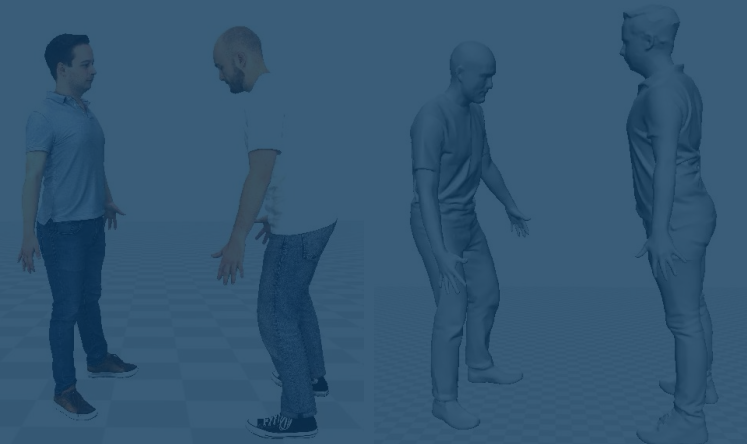
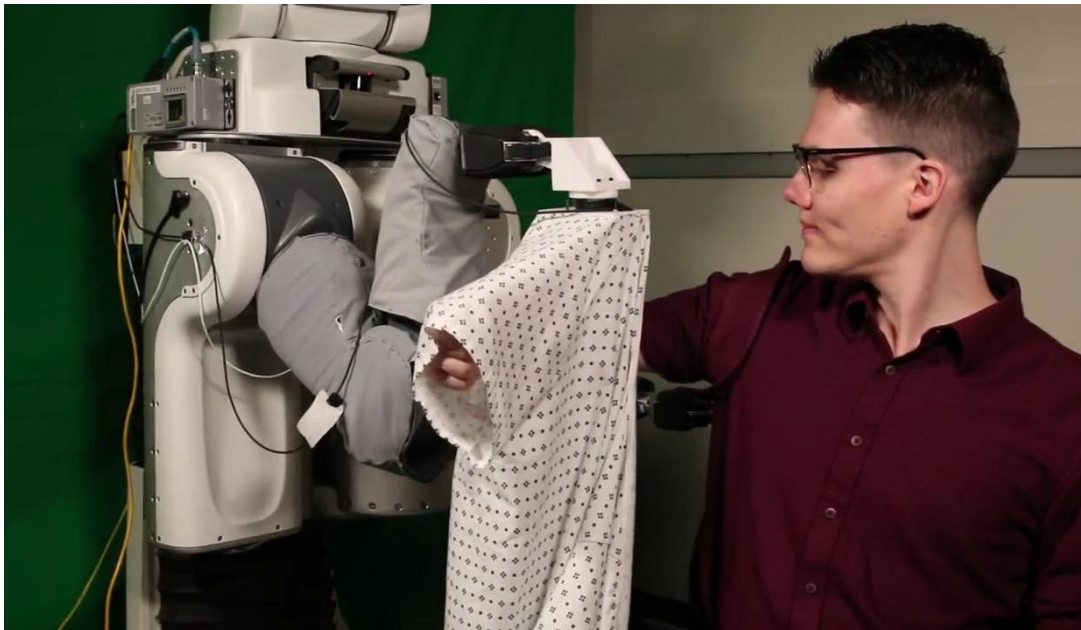
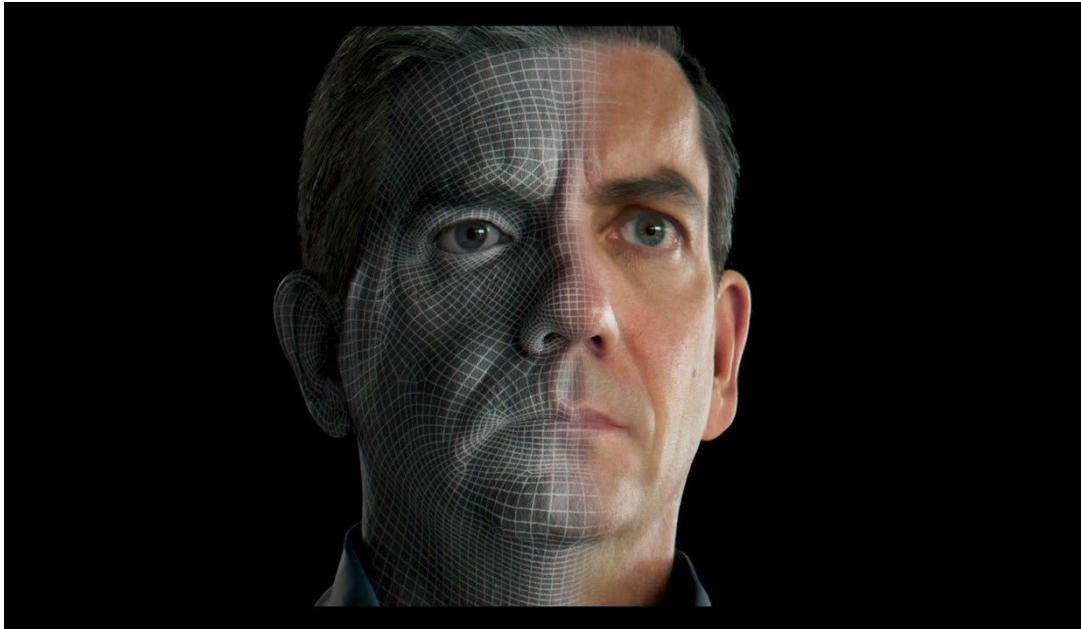


# Digitalization of Clothed Humans with Expressive Behaviors from Videos

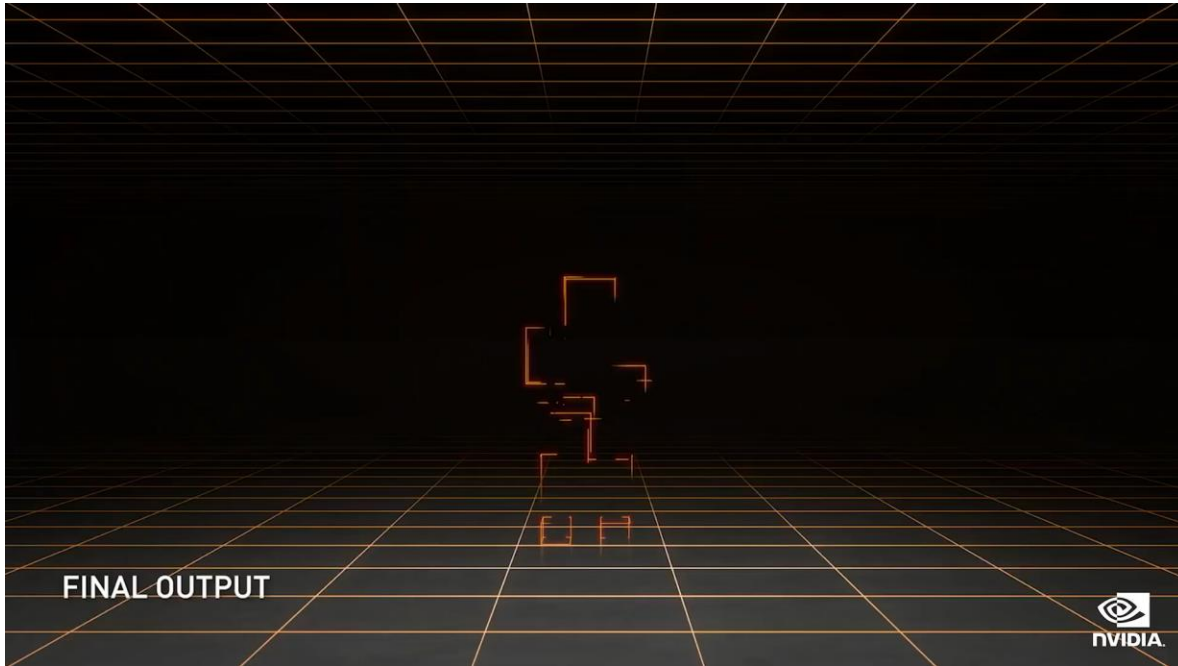
Chen Guo

Wed., 10.04.2024





Goal: Automatic reconstruction of detailed and expressive human model in clothing from videos



Immersive Telepresence



- High requirement of devices
- Tons of manual efforts



# Vid2Avatar: 3D Avatar Reconstruction from Video in the Wild via Self-supervised Scene Decomposition

Chen Guo Tianjian Jiang Xu Chen Jie Song Otmar Hilliges

CVPR 2023

**ETH** zürich



# Vid2Avatar: 3D Avatar Reconstruction from Video in the Wild via Self-supervised Scene Decomposition

Chen Guo Tianjian Jiang Xu Chen Jie Song Otmar Hilliges

CVPR 2023

**ETH** zürich



Advanced Interactive  
Technologies



# Vid2Avatar: 3D Avatar Reconstruction from Video in the Wild via Self-supervised Scene Decomposition

Chen Guo Tianjian Jiang Xu Chen Jie Song Otmar Hilliges

CVPR 2023

**ETH** zürich



# Problem Definition

# Problem Definition



Input: RGB sequence



# Problem Definition



Input: RGB sequence



Reconstruction

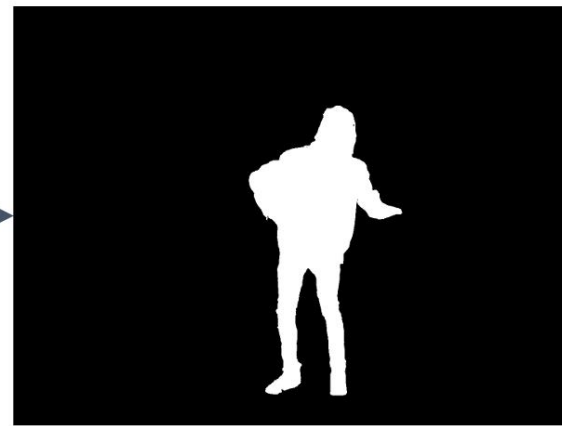
- Accurate separation
- Detailed 3D surfaces

# Motivation



Image

Off-the-shelf  
Segmentation Tool



Mask[1]

Reconstruction



SelfRecon[2]

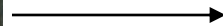
[1] RVM: Lin et al. '21

[2] SelfRecon: Jiang et al. '22

# Method Overview



Input: RGB sequence



Segmentation



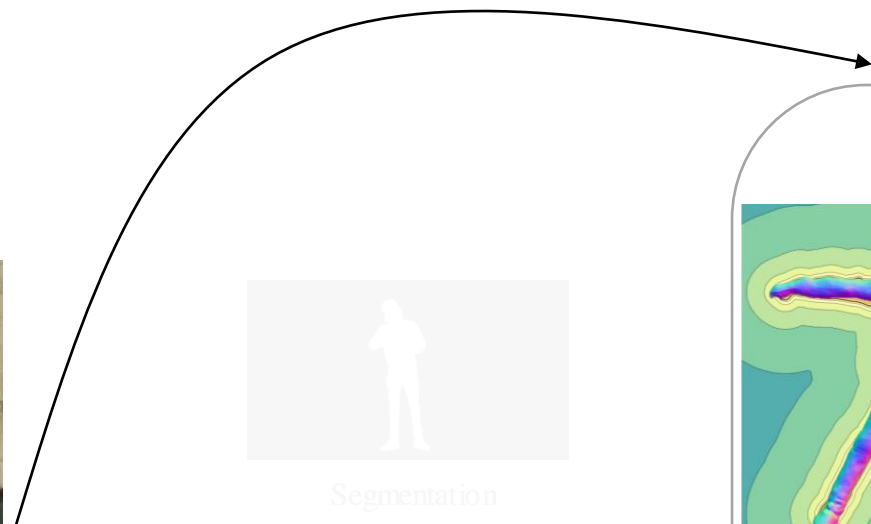
Surface Reconstruction



# Method Overview



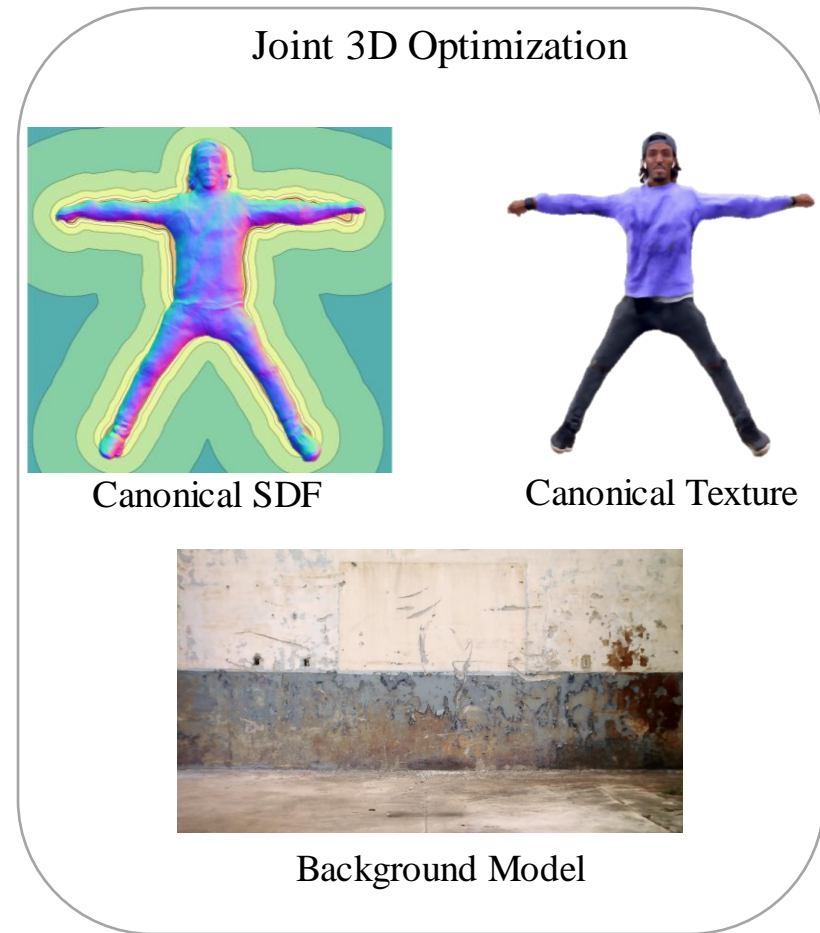
Input: RGB sequence



Segmentation



Surface Reconstruction



Joint 3D Optimization



Canonical SDF



Canonical Texture



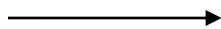
Background Model



# Method Overview



Input: RGB sequence

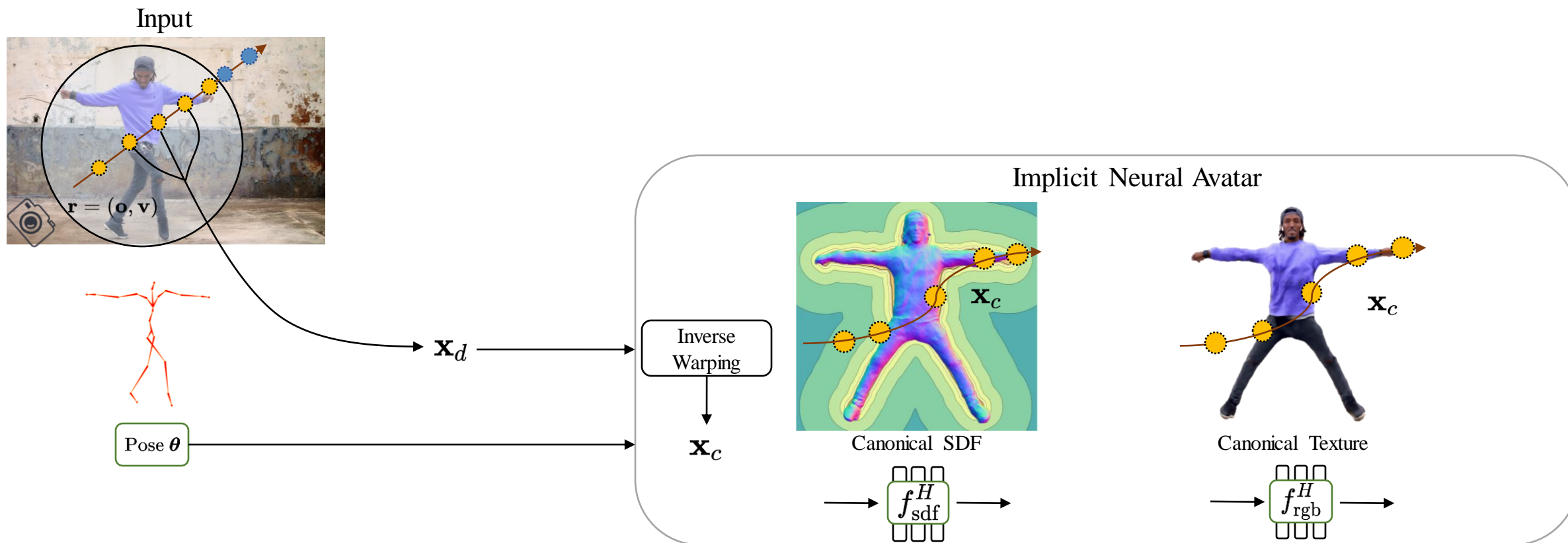


3D Avatar

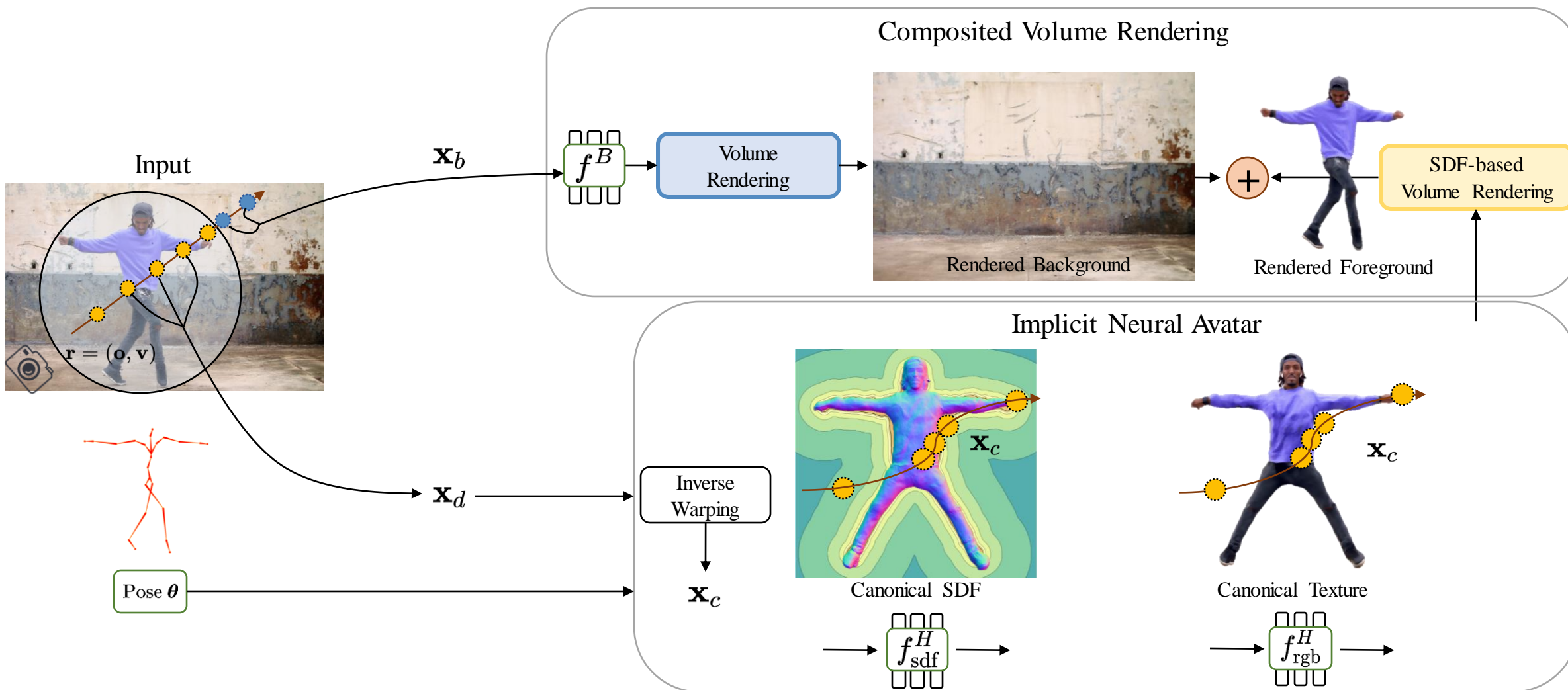


Background

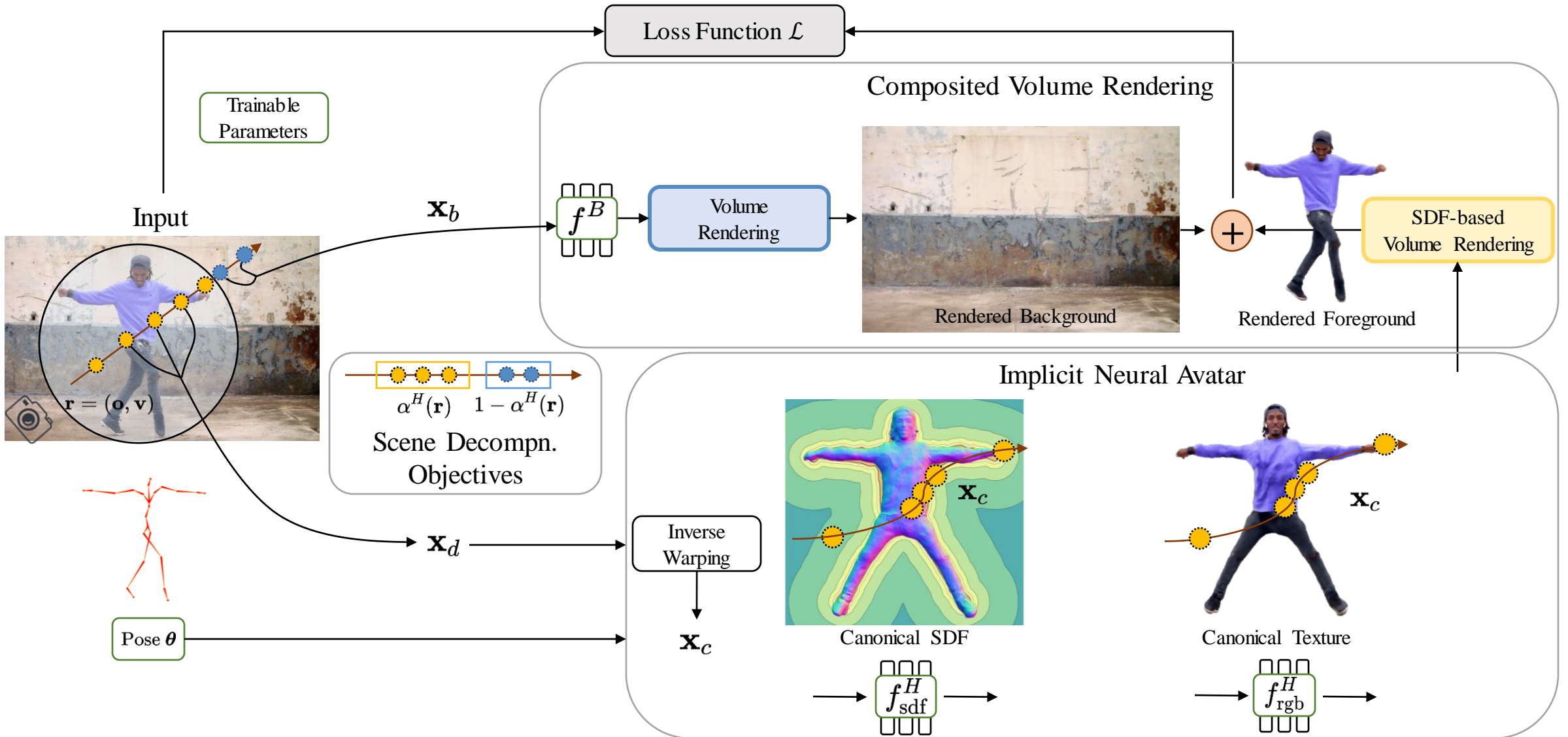
# Method



# Method



# Method

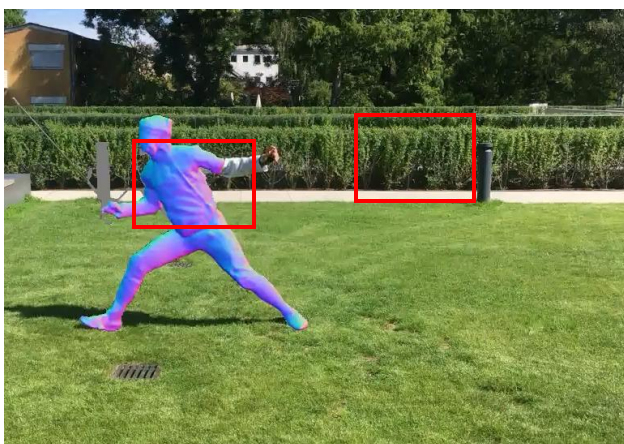




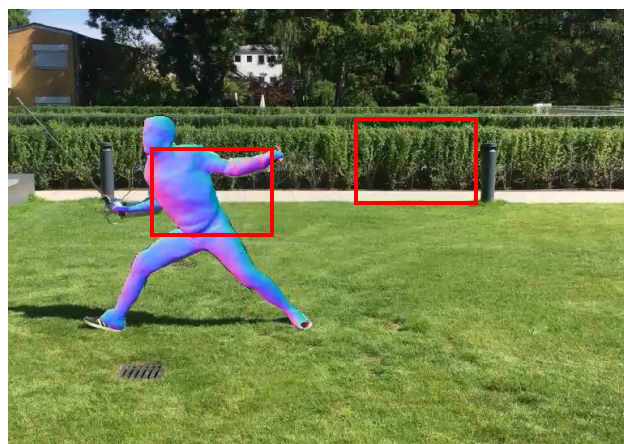
# Reconstruction Comparison on 3DPW



Video



ICON



SelfRecon



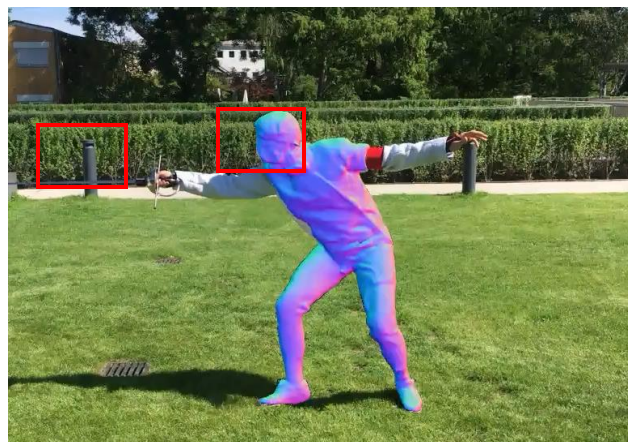
Ours



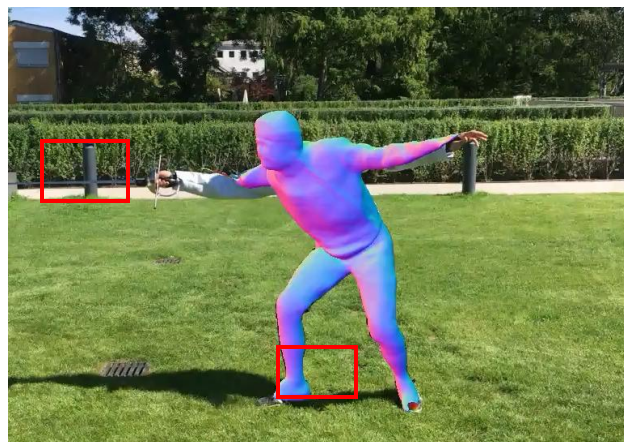
# Reconstruction Comparison on 3DPW



Video



ICON



SelfRecon



Ours

# Reconstruction Comparison on Online Video



Video



ICON



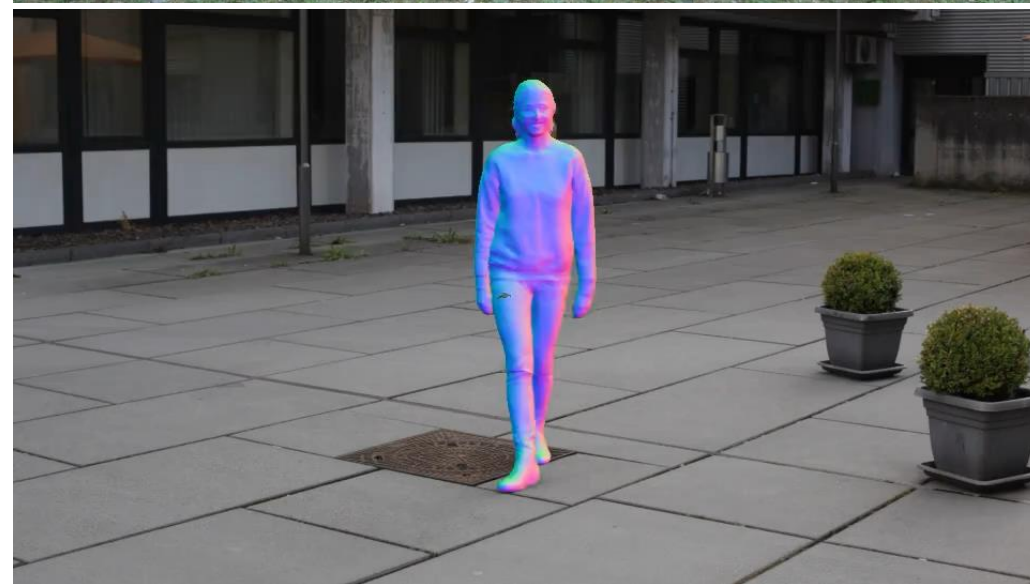
SelfRecon



Ours



# Qualitative Results on Monocular in-the-wild Videos – Datasets

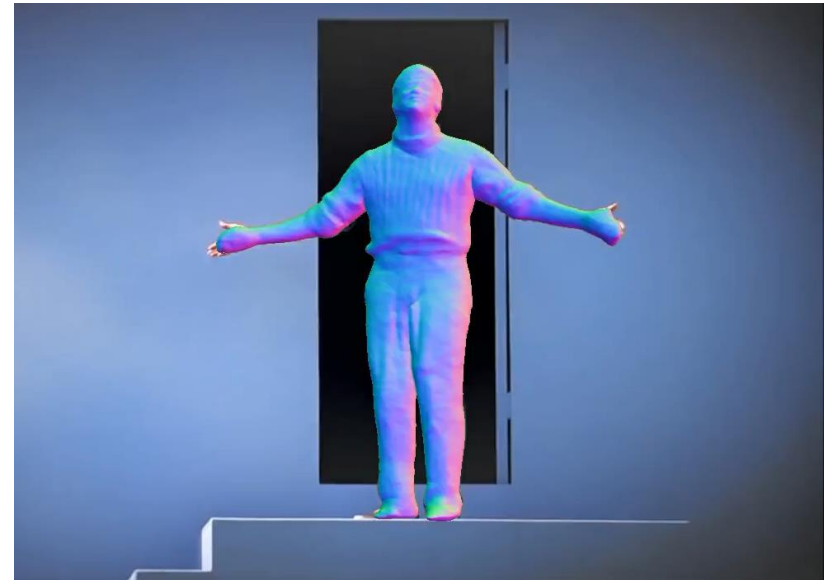


Video

Reconstruction



# Qualitative Results on Monocular in-the-wild Videos – Online Videos



Video

Reconstruction

# Qualitative Results on Monocular in-the-wild Videos – Self-captured Videos



Video



Reconstruction



# Qualitative Results on Monocular in-the-wild Videos – Self-captured Videos



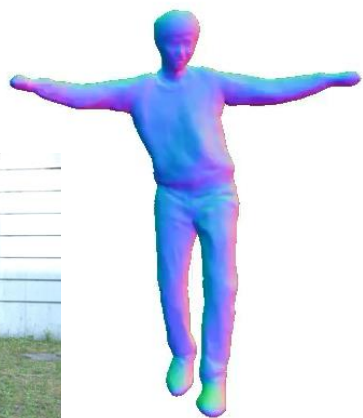
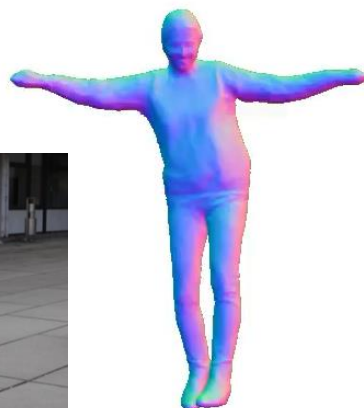
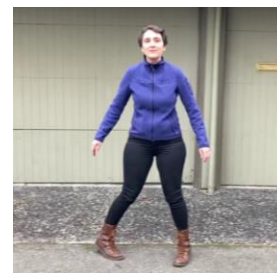
Video



Reconstruction

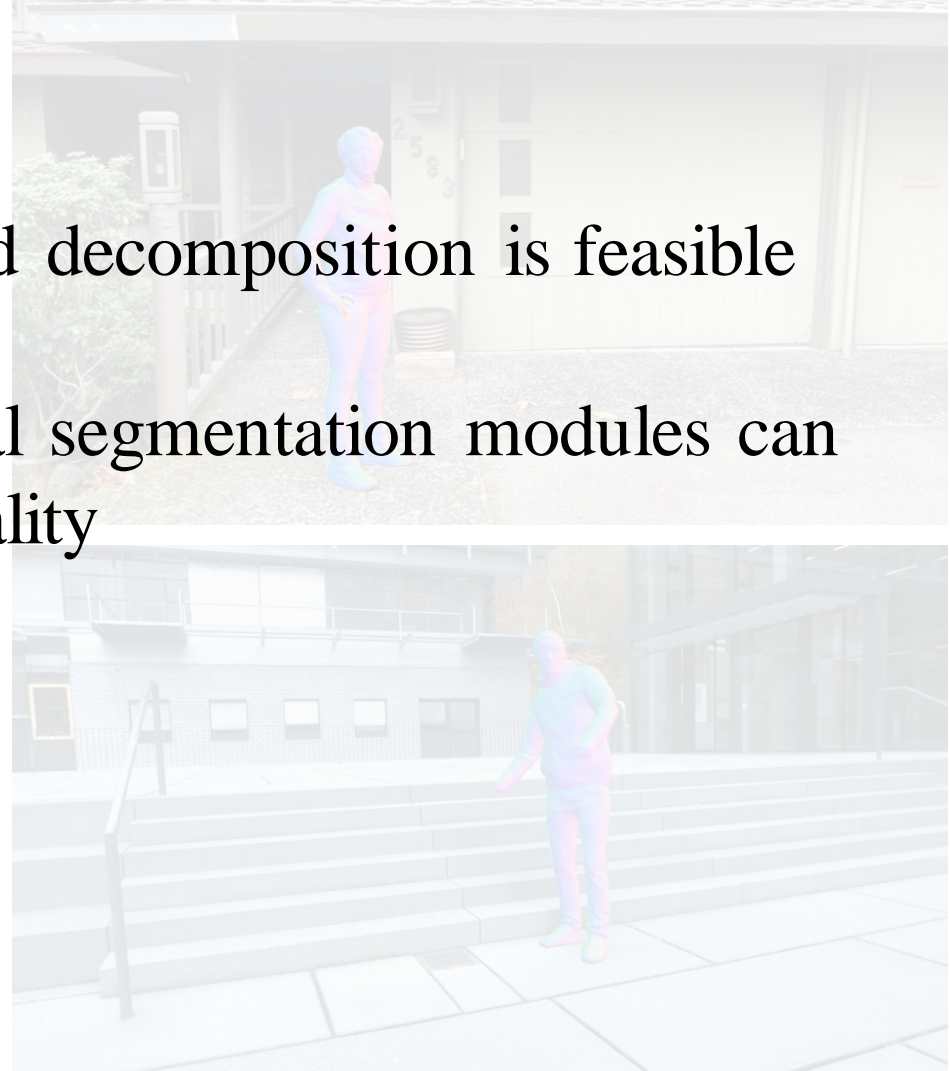
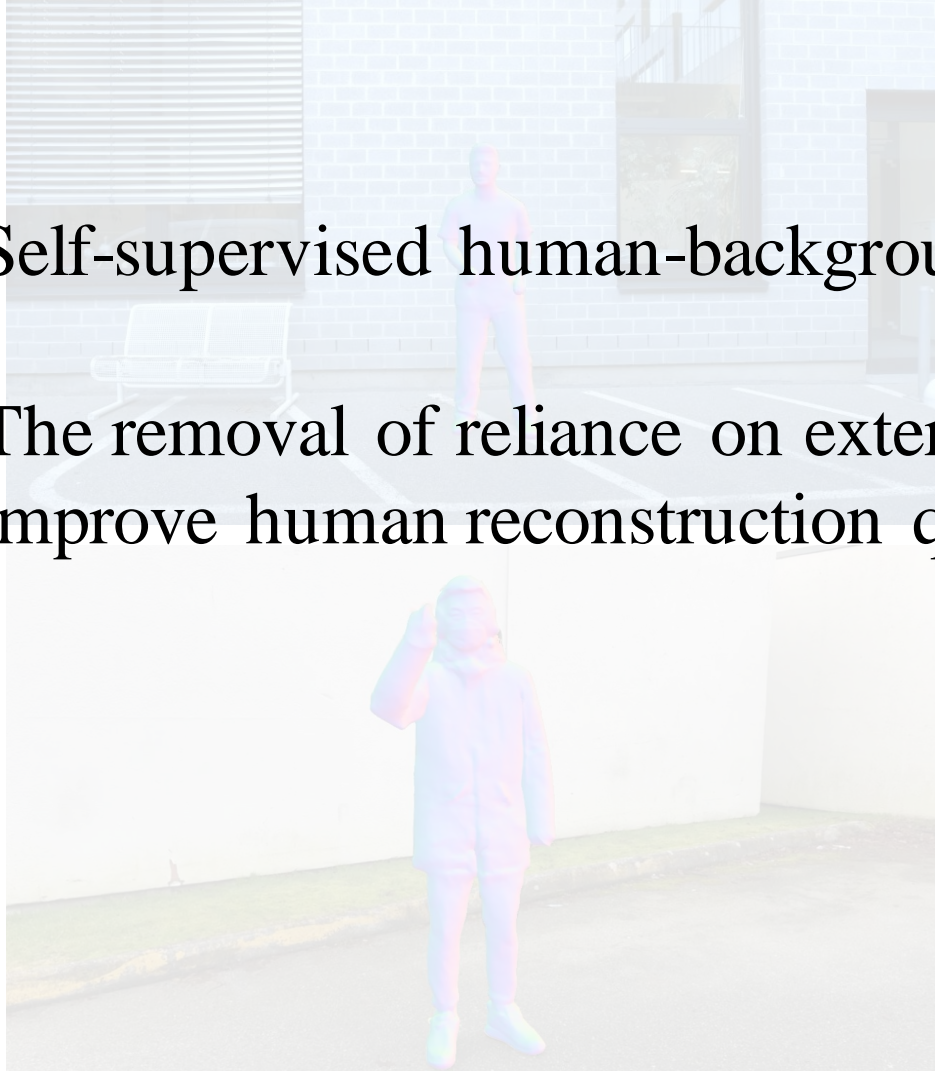


# Animation



# Summary

- Self-supervised human-background decomposition is feasible
- The removal of reliance on external segmentation modules can improve human reconstruction quality



# X-Avatar

## Expressive Human Avatars

Kaiyue Shen<sup>1\*</sup>

Chen Guo<sup>1\*</sup>

Manuel Kaufmann<sup>1</sup>

Juan Jose Zarate<sup>1</sup>

Julien Valentin<sup>2</sup>

Jie Song<sup>1</sup>

Otmar Hilliges<sup>1</sup>

\*Equal Contribution

<sup>1</sup>ETH Zurich

<sup>2</sup>Microsoft

CVPR 2023

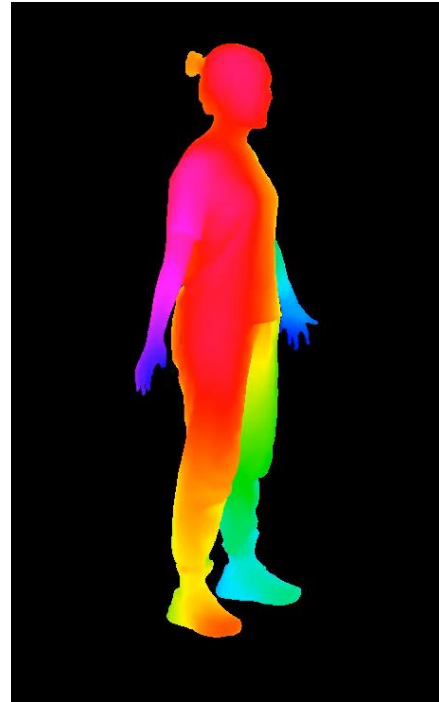


# Problem Setting



Scan input

OR



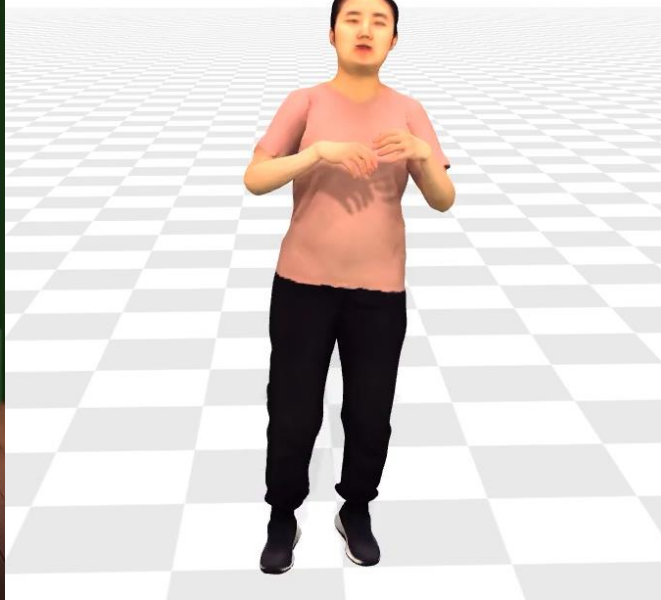
RGB-D input



Built X-Avatar



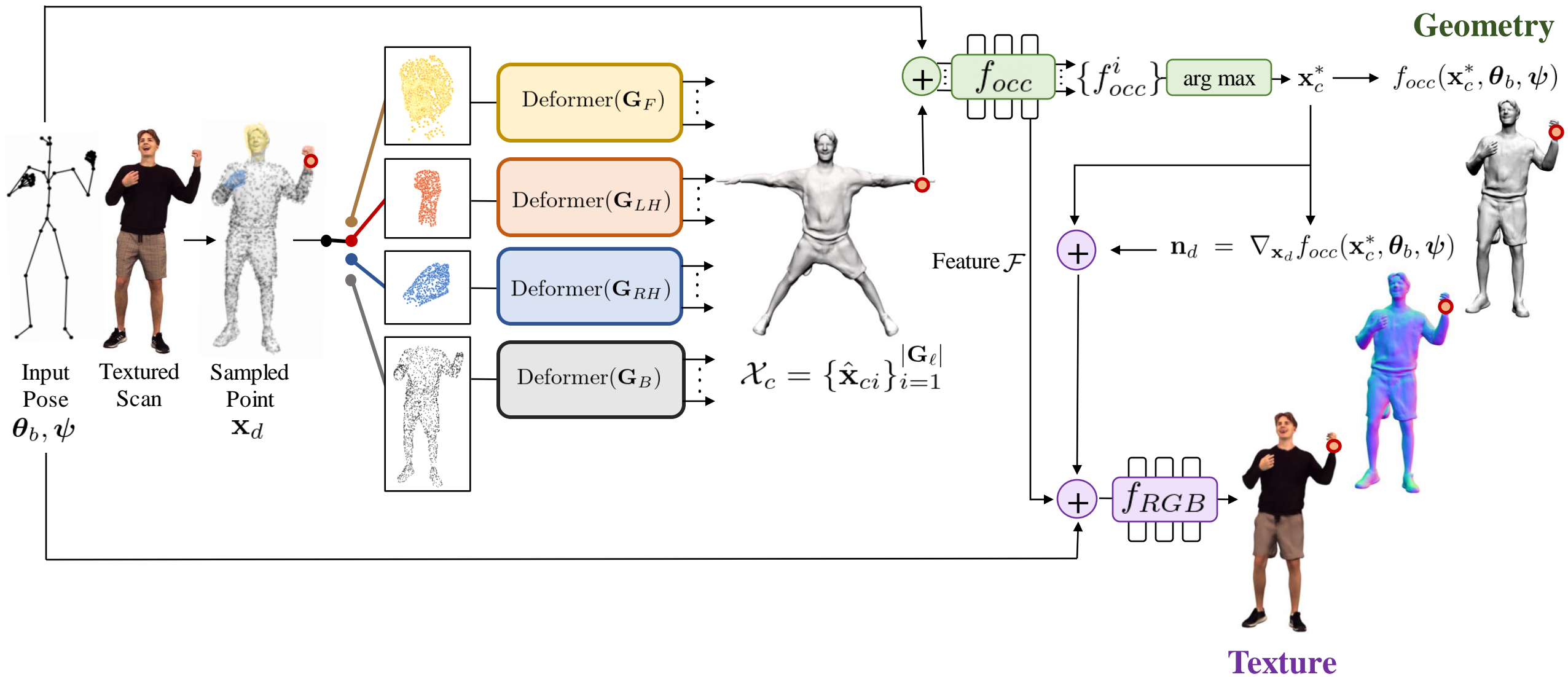
Driving pose



Animation



Zoom in



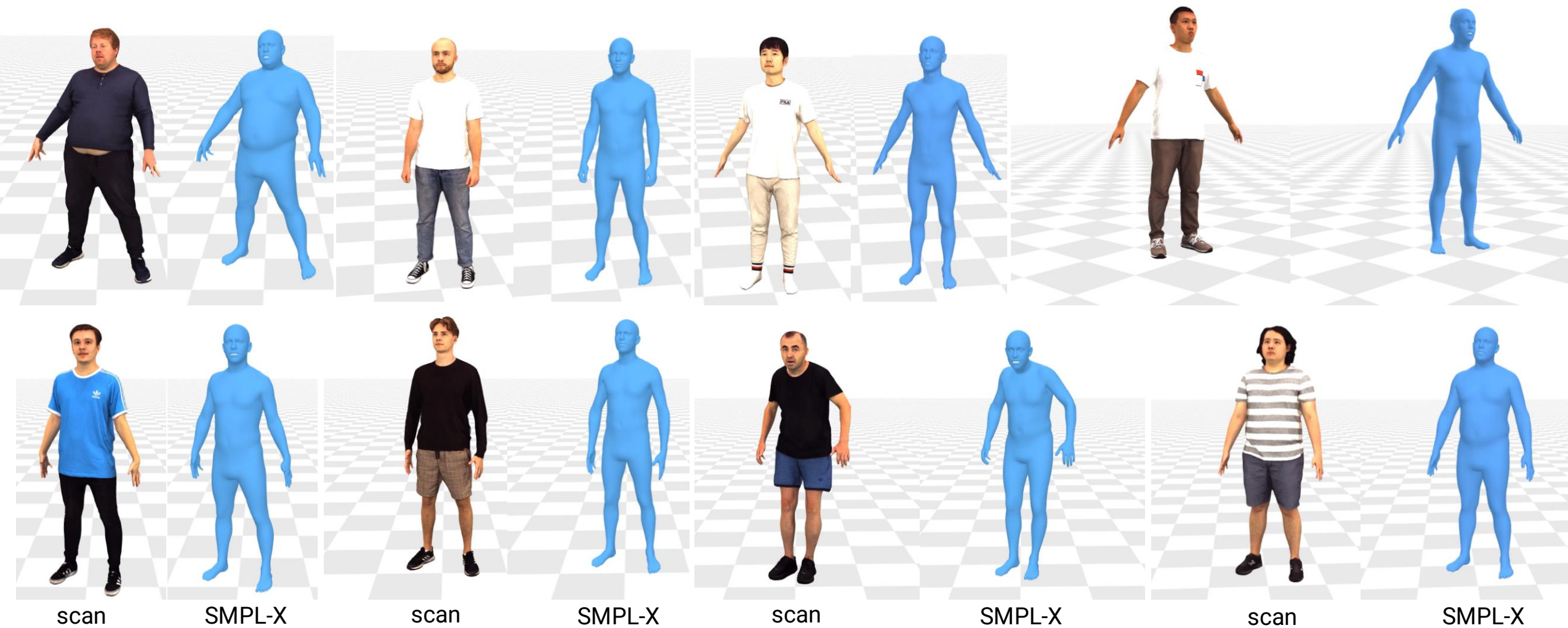


# X-Humans



# Male Subjects (shown 8 of 11)

- 20 subjects, 234 sequences, 34,663 frames
- Textured scans, SMPL[-X] registrations

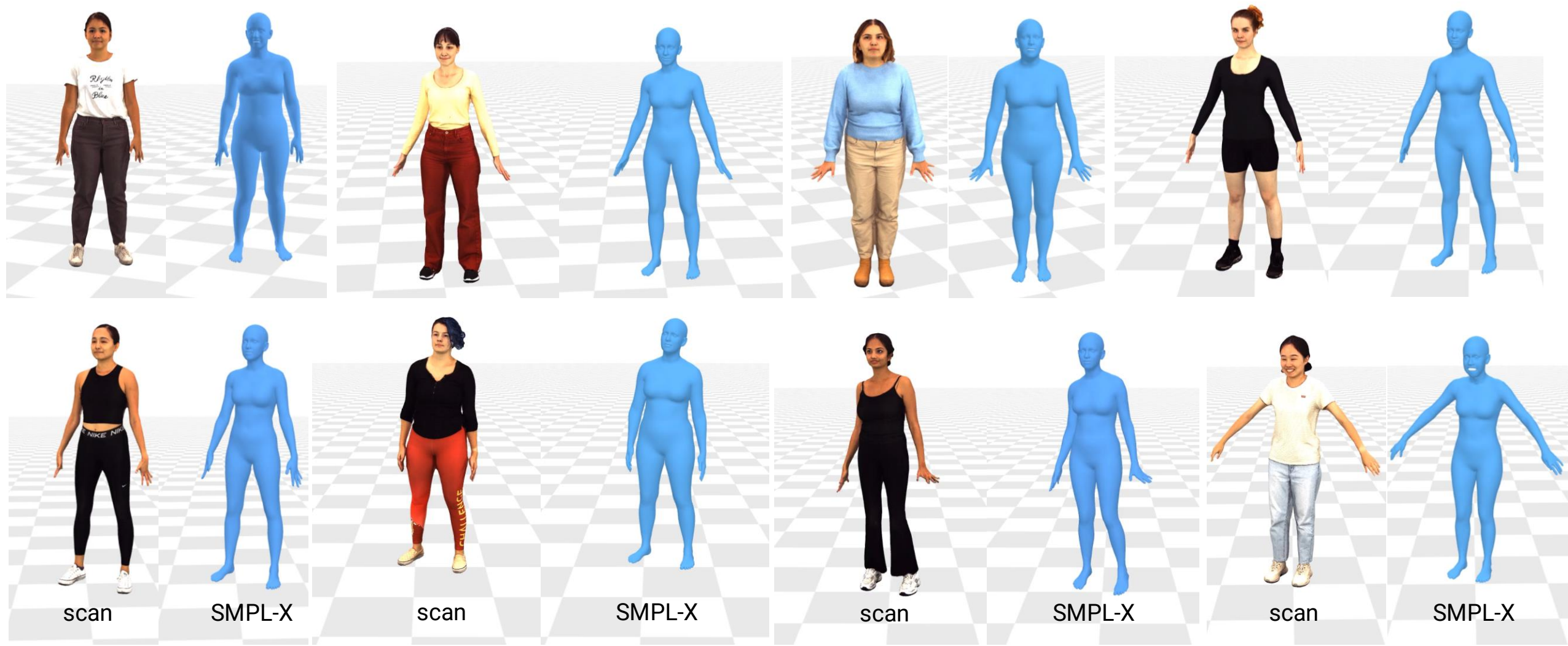


- 20 subjects, 234 sequences, 34,663 frames
- Textured scans, SMPL[-X] registrations

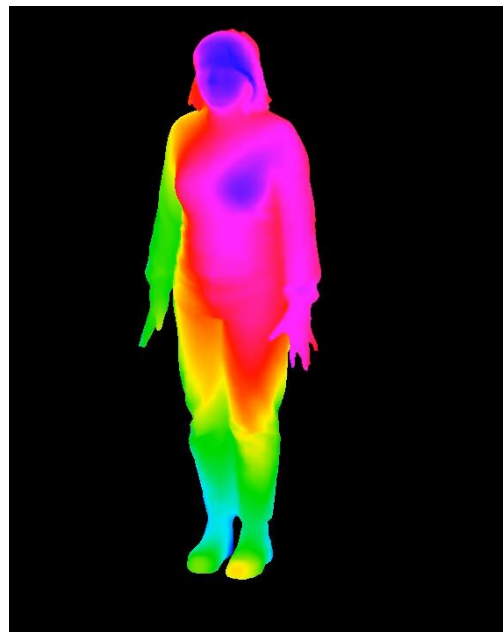


# Female Subjects (shown 8 of 9)

- 20 subjects, 234 sequences, 34,663 frames
- Textured scans, SMPL[-X] registrations
- Body pose + hand gesture + facial expression
- Various clothing types, hair styles, genders and ages



# Creating X-Avatars from RGB-D



RGB-D input

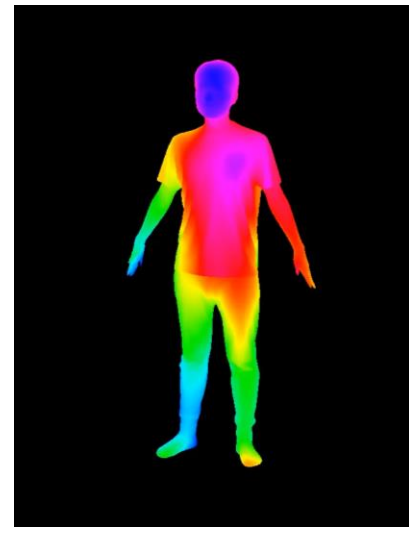
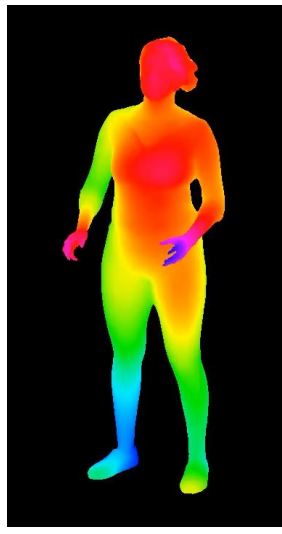
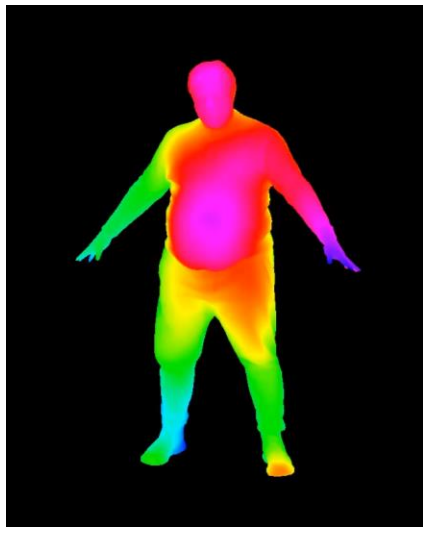


Built X-Avatar



Animation

RGB-D input



Animation





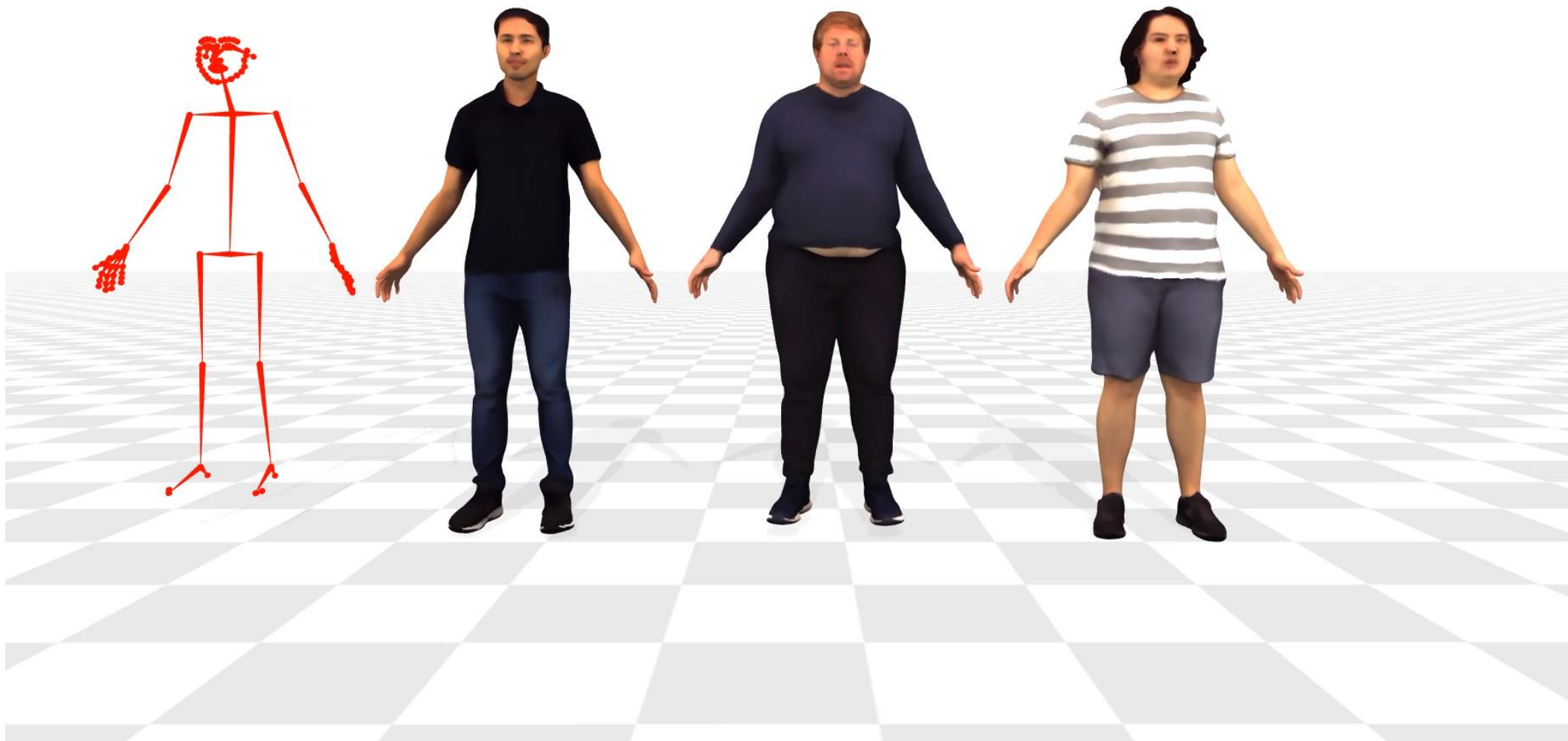
# Re-targeting (tennis)



# Re-targeting (dance)



# Re-targeting with 3D consistency





# Summary

- **X-Avatar**, one of the first expressive implicit human avatar models.
- An unified approach to efficiently build X-Avatars from scans and RGB-D data.
- **X-Humans**, a dataset of high-quality textured scans of clothed people performing varied body and hand movements and facial expressions.



# Future Work



# Reconstruct Anything

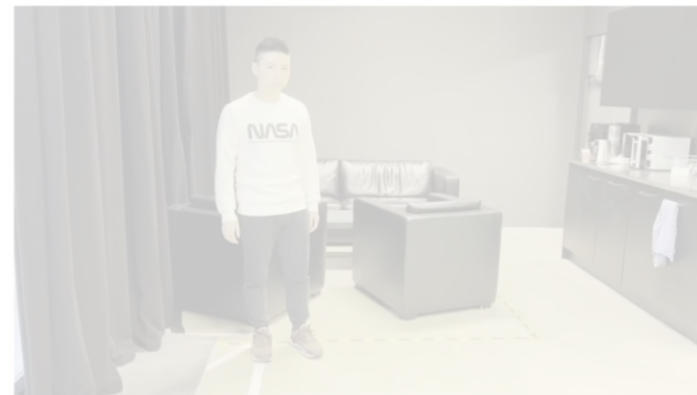
Vid2Avatar



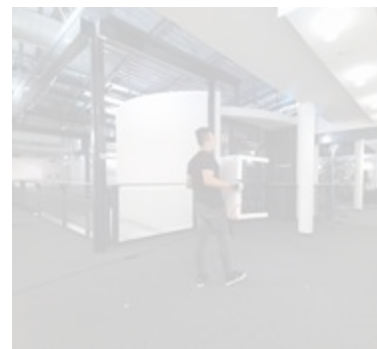
Vid2Anything



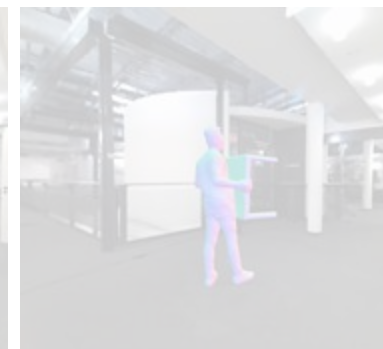
humans, scene, objects, etc.



Input: RGB video



Input



Reconstructions



Input



Reconstructions

Accepted as  
CVPR 2024 Oral



# Thank you!



Co-authors: Otmar Hilliges, Kaiyue Shen, Yifei Yin, Jie Song, Julien Valentin, Xu Chen,, Tianjian Jiang, Manuel Kaufmann, Juan Zarate