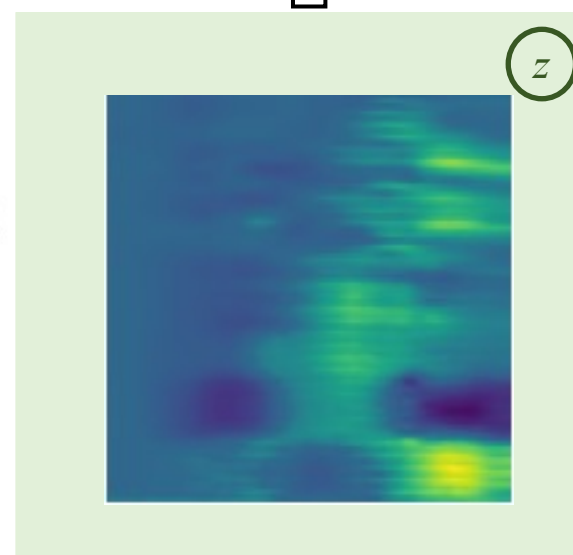
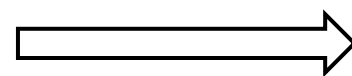


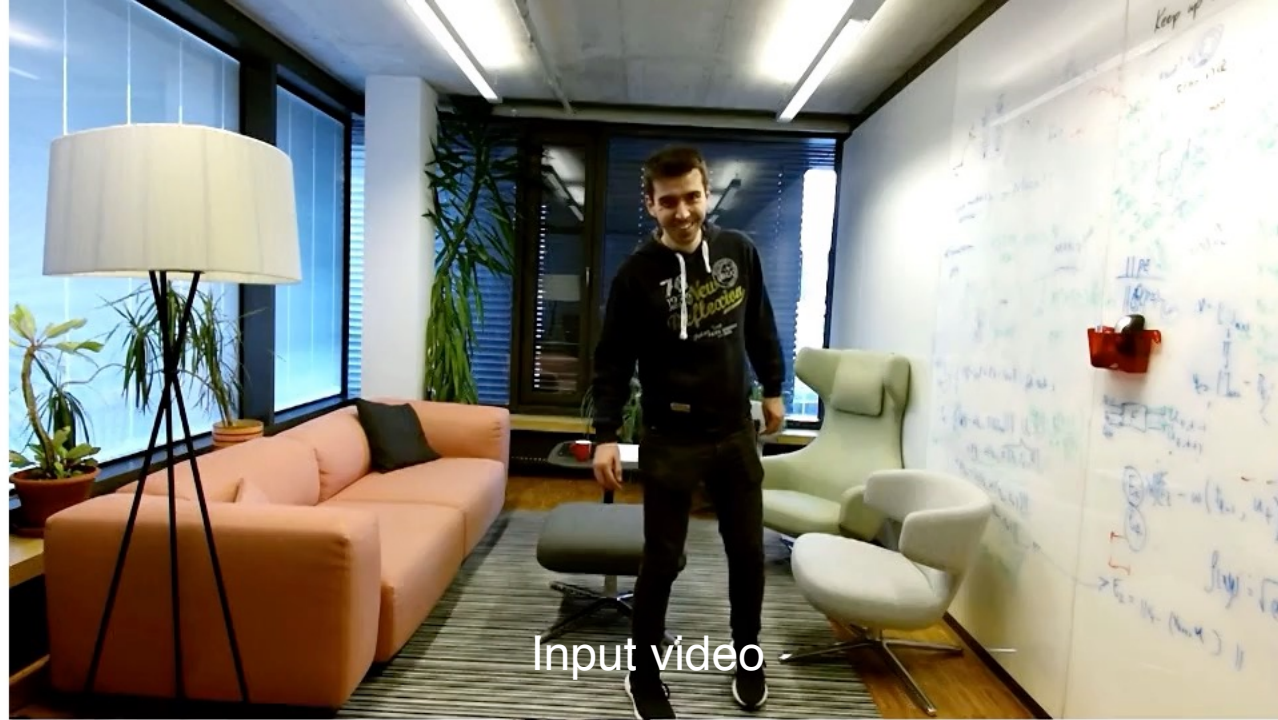






$$\min \sum_{t=1}^{T-2} |z_{t+1} - z_t|^2$$





Input video



PROX (Hassan et al.)



LEMO (Ours)

Outline

- **LEMO: Learning Motion Priors for 4D Human Body Capture in 3D Scenes**

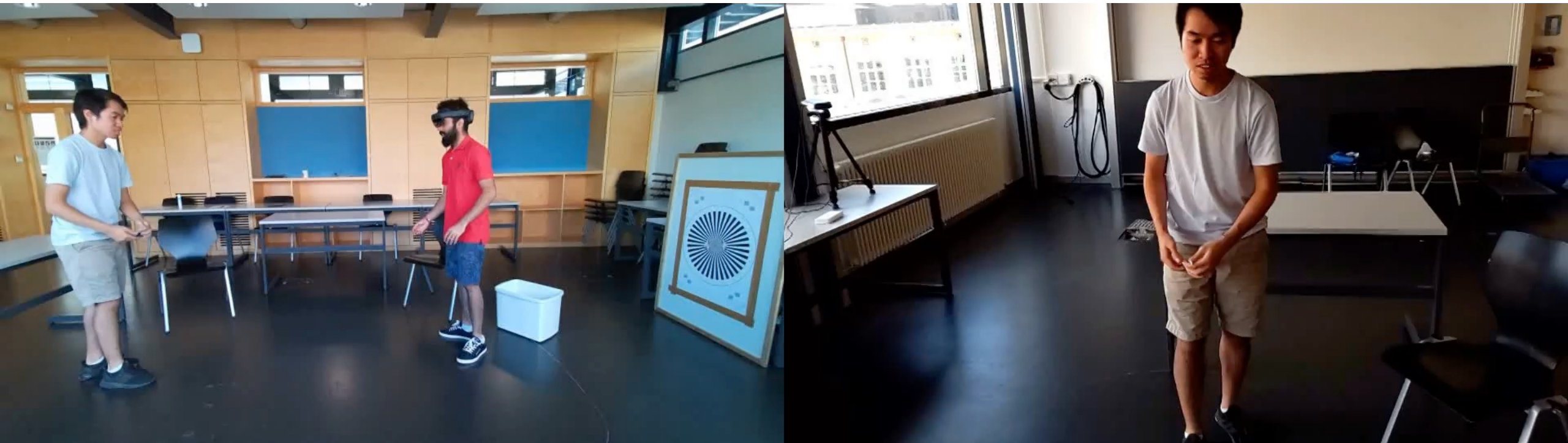
Siwei Zhang, Yan Zhang, Federica Bogo, Marc Pollefeys, Siyu Tang

ICCV 2021, Oral presentation

- **EgoBody: Human Body Shape and Motion of Interacting People from Head-Mounted Devices**

Siwei Zhang, Qianli Ma, Yan Zhang, Zhiyin Qian, Taein Kwon, Marc Pollefeys, Federica Bogo, Siyu Tang

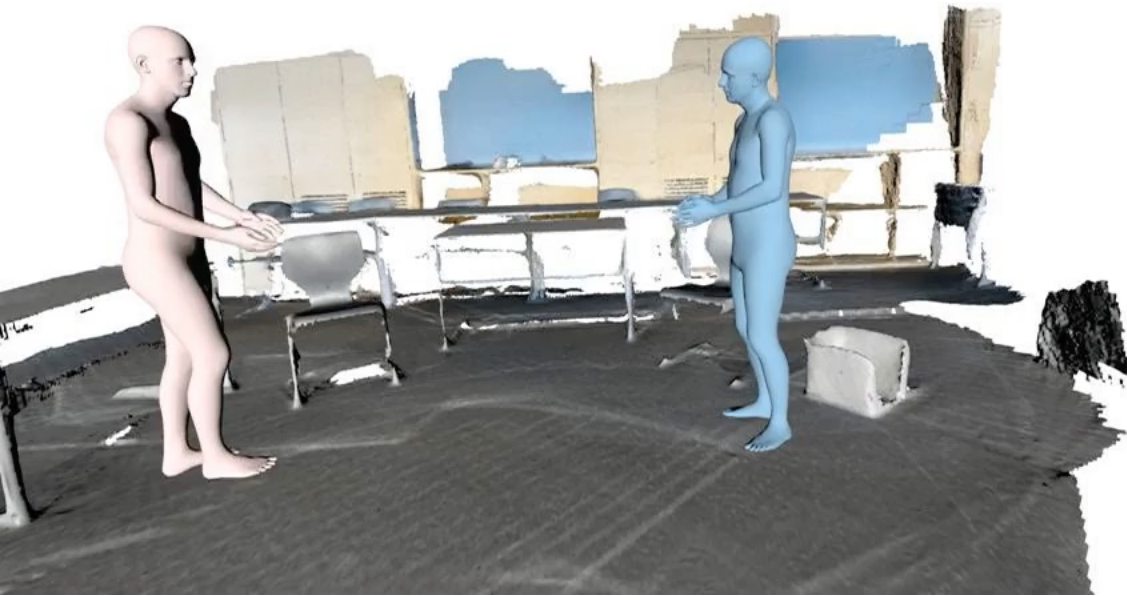
Introducing: EgoBody Dataset



**Third-Person Views,
taken with Kinects**

**Egocentric View,
taken with HoloLens2**

Introducing: EgoBody Dataset



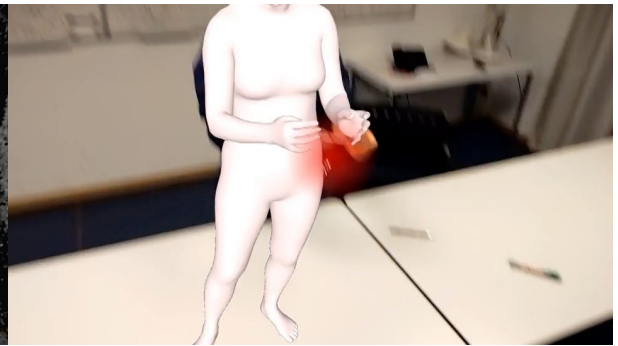
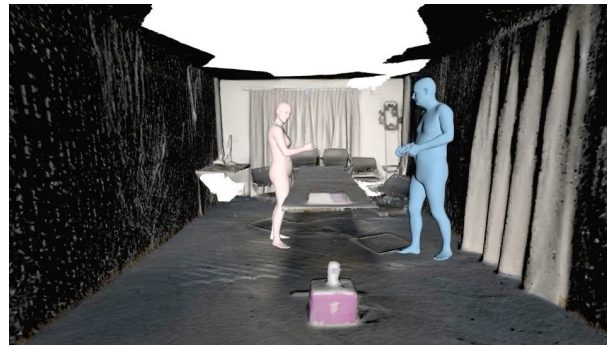
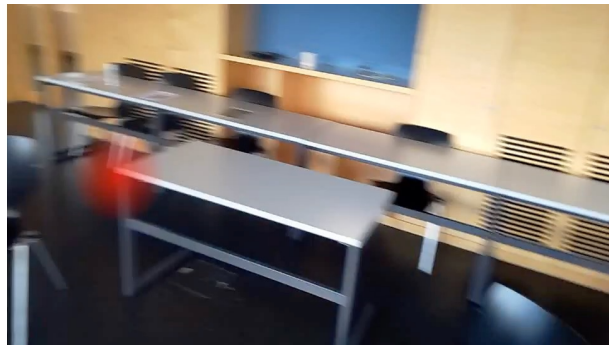
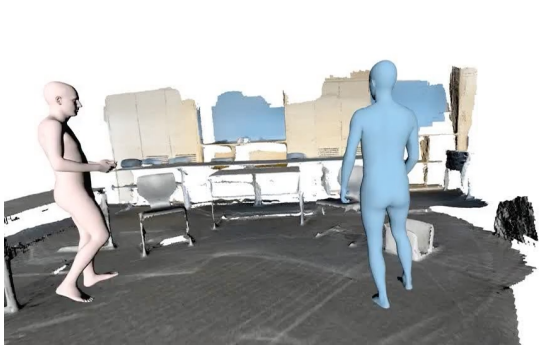
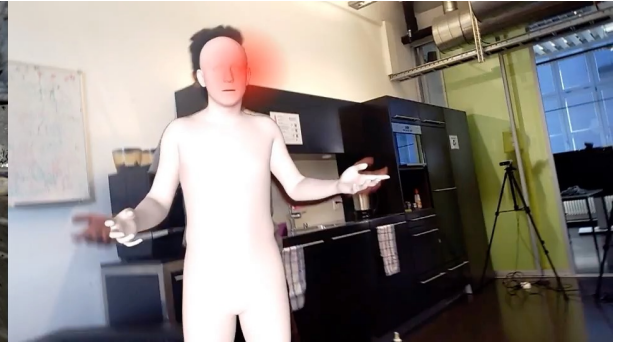
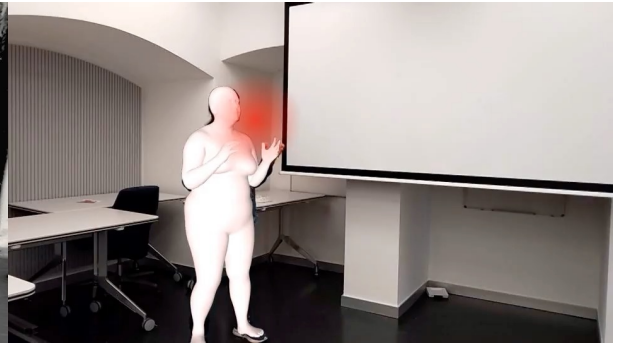
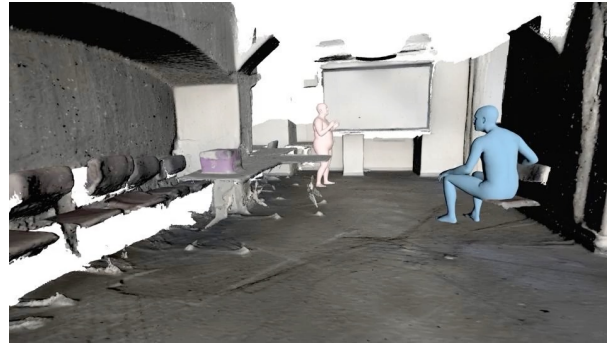
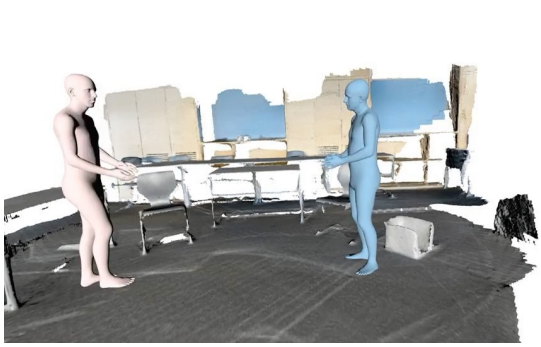
Third-Person View

+ ground truth body annotation of
the camera wearer & the second person
+ scene reconstruction



Egocentric View

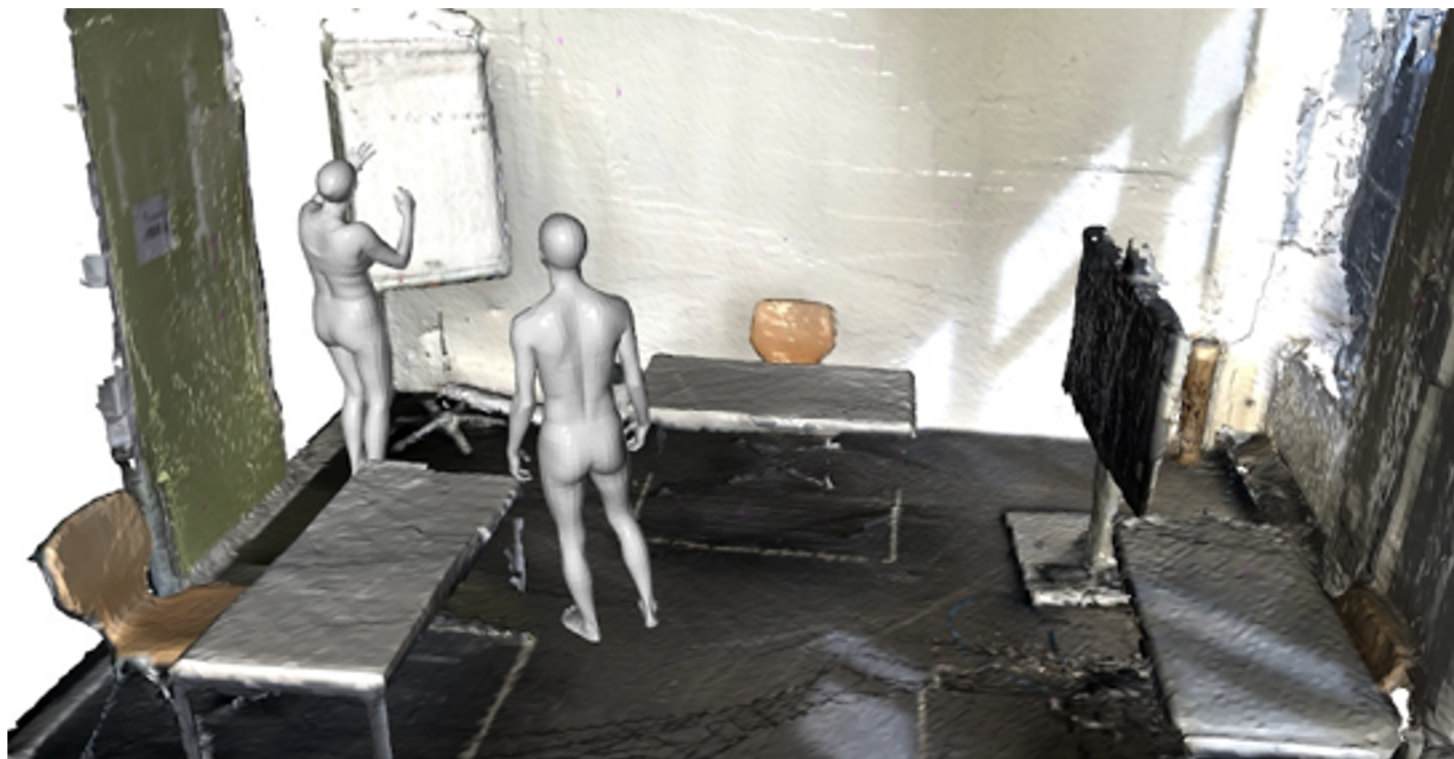
+ ground truth body annotation
+ eye gaze / attention

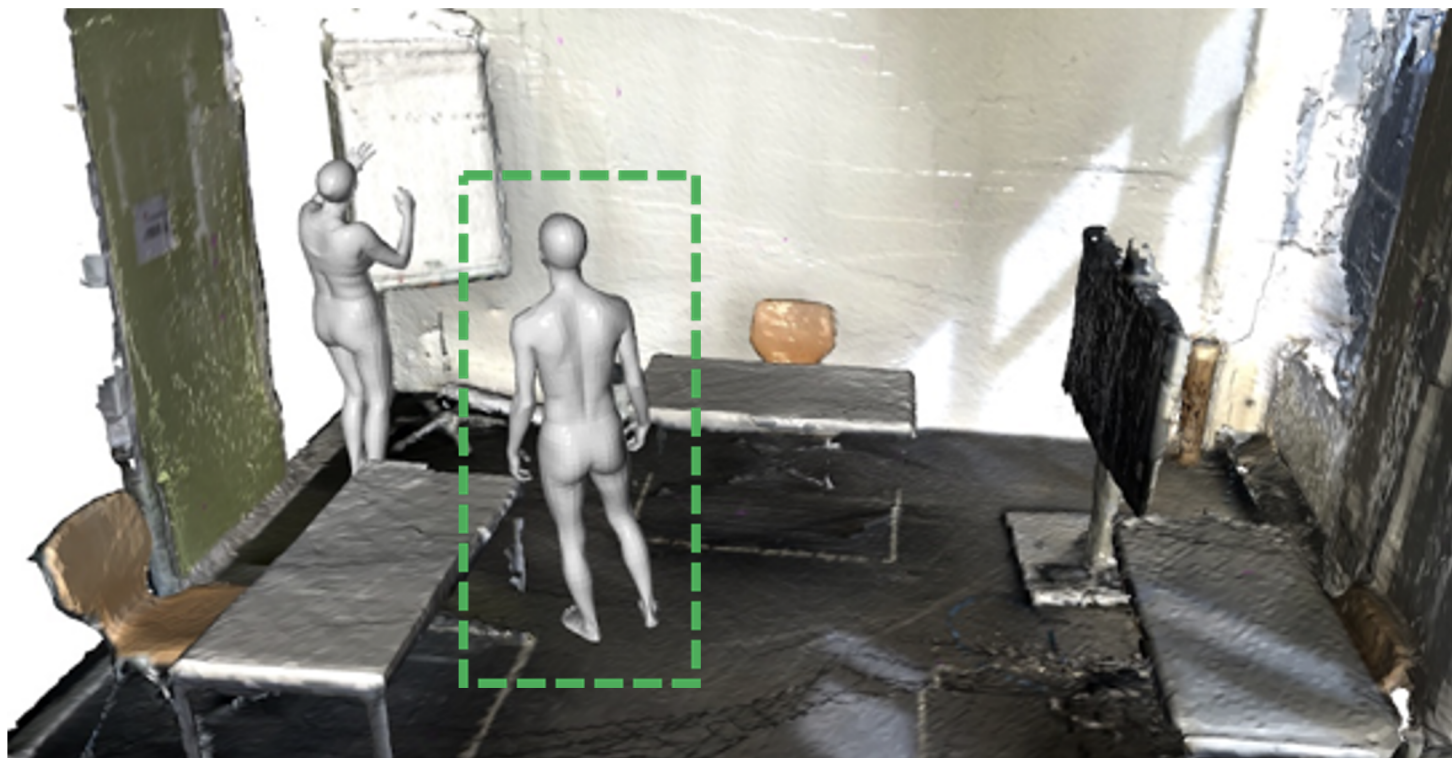


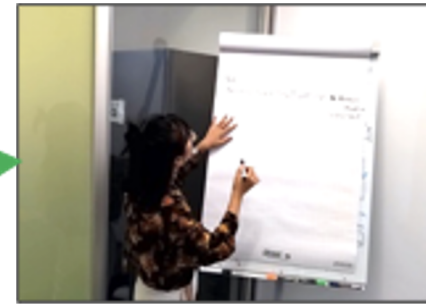
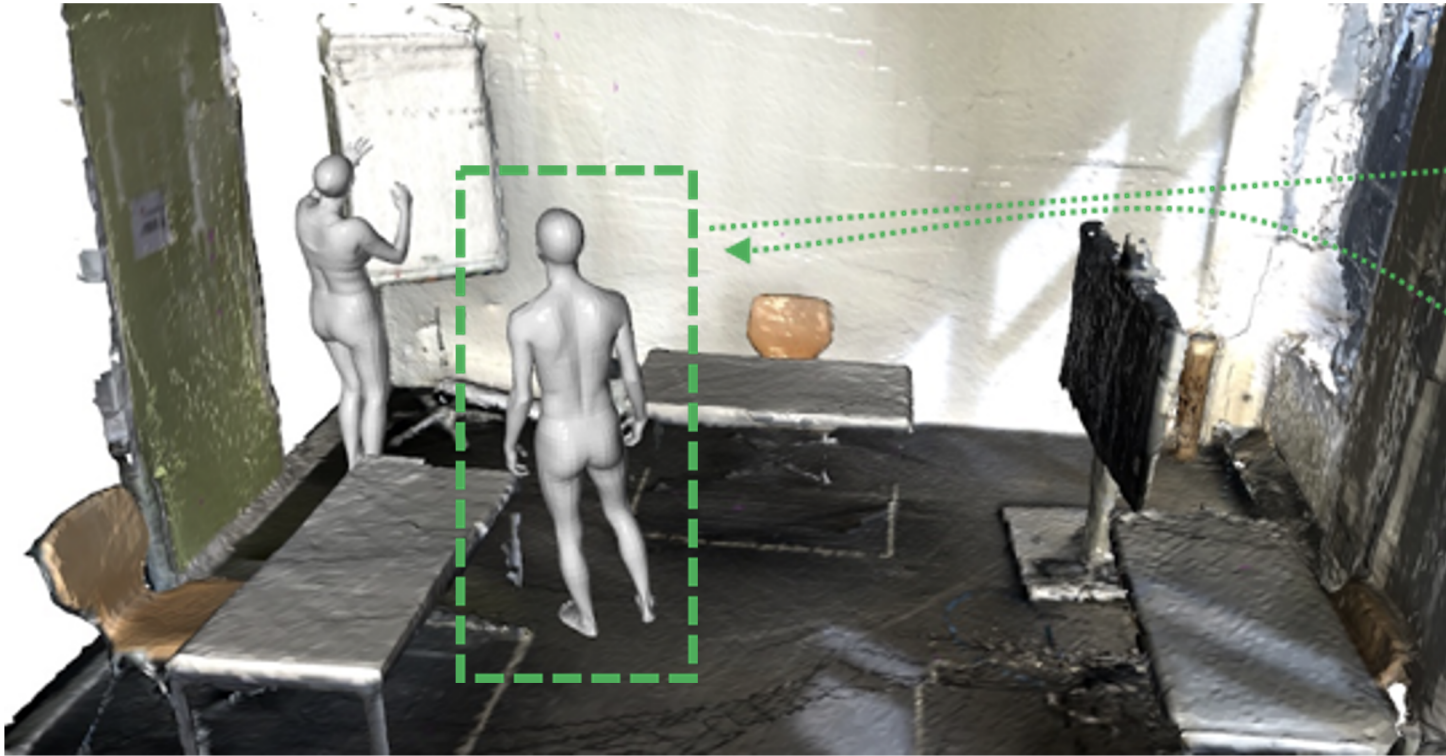
Dataset Overview

- **68 sequences**
- **20 subjects**
- **9 indoor scenes**
- **153k multi-view third-person view RGBD frames** from Azure Kinects
- **139k egocentric view RGB frames** from HoloLens2
- **Eye gaze, hand/head tracking** from HoloLens2
- **3D human shape and motion annotations** for both interacting subjects

Capture Setup



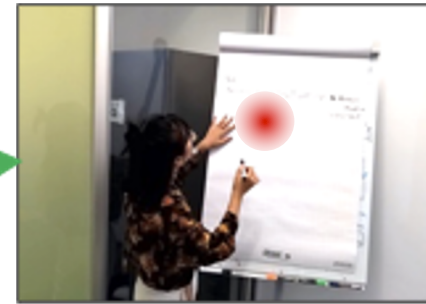
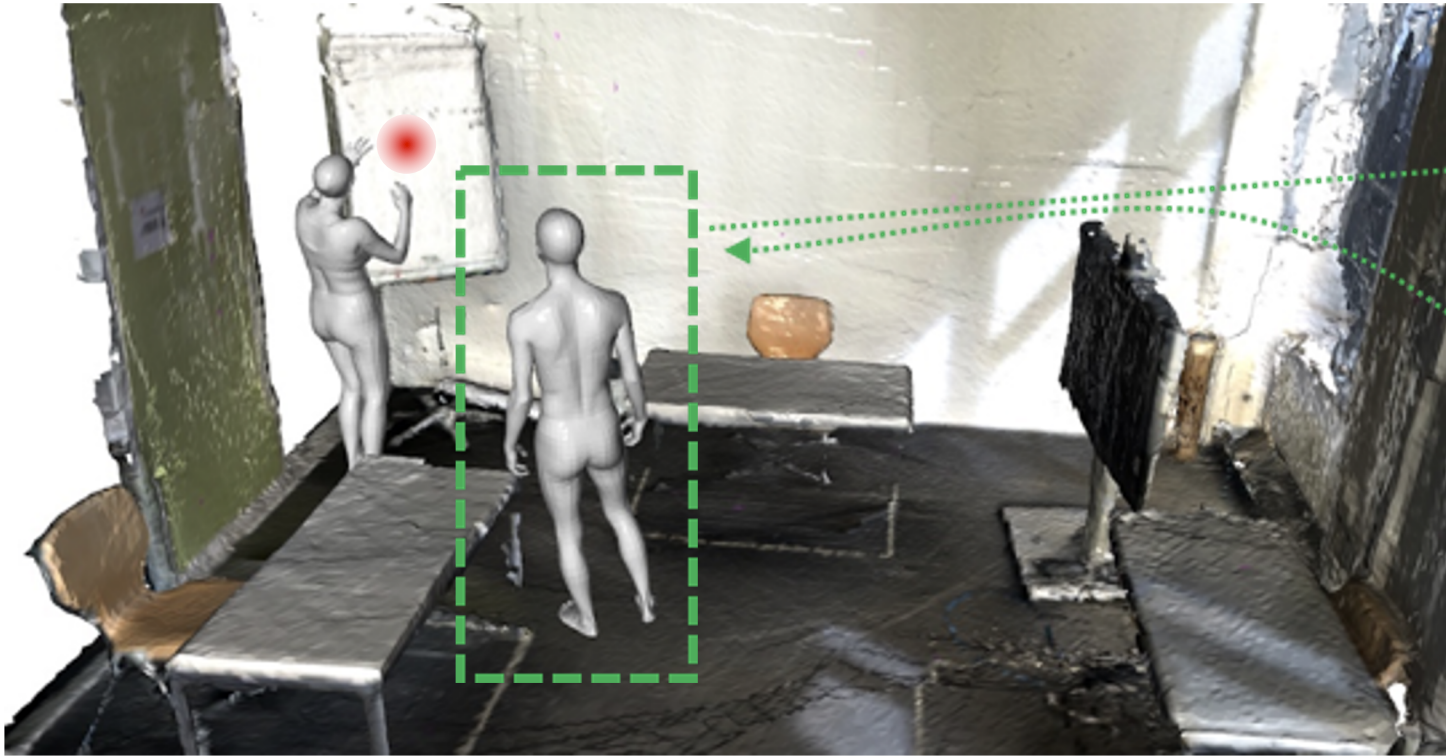




First-Person View



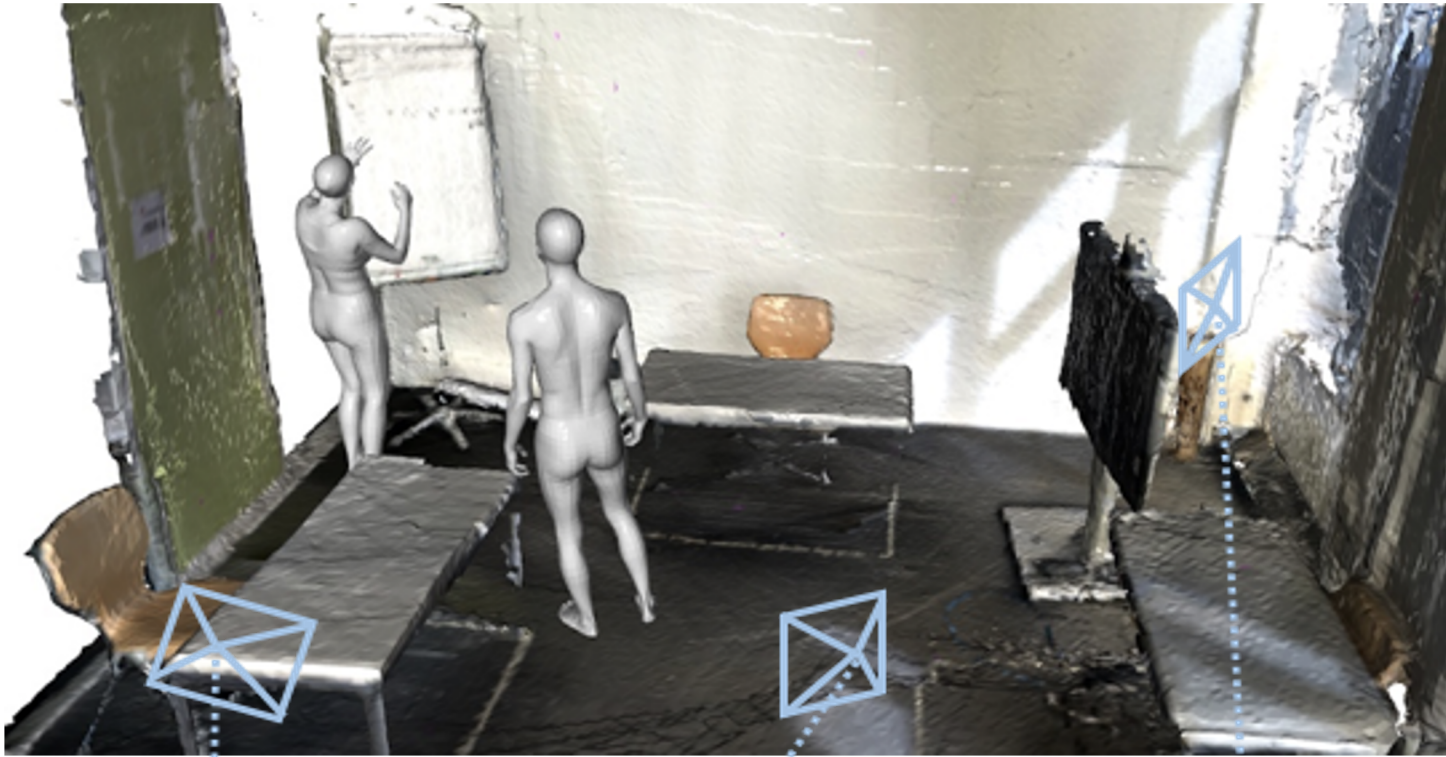
HoloLens2



First-Person View



HoloLens2



Kinect view 1



Kinect view 2

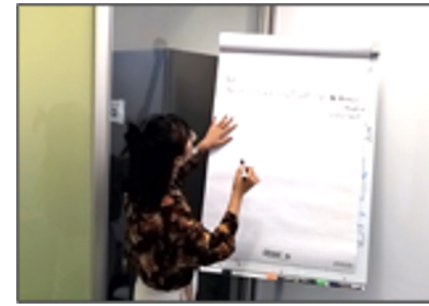
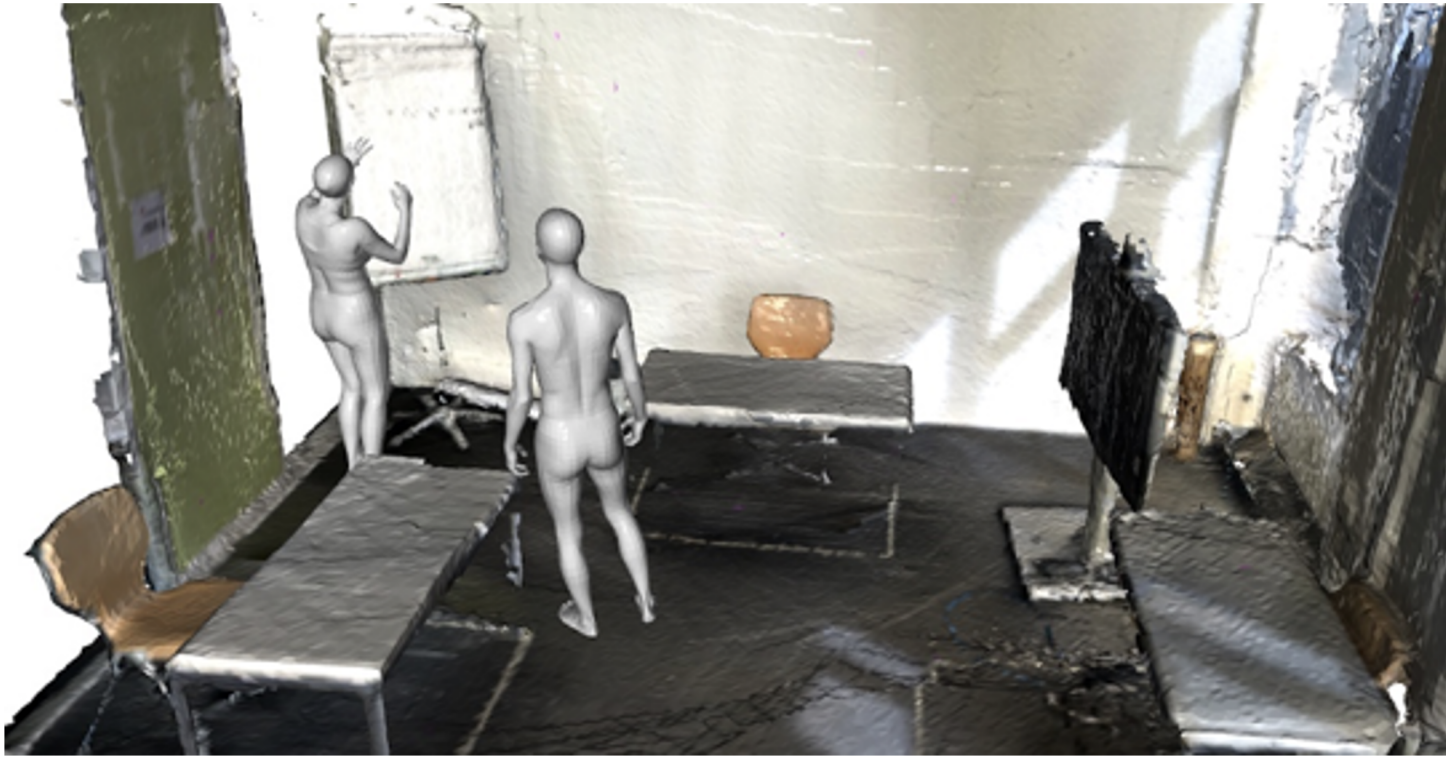


Kinect view 3



Azure Kinect

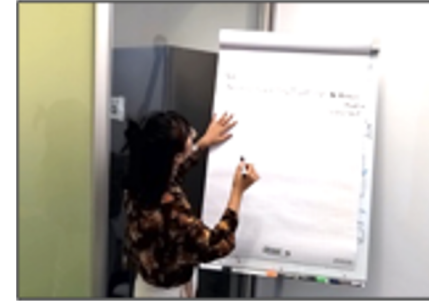
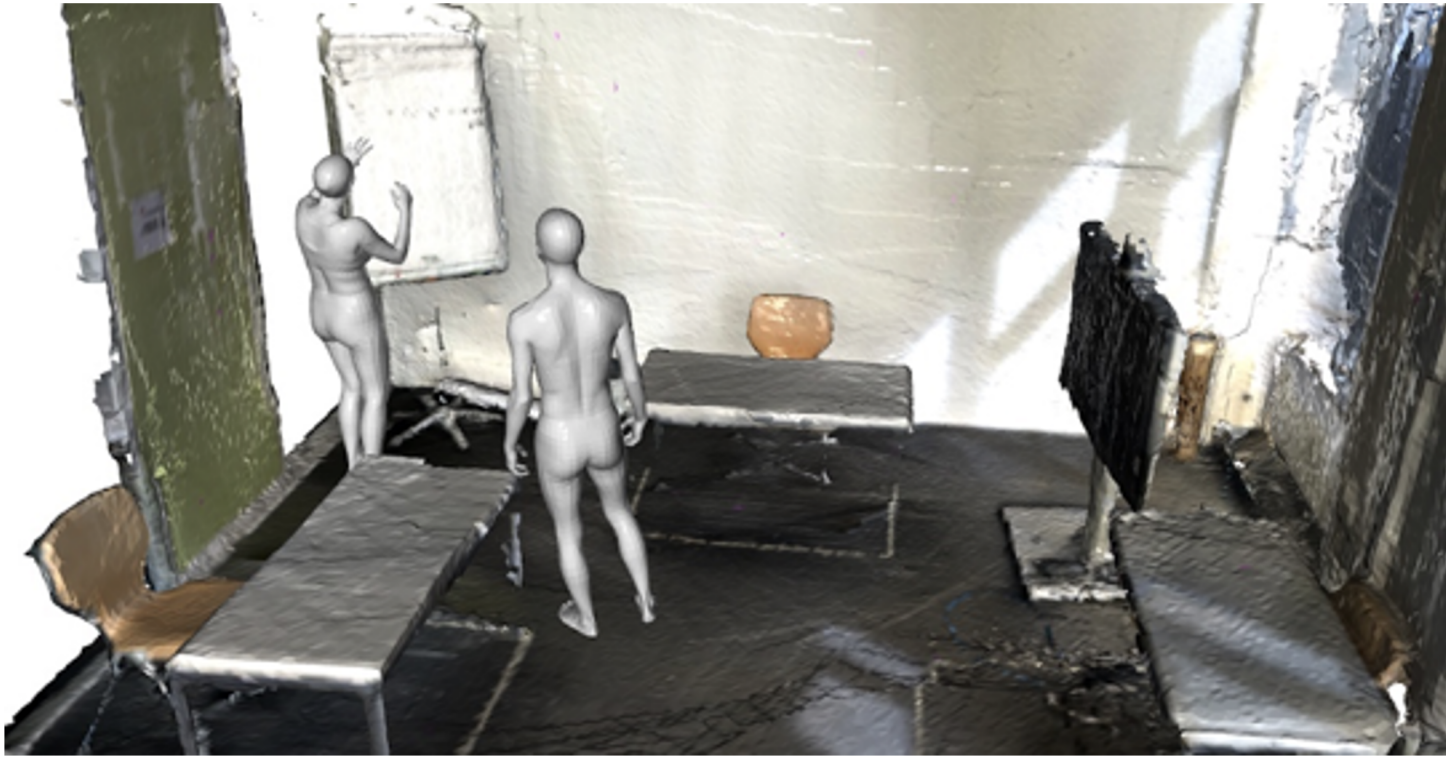
Third-Person View



Kinect view 1

Kinect view 2

Kinect view 3



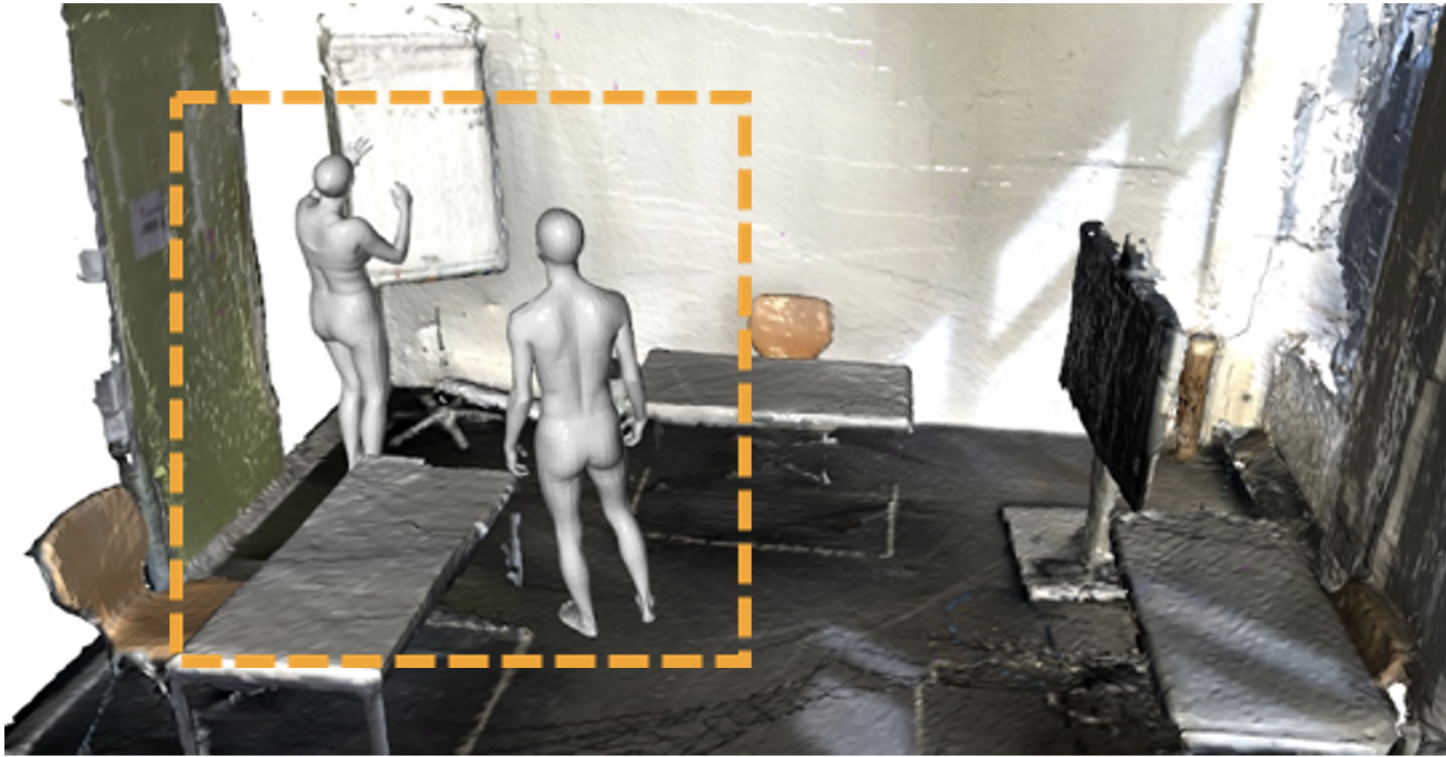
Kinect view 1



Kinect view 2



Kinect view 3



Kinect view 1

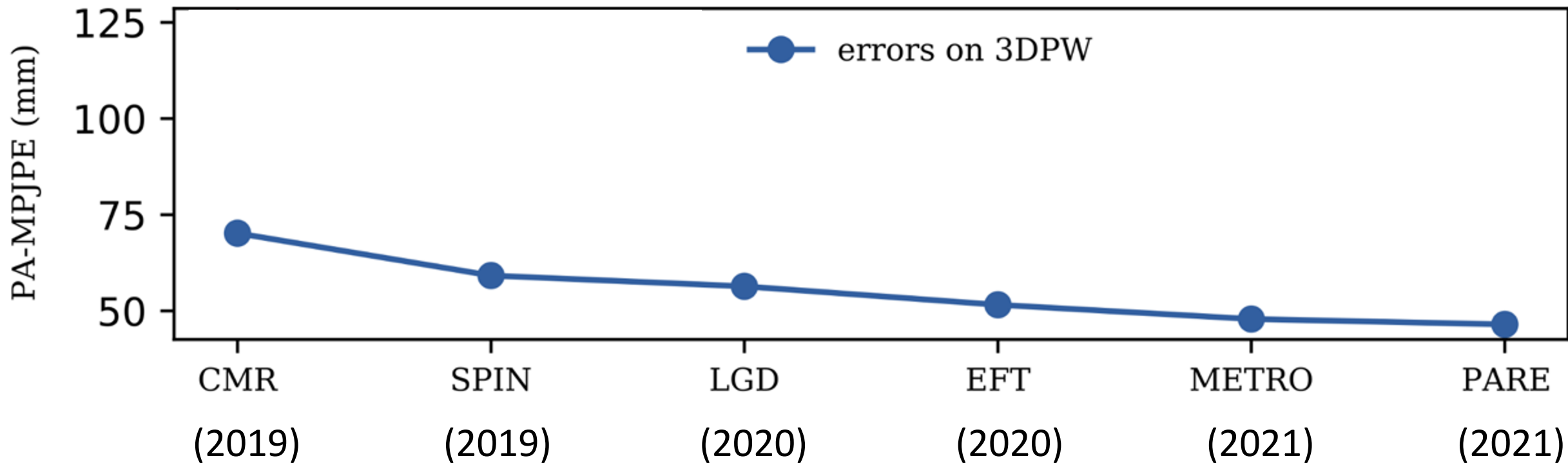


Kinect view 2

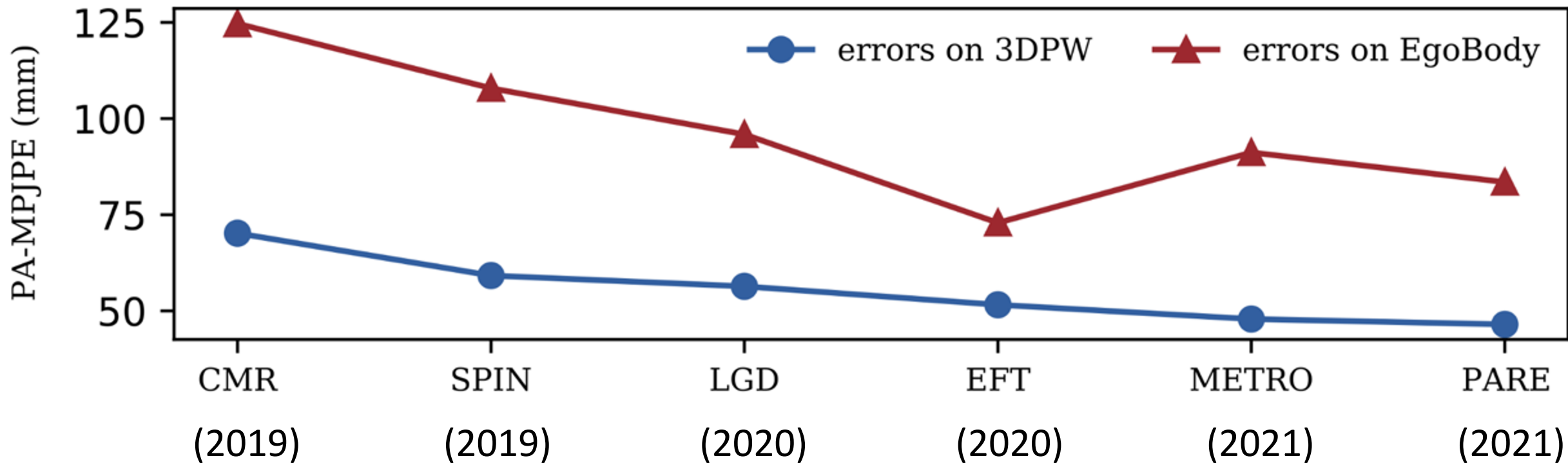


Kinect view 3

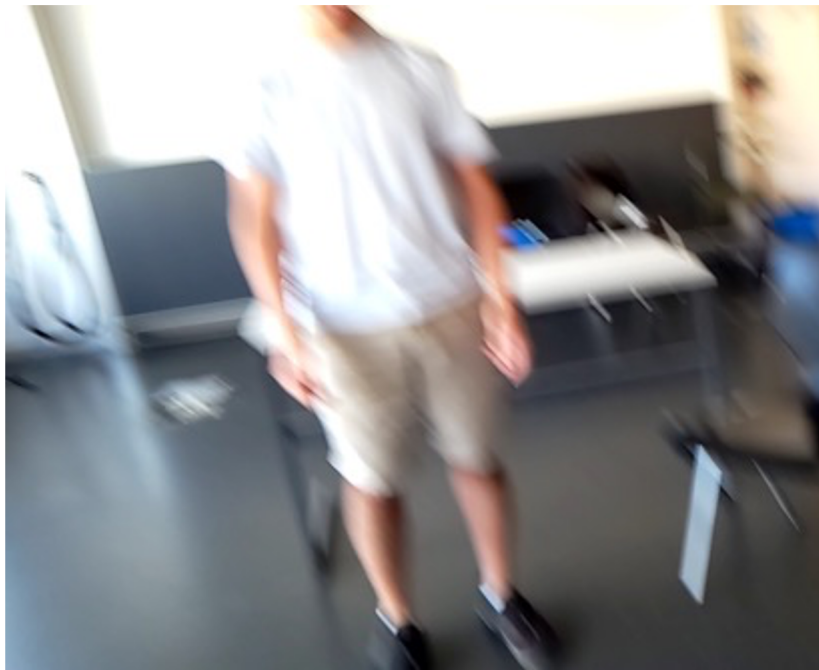
Benchmark:
3D Human Pose and Shape Estimation
From Egocentric Images



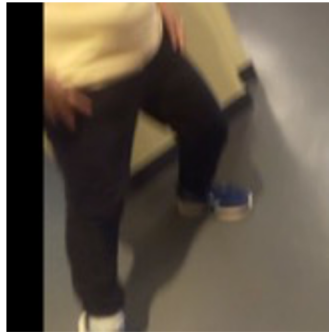
Advance of 3D Human Pose and Shape Estimation Methods



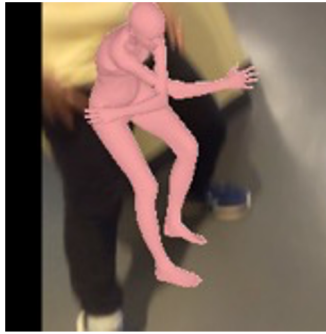
Advance of 3D Human Pose and Shape Estimation Methods



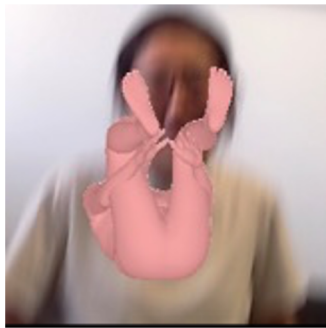
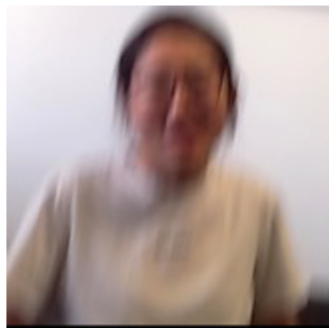
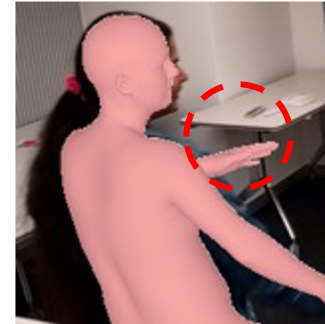
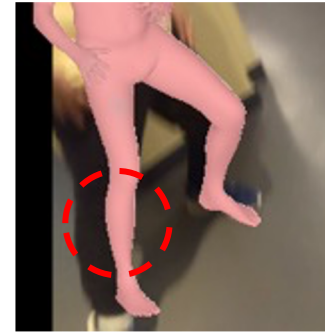
Input Image



SPIN
(Kolotouros et al.)



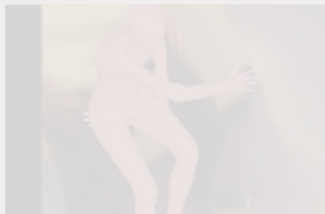
EFT
(Joo et al.)



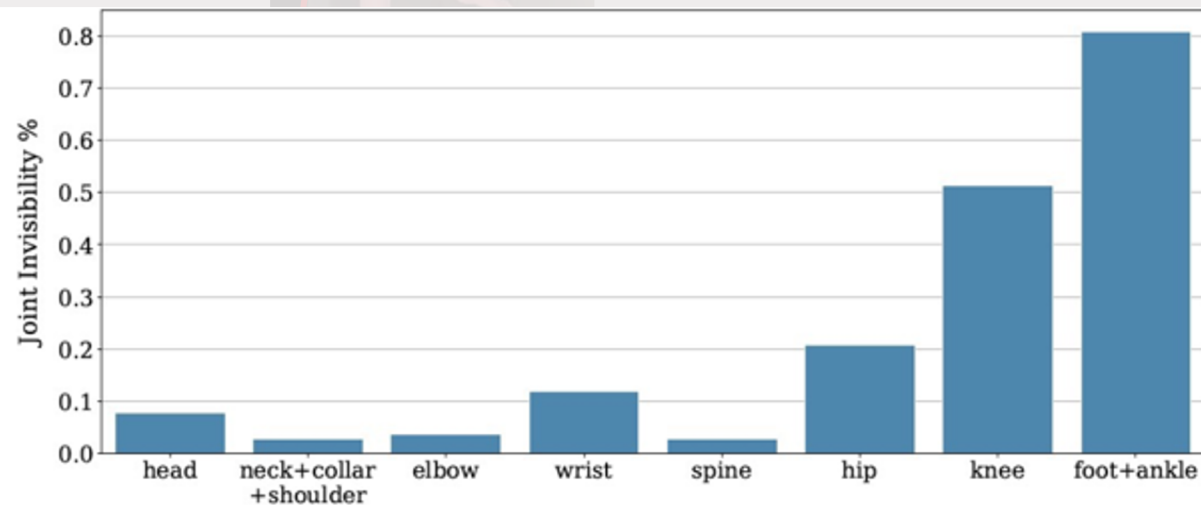
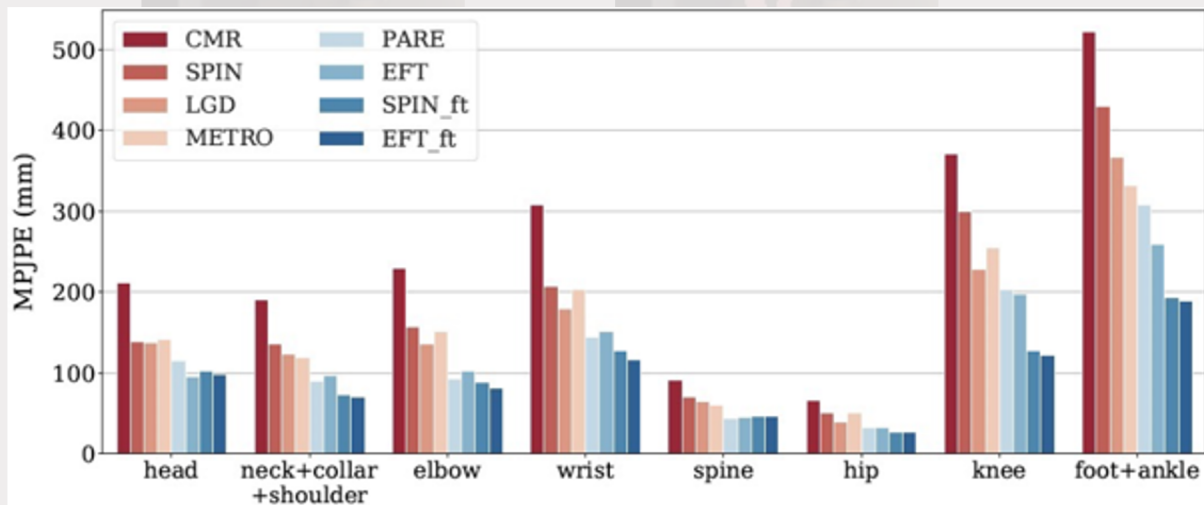
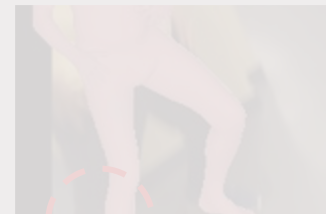
Input Image



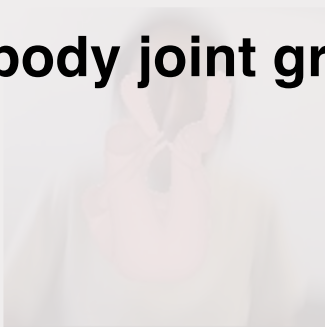
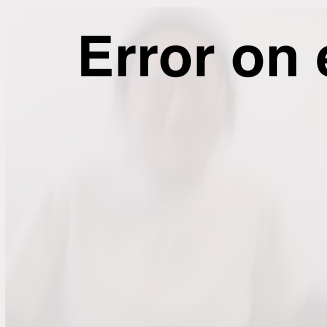
SPIN
(Kolotouros et al.)



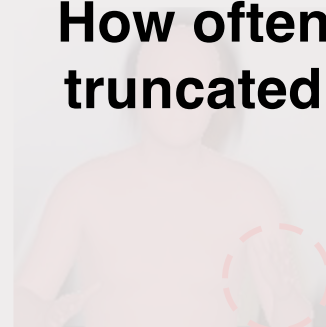
EFT
(Joo et al.)



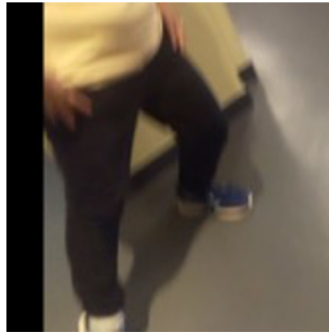
Error on each body joint group



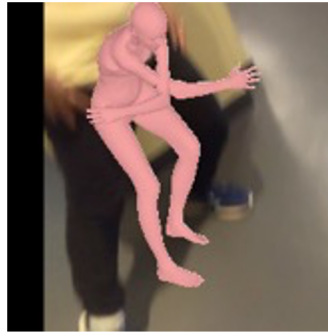
How often a joint group is truncated from the image



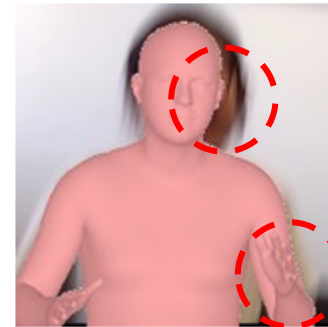
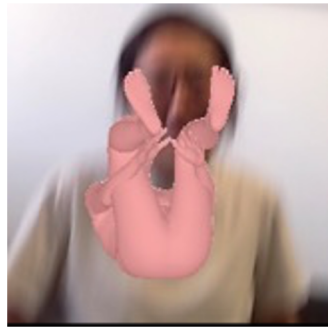
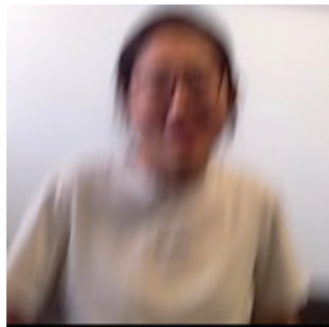
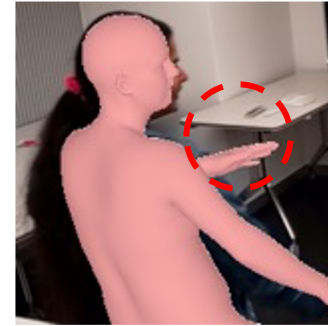
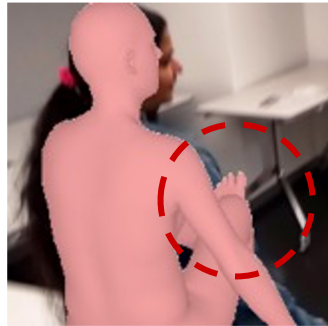
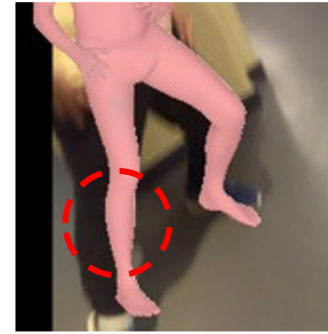
Input Image



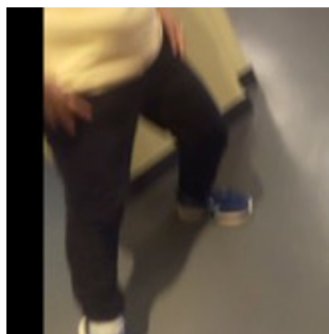
SPIN
(Kolotouros et al.)



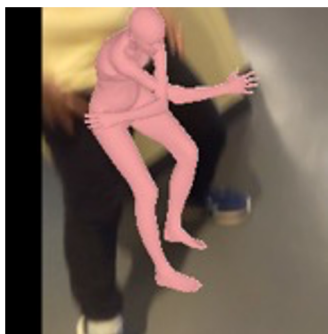
EFT
(Joo et al.)



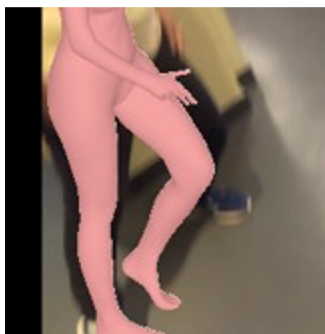
Input Image



SPIN
(Kolotouros et al.)



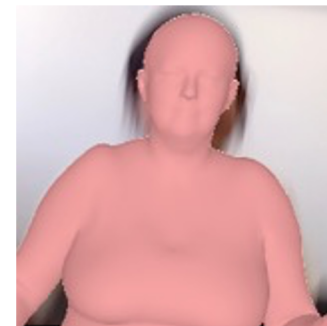
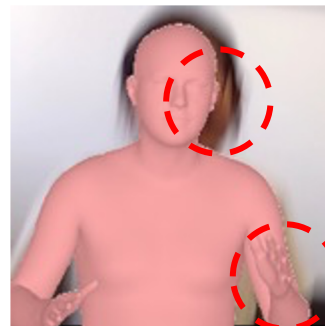
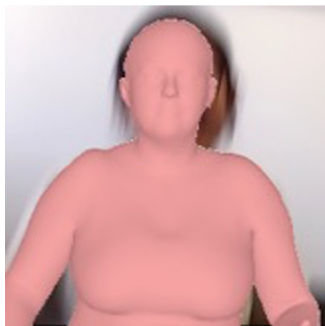
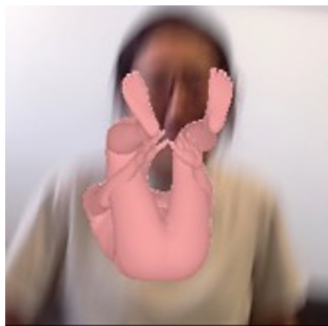
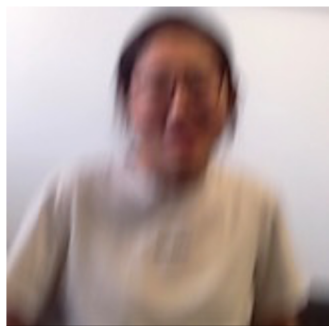
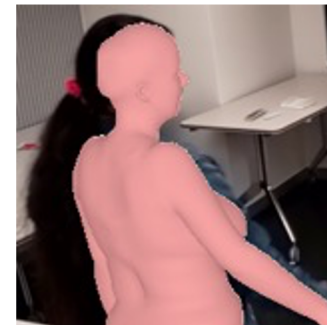
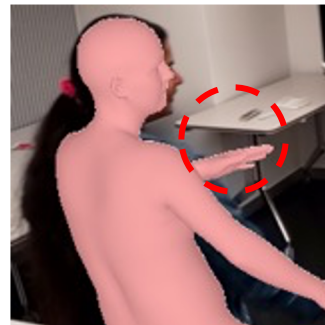
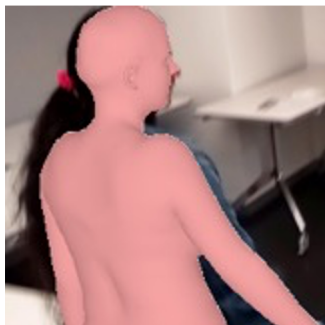
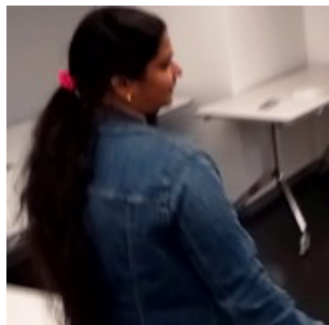
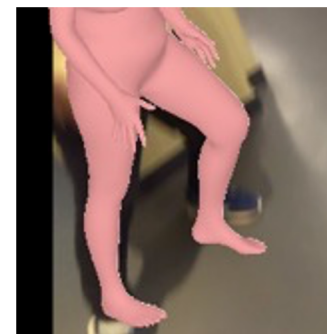
SPIN
fine-tuned on our
training set



EFT
(Joo et al.)



EFT
fine-tuned on our
training set



Pose and Shape Estimation Errors on **EgoBody** Test Set

Method	MPJPE ↓	PA-MPJPE ↓	V2V ↓	PA-V2V ↓
SPIN	189.9		210.5	
SPIN-ft (Ours)	96.2	-49%	122.3	-42%
METRO	161.5		187.5	
METRO-ft (Ours)	105.4	-35%	105.9	-44%
EFT	123.3		143.3	
EFT-ft (Ours)	92.5	-25%	115.1	-20%

Pose and Shape Estimation Errors on **EgoBody** Test Set




Method	MPJPE ↓	PA-MPJPE ↓	V2V ↓	PA-V2V ↓
SPIN	189.9		210.5	
SPIN-ft (Ours)	96.2	-49%	122.3	-42%
METRO	161.5		187.5	
METRO-ft (Ours)	105.4	-35%	105.9	-44%
EFT	123.3		143.3	
EFT-ft (Ours)	92.5	-25%	115.1	-20%

Pose accuracy

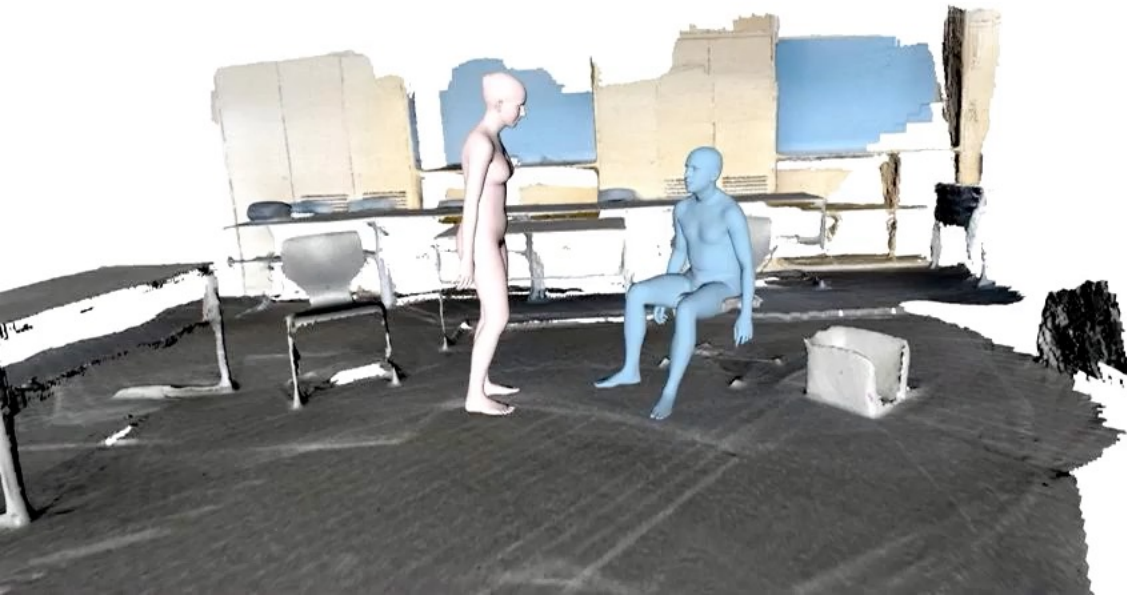
Shape accuracy

Cross-dataset Evaluation on You2Me

Finetuned on **EgoBody** training set
Pose and Shape Estimation Errors on the **You2Me** Dataset

Method	PA-MPJPE ↓		
SPIN	155.0		-42%
SPIN-ft (Ours)	89.4		
METRO	117.7		-23%
METRO (Ours)	90.1		
EFT	96.0		-8%
EFT-ft (Ours)	88.7		

More results: EgoBody Dataset



Third-Person View

+ ground truth body annotation of
the camera wearer & the second person
+ scene reconstruction



Egocentric View

+ ground truth body annotation
+ eye gaze / attention

More results: LEMO



Egocentric Interaction Capture for Mixed Reality



Project page (LEMO):

<https://sanweiliti.github.io/LEMO/LEMO.html>



Project page (EgoBody):

<https://sanweiliti.github.io/egobody/egobody.html>

ETH zürich



VLG

Computer Vision
and Learning
Group



Microsoft

siwei.zhang@inf.ethz.ch

Contact-aware Motion Infilling Prior

