

Human-Understanding AI

Gyeongsik Moon



Visual
Computing
and AI Lab

- Gyeongsik Moon (문경식)
- Department of CSE (컴퓨터 학과)
- Computer Vision, Computer Graphics, Robotics
- Fresh faculty (first semester!)
- 1 year at DGIST, 2 years at Meta after finishing Ph.D. at SNU in 2021
- Website: www.vcai.korea.ac.kr



Home News Members Publications Joining Us

Visual Computing and AI Lab

Designing artificial intelligence
to perceive, represent, and model human-centric 3D worlds
through computer vision, computer graphics, and robotics

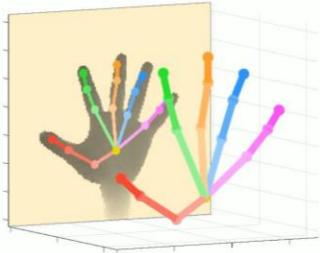
We have openings for graduate students and undergraduate interns: click [here](#) for information

We study three main directions—*human perception, human representation, and human behavior modeling*—as outlined below.

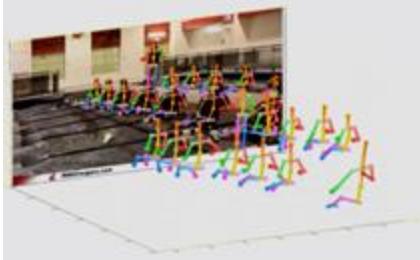


Computer**Vision**Lab
Seoul National University

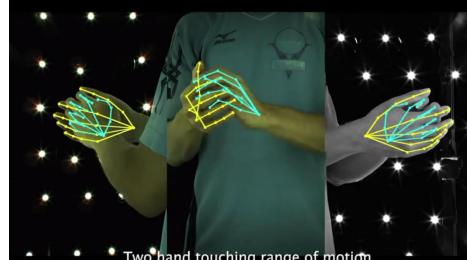
Computer Vision



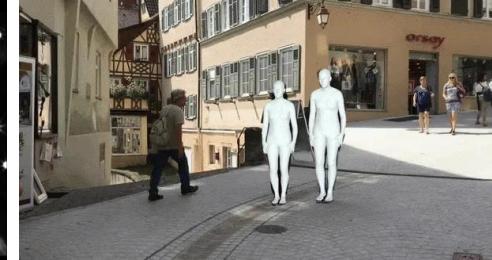
3D hand pose estimation
(CVPR 2018)



3D multi-person pose estimation
(ICCV 2019)



3D interacting hand pose estimation
(ECCV 2020)



3D body shape estimation
(ECCV 2020)



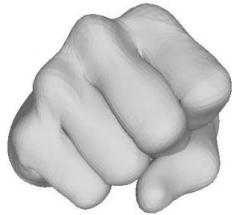
3D body shape estimation
(CVPR 2021)



Hand-object interaction
(CVPR 2022)



3D whole-body pose estimation
(CVPR 2022)



(b) DeepHandMesh (ours)



(c) 3D reconstruction

High-fidelity 3D hand geometry model
(ECCV 2020)



Relighted 3D interacting hands
(NeurIPS 2023)



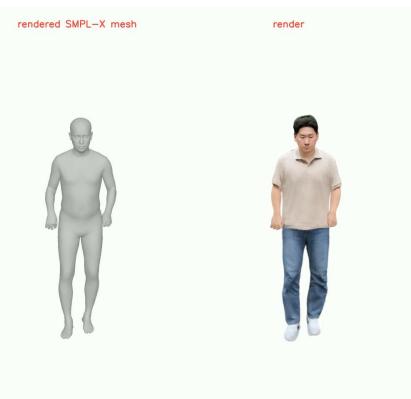
Authentic 3D hand avatar
(CVPR 2024)



Universal relightable hand model
(CVPR 2024)



Expressive Whole-Body 3D Gaussian Avatar
(ECCV 2024)



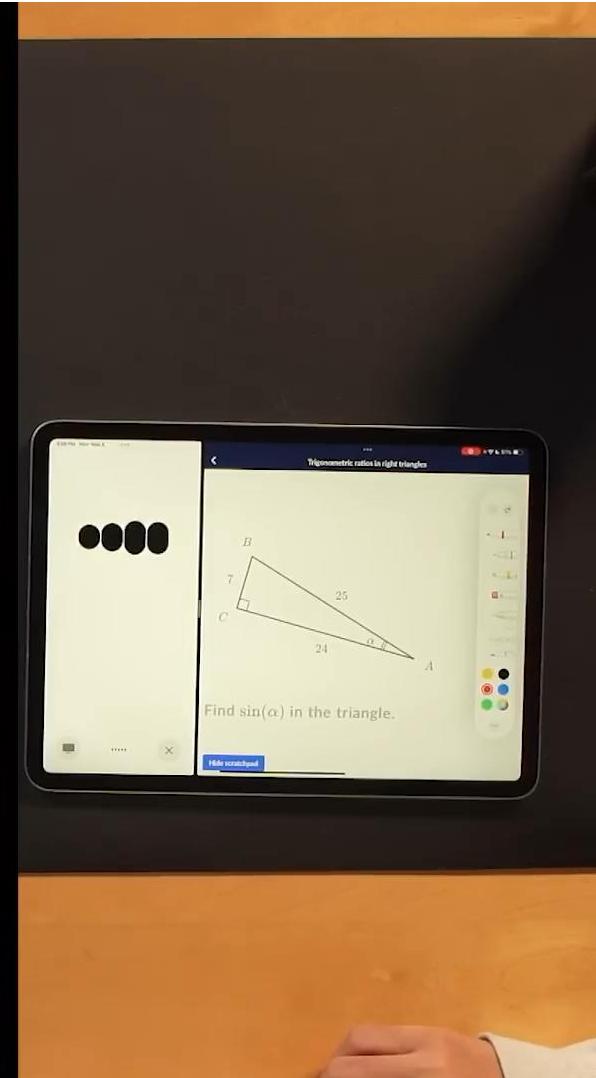
```
graph TD; A((Artificial Intelligence)) --- B((Computer Vision)); A --- C((Computer Graphics)); A --- D((Machine Learning))
```

Artificial
Intelligence

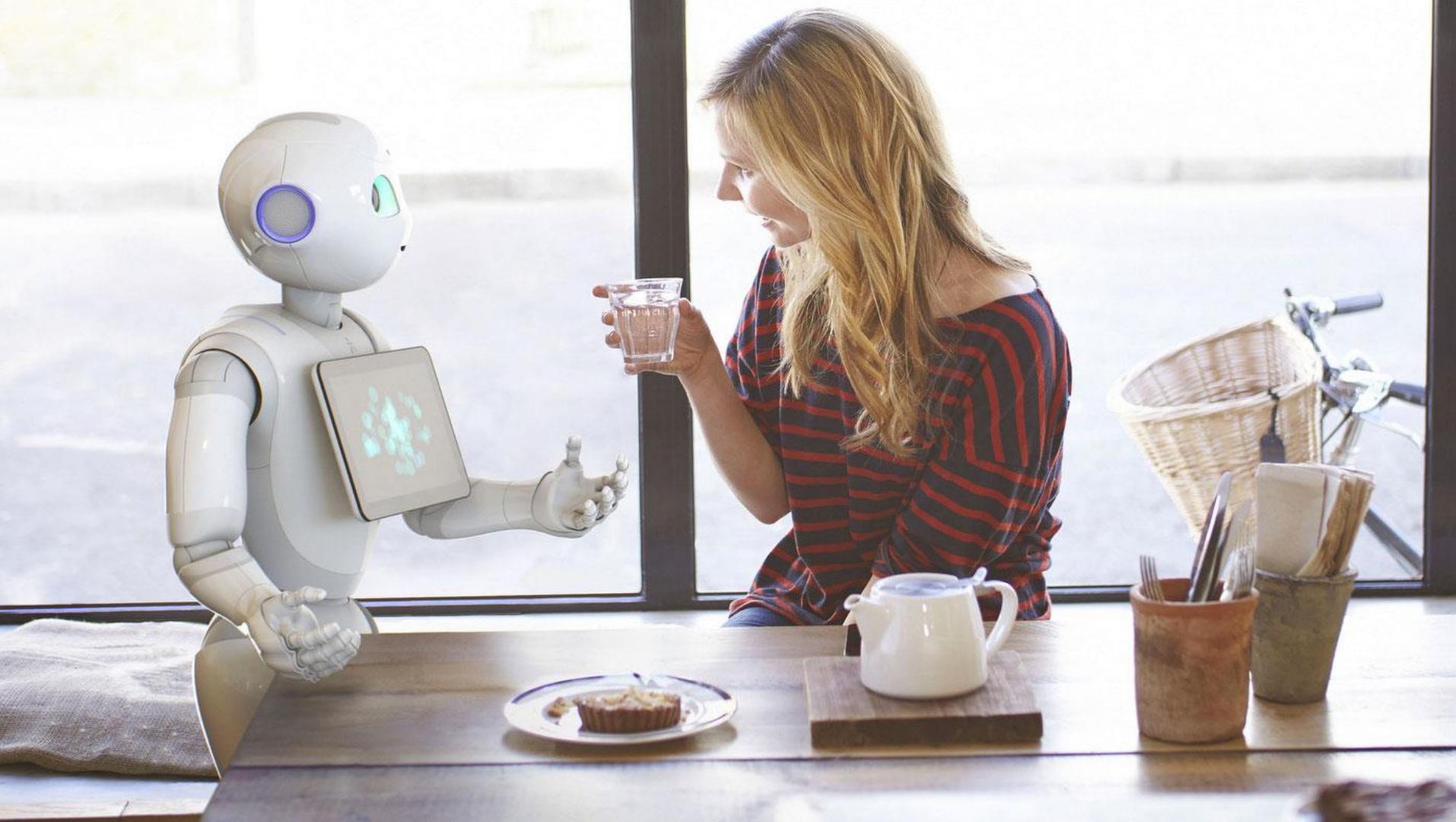
Computer
Vision

Computer
Graphics

Machine
Learning

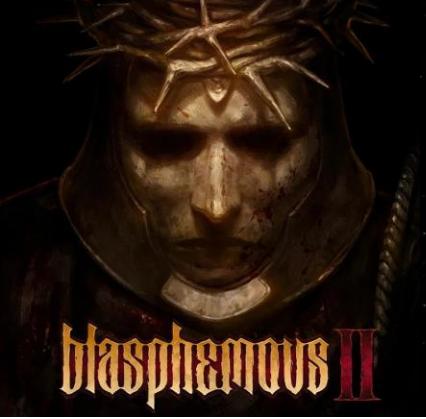






A woman with long blonde hair, wearing a red and black patterned dress, is seated at a table, facing a white humanoid robot. The robot has a large, round head with a blue and white sensor array. They are both looking at a small white mug on the table. On the table, there is also a white teapot, a small plate with a piece of cake, and a glass of water. The background shows a bright room with a window.

**Understand humans, and
eventually, interact with humans**



Humans are at **EVERYWHERE**



What makes it challenging?



“A Korean guy is breakdancing” with SORA of OpenAI



LUMA





Human image: Lots of physical attributes

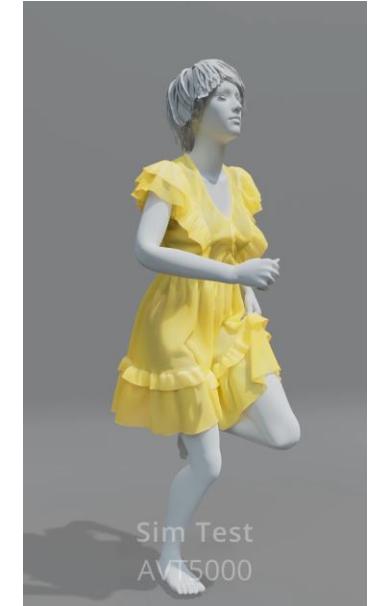
Pose



Facial expression



Cloth dynamics



Viewpoint



Illumination



Human image: Lots of physical attributes

- Human is the **most difficult** object to represent
- In particular, it is **dynamic** (e.g., diverse body poses) in contrast to **static scenes or objects**

Dynamic humans

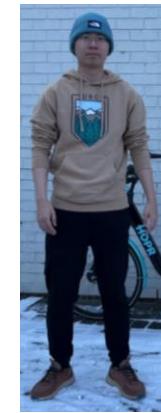


Static scenes

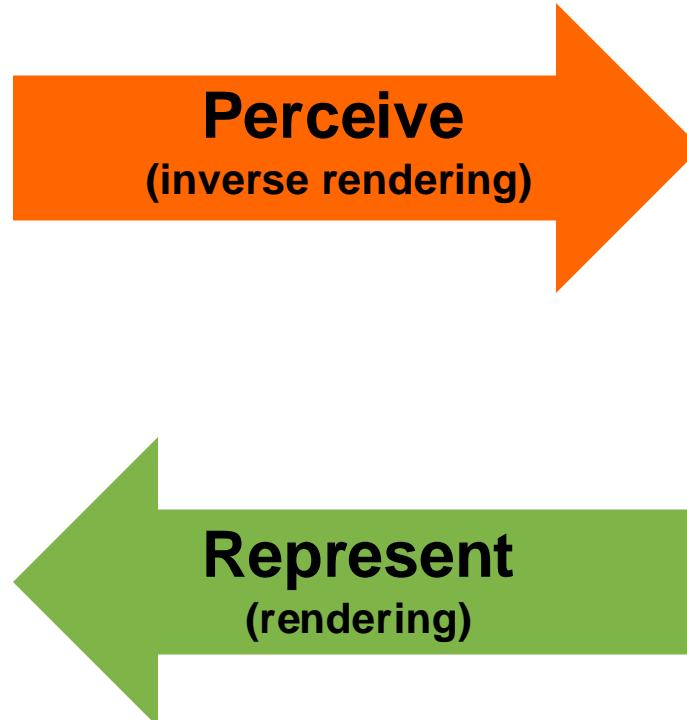


>
More
challenging

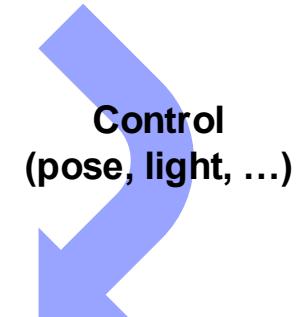
Perceive and Represent Human Images



Human images

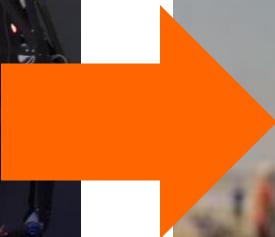


Human features
(e.g., geometry, texture, light, latent features, ...)

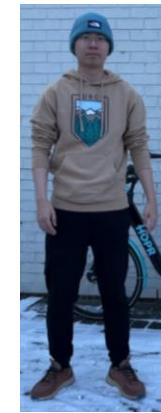


Understanding Humans in Casual Environments

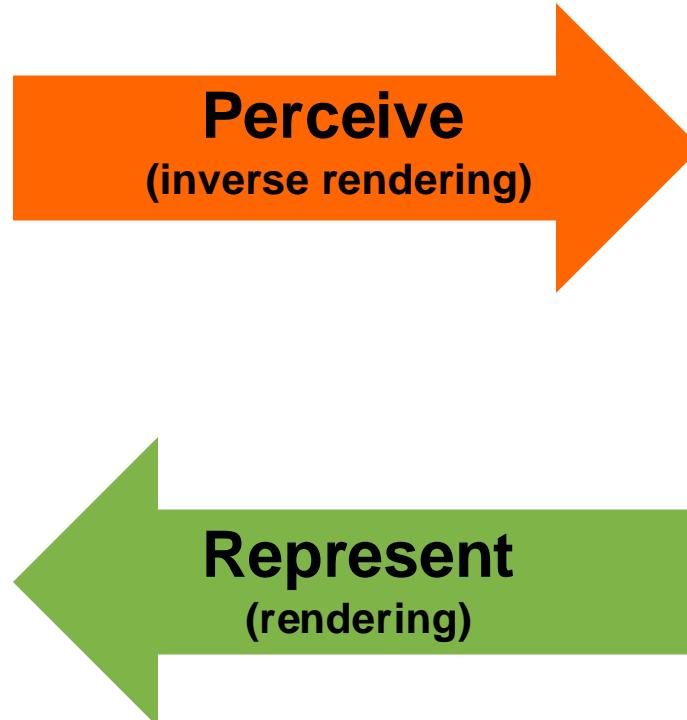
- Casual environments have much less available observations
- Need to use **priors to complement limited observations**



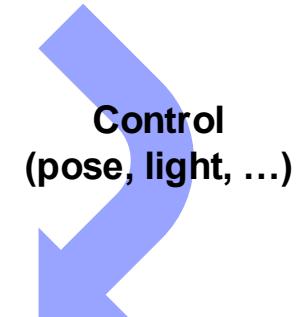
Perceive and Represent Human Images



Human images



Human features
(e.g., geometry, texture, light, latent features, ...)



URHand: Universal Relightable Hands

Zhaoxi Chen^{1,2}, Gyeongsik Moon¹, Kaiwen Guo¹, Chen Cao¹, Stanislav Pidhorskyi¹,
Tomas Simon¹, Rohan Joshi¹, Yuan Dong¹, Yichen Xu¹, Bernardo Pires¹,
He Wen¹, Lucas Evans¹, Bo Peng¹, Julia Buffalini¹, Autumn Trimble¹,
Kevyn McPhail¹, Melissa Schoeller¹, Shou-I Yu¹, Javier Romero¹,
Michael Zollhöfer¹, Yaser Sheikh¹, Ziwei Liu², Shunsuke Saito¹

¹**Codec Avatars Lab, Meta**

²**Nanyang Technological University**



Why Relighting?

Generalize to novel illuminations





Given large-scale light-stage data...





Light-stage Data



URHand: Universal Prior for Relightable Hands



Generalize to Novel Views, Poses,
Identities, and Illuminations

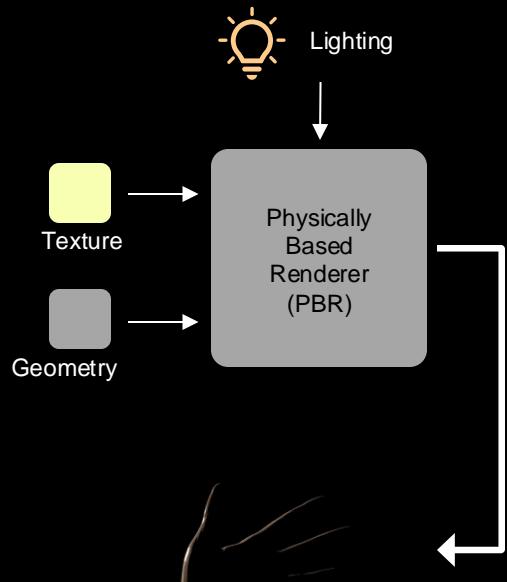


URHand: Universal Prior for Relightable Hands



URHand: Universal Prior
for Relightable Hands

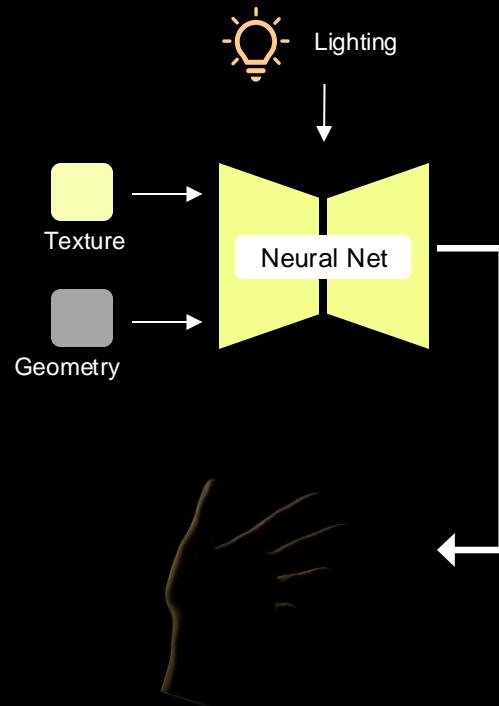
→ Your Hand Relightable
Quick Personalization from a Phone Scan



Physically Based Relighting

+ Generalizable in extrapolation

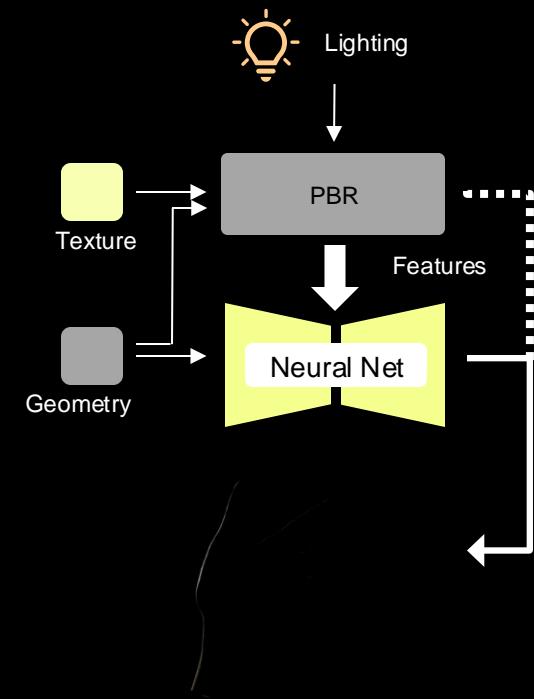
- Slow synthesis and low photorealism



Neural Relighting

+ Fast synthesis and high photorealism

- Hard to generalize in extrapolation

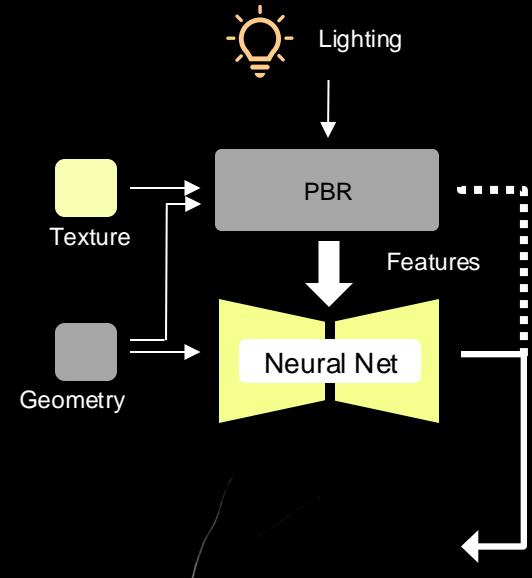
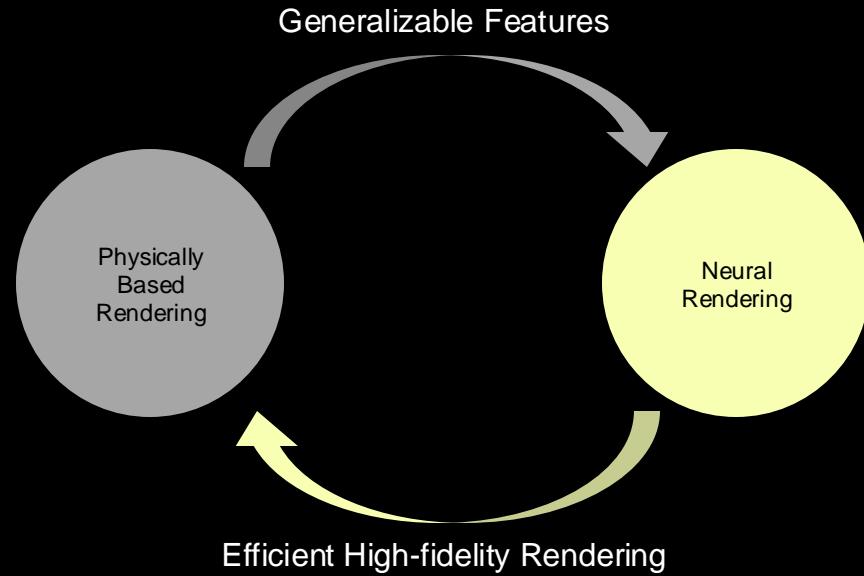


Neural-Physical Relighting

+ Generalizable in extrapolation

+ Fast synthesis and high photorealism

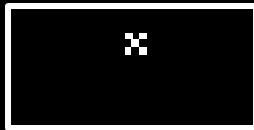
Our Solution: Neural-Physical Relighting



Neural-Physical Relighting

- + Generalizable in extrapolation
- + Fast synthesis and high photorealism

A Gap between Training and Inference



Training

One Light At a Time (OLAT)



Inference

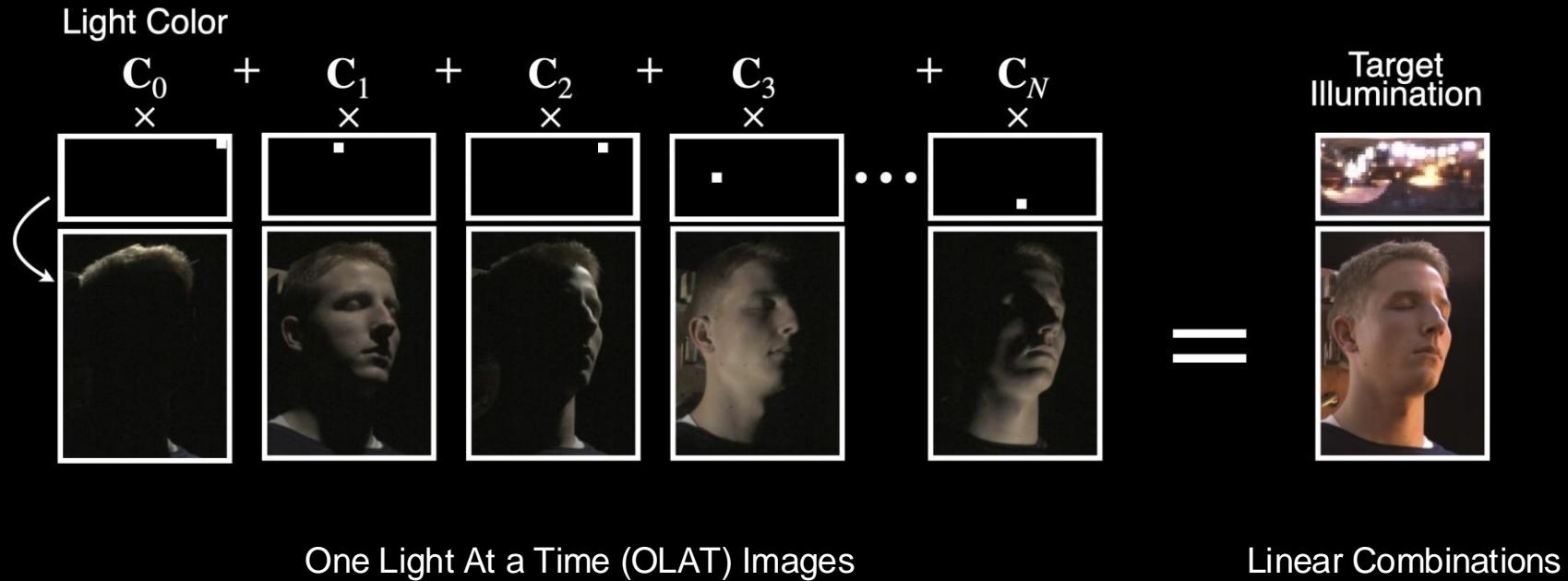
Combination of multiple natural lights



Generalize?

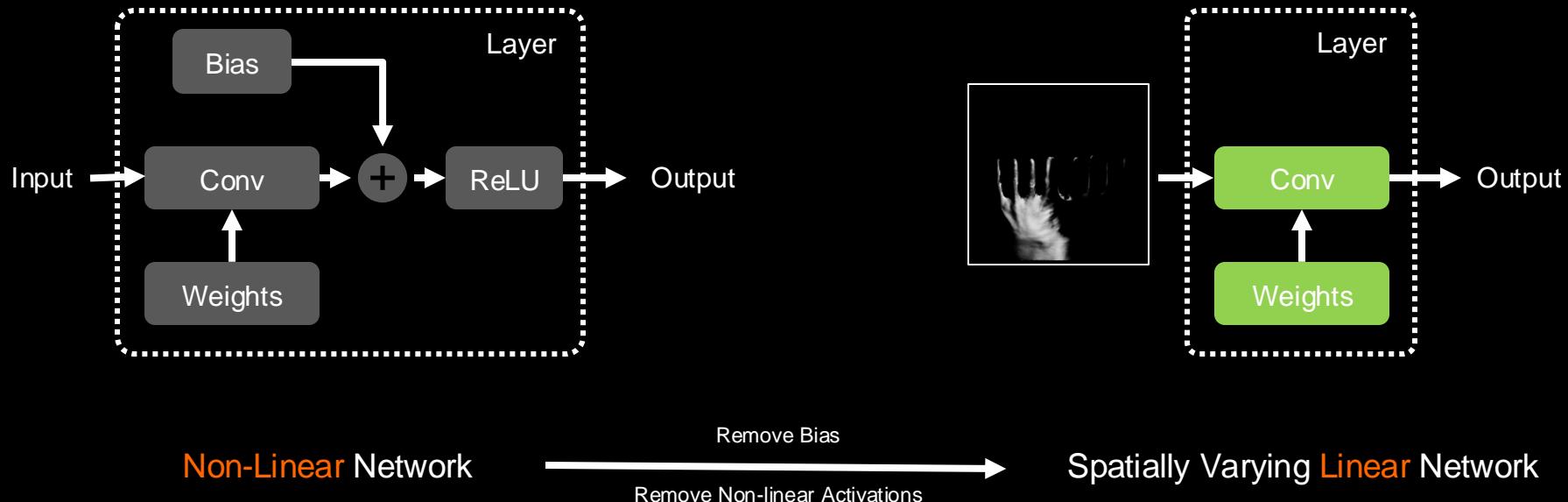


Recap: Linearity of Light Transport

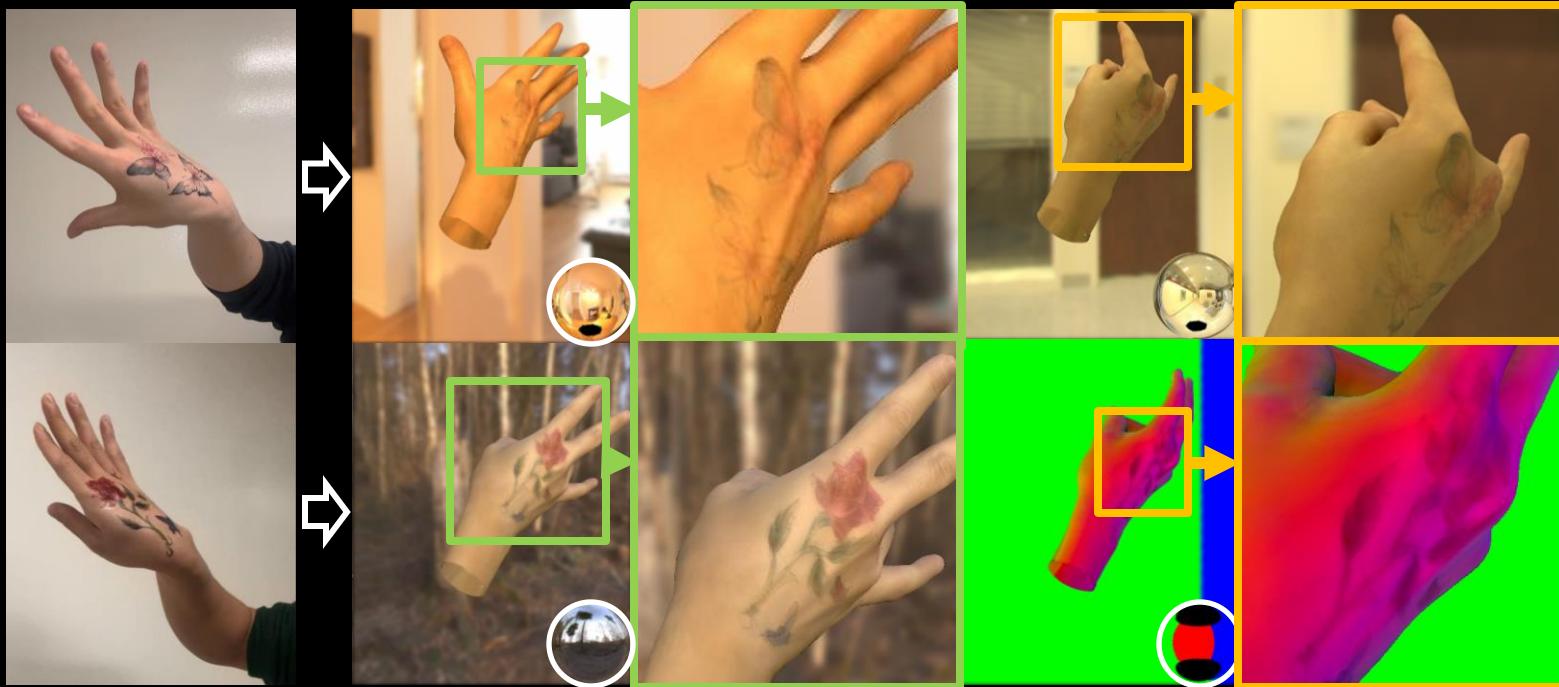


Key Idea: Linear Lighting Model

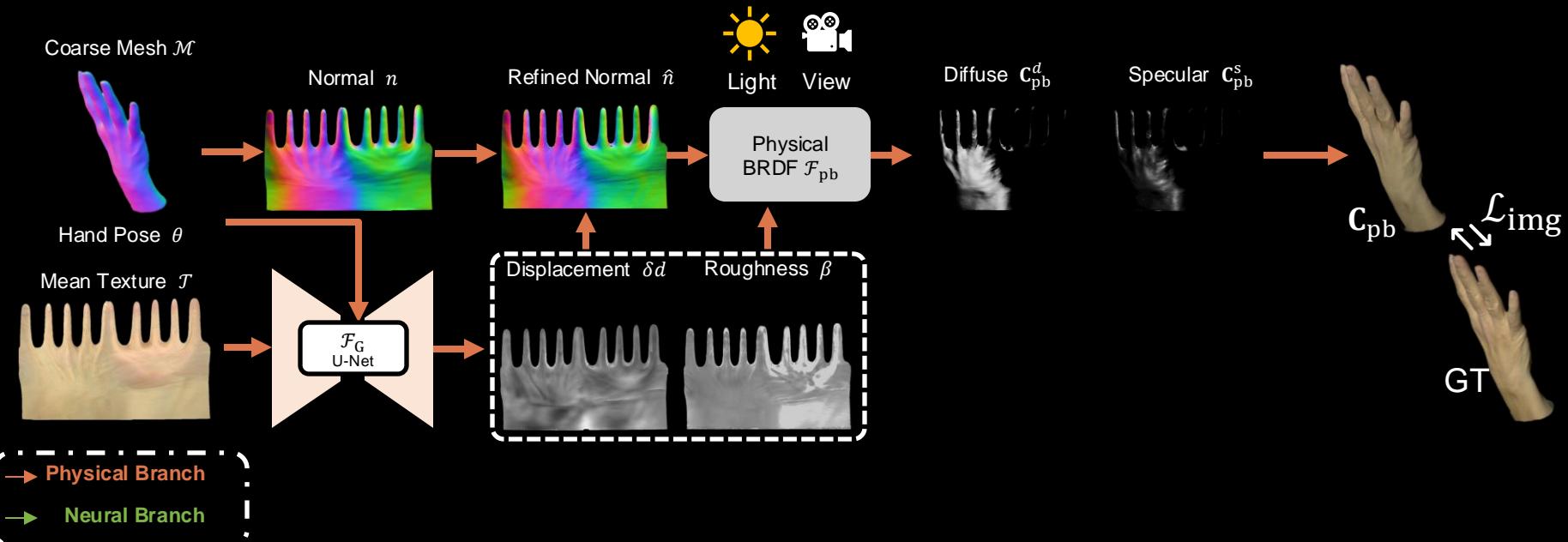
$$\sum \mathbf{C}_i f_{\text{linear}}(\mathbf{b}_i) = f_{\text{linear}}(\sum \mathbf{C}_i \mathbf{b}_i)$$



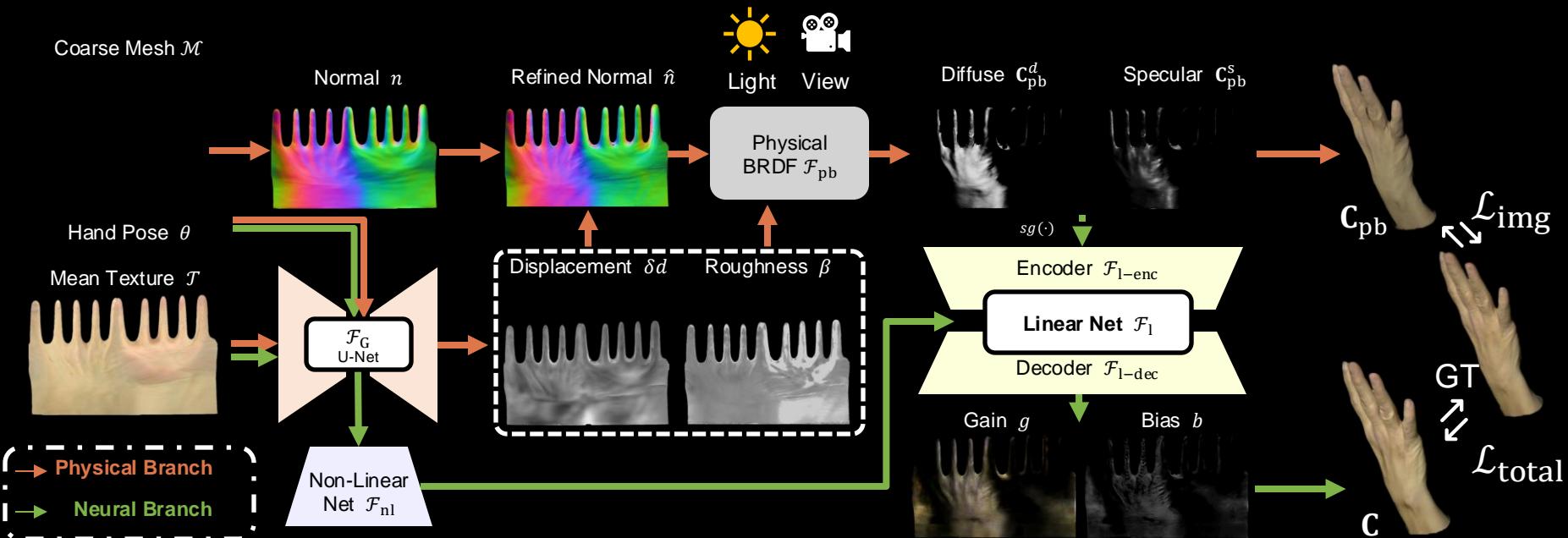
Key Idea: Linear Lighting Model



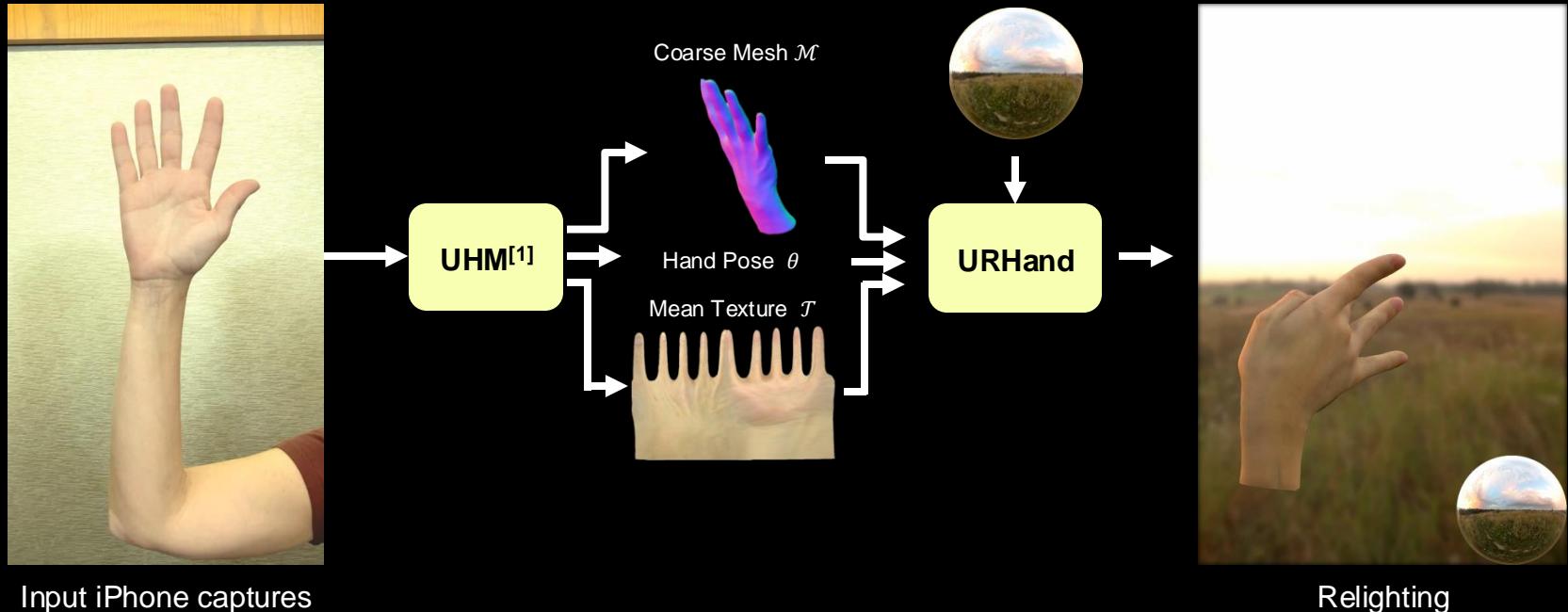
URHand: Physical Branch



URHand: Neural Branch



URHand: Inference



Input iPhone captures

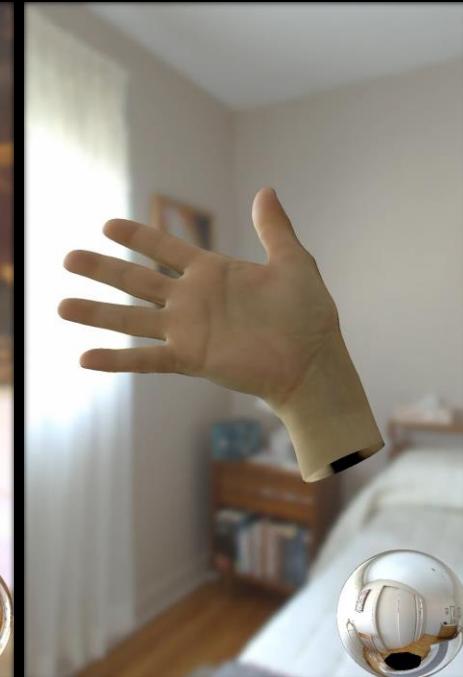
Relighting

Relight with Environment Map



URHand \Rightarrow Your Hand

Quick Personalization of a Relightable Hand from Phone Scan



Input Phone Scan

Photorealistic Rendering under Arbitrary Illuminations

Expressive Whole-Body 3D Gaussian Avatar

<https://mks0601.github.io/ExAvatar/>



Gyeongsik Moon



Takaaki Shiratori



Shunsuke Saito

Expressive Whole-Body 3D Avatar from a Monocular Video

A monocular video



Driving with novel **body poses, hand poses, and facial expressions**

rendered SMPL-X mesh

render



What makes it challenging?



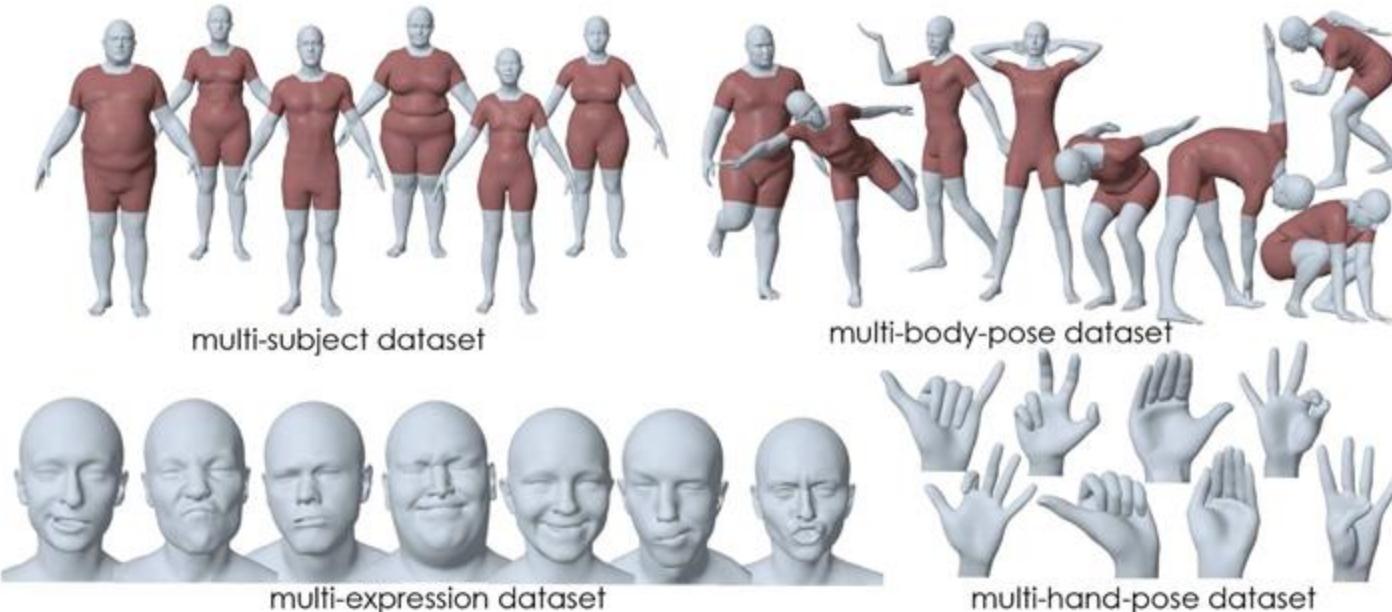
Limited training frames

- Limited body poses
- Limited hand poses
- Limited facial expressions

How can we make it generalize well
to novel poses and facial
expressions?

Strong Priors from Mesh-Based Models

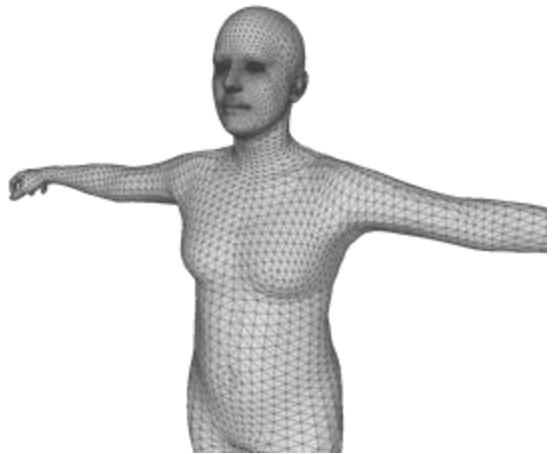
- 3D mesh-based models (e.g., SMPL-X [1] and GHUM [2]) already support animation with novel body poses, hand poses, and facial expressions
- They provide **strong geometric priors for the 3D human animations**



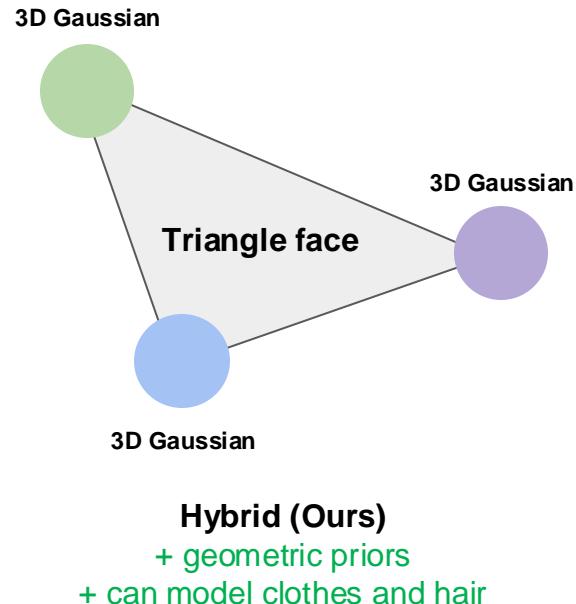
[1] Pavlakos et al. "Expressive body capture: 3D hands, face, and body from a single image." CVPR. 2019.

[2] Xu et al. "Ghum & ghum: Generative 3D human shape and articulated pose models." CVPR. 2020.

Hybrid Representation of 3DGS and Surface Mesh



Surface mesh
+ geometric priors
- hard to model clothes and hair



Hybrid Representation of 3DGS and Surface Mesh

We can utilize useful surface-based regularizers
thanks to the hybrid representation!



(a) With Lap. reg. (Ours)



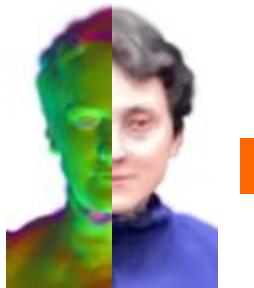
(b) Without Lap. reg.



(c) Without Lap. reg. + strong L2 reg.

Hybrid Representation of 3DGS and Surface Mesh

We can utilize useful **surface-based face loss**
thanks to the **hybrid representation!**



(a) With the face loss (Ours)

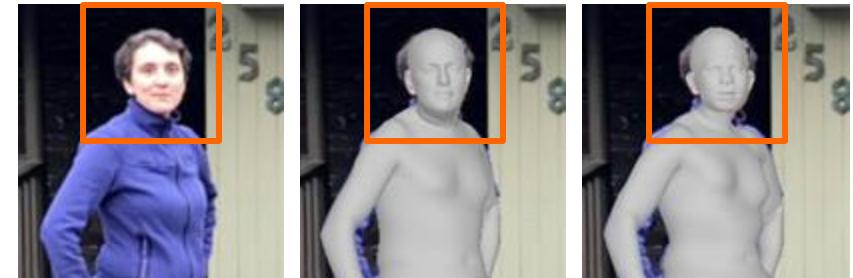
Mouth geometry is at the correct position



(b) Without the face loss

Mouth geometry is below the lower lip

Co-Registration of Body, Hands, and Face



(a) Image

(b) With
joint offset (Ours)

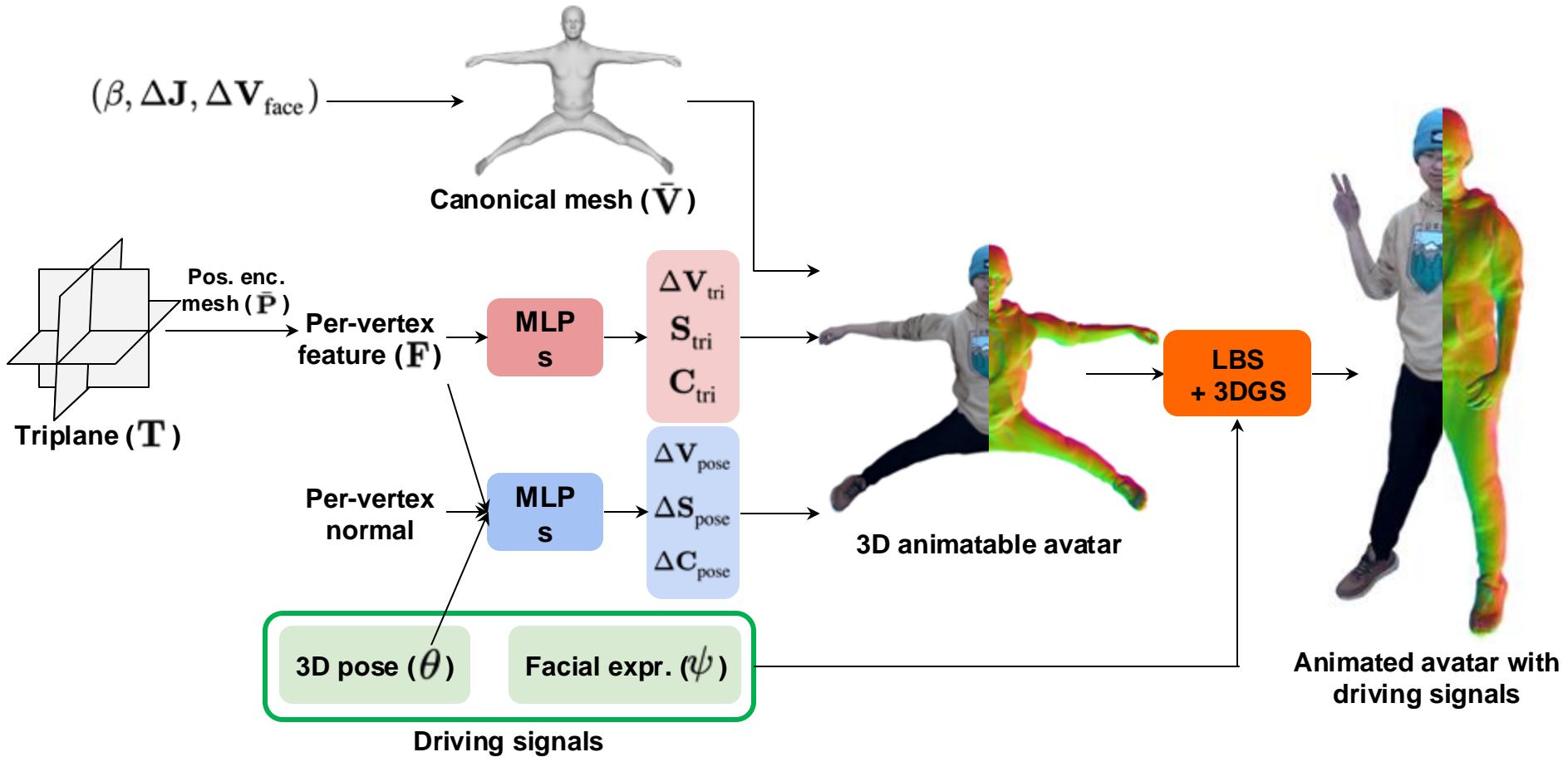
(c) Without
joint offset

(a) Image

(b) With
face offset (Ours)

(c) Without
face offset

Architecture



Comparison to previous works



ExAvatar (Ours)



3DGS-Avatar



ExAvatar (Ours)

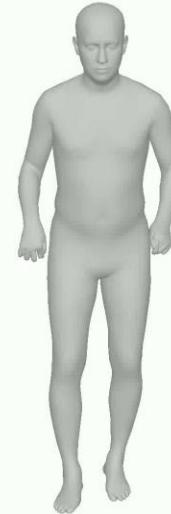


3DGS-Avatar

Motion transfer from in-the-wild videos



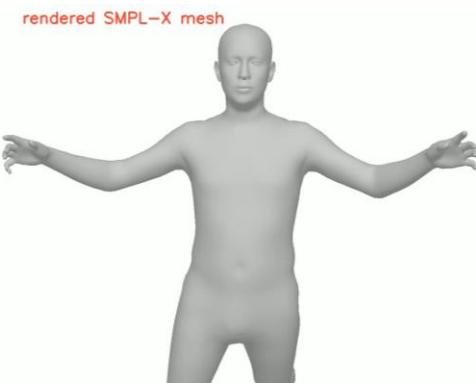
rendered SMPL-X mesh



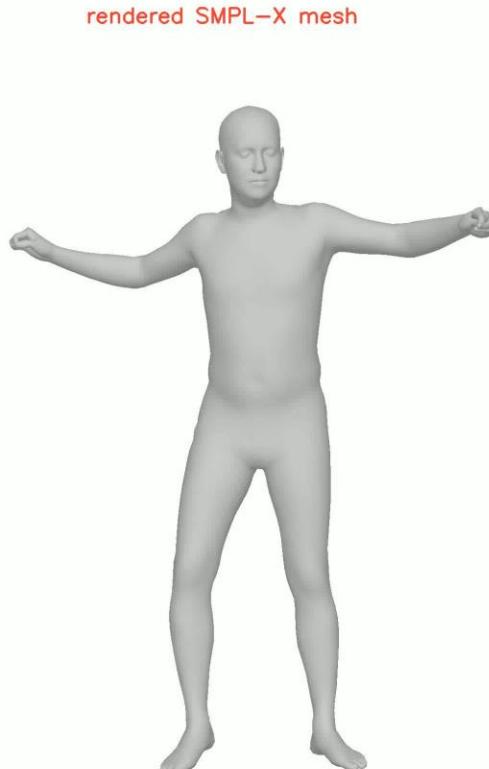
render



Motion transfer from in-the-wild videos



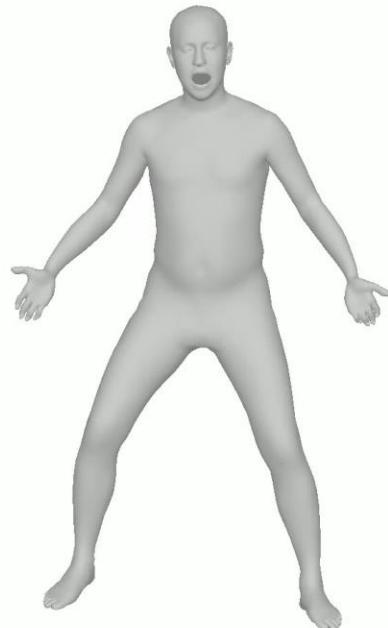
Motion transfer from in-the-wild videos



Motion transfer from in-the-wild videos



rendered SMPL-X mesh



render



Motion transfer from in-the-wild videos



KOREA UNIVERSITY
Department of
Computer Science and Engineering

고려대학교 응원가 – 민족의 아리아



Render from any viewpoints with 3D avatar



rendered SMPL-X mesh



render



We just need *a casually captured video*



rendered SMPL-X mesh

render



Ph.D. DISSERTATION

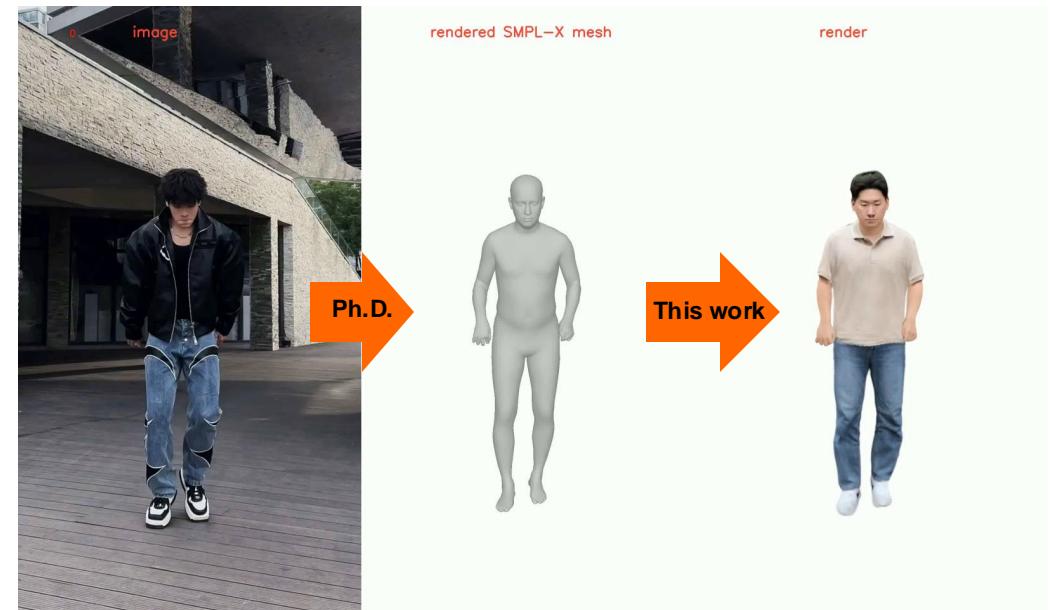
Expressive Whole-Body 3D Multi-Person Pose and Shape Estimation from a Single Image

단일 이미지로부터 여러 사람의
표현적 전신 3D 자세 및 형태 추정

BY
GYEONGSIK MOON

February 2021

DEPARTMENT OF ELECTRICAL ENGINEERING AND
COMPUTER SCIENCE
COLLEGE OF ENGINEERING
SEOUL NATIONAL UNIVERSITY



PERSONA: Personalized Whole-Body 3D Avatar with Pose-Driven Deformations from a Single Image

Work in progress

How about a *single image*?

Artificially captured video



A single image



A single image is **much more accessible** than artificially captured videos

How about a *single image*?

Directly applying existing avatar works -> bad results due to limited visibility



(a) Reference image



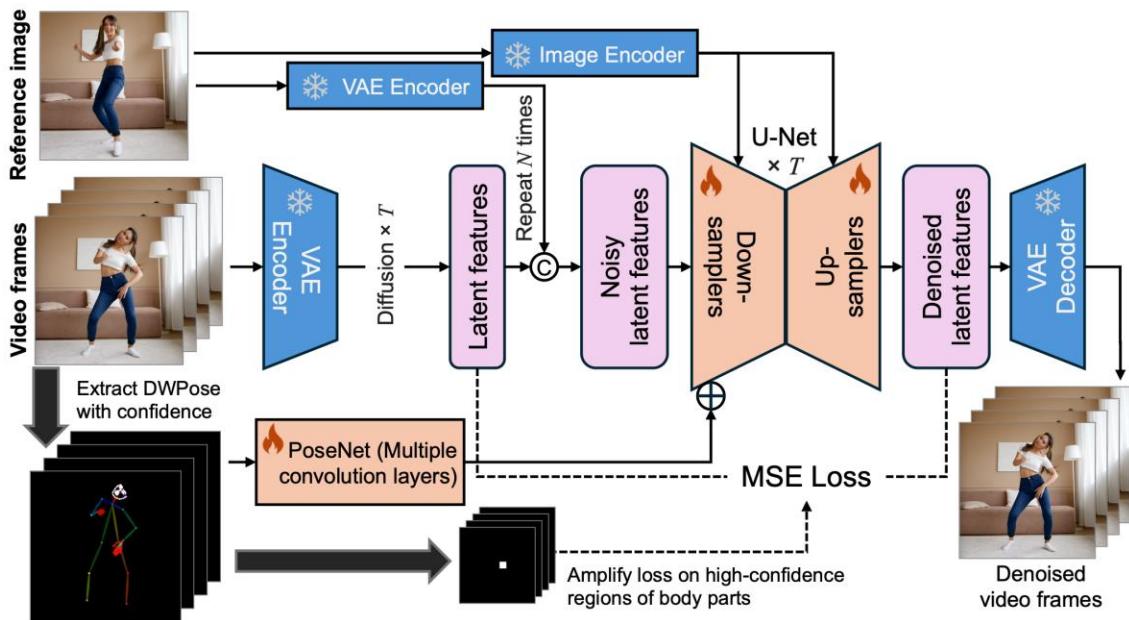
(c) ExAvatar

(b) Target pose

Prior from generative models

MagicAnimate [1]/AnimateAnyone [2]/MimicMotion [3]:

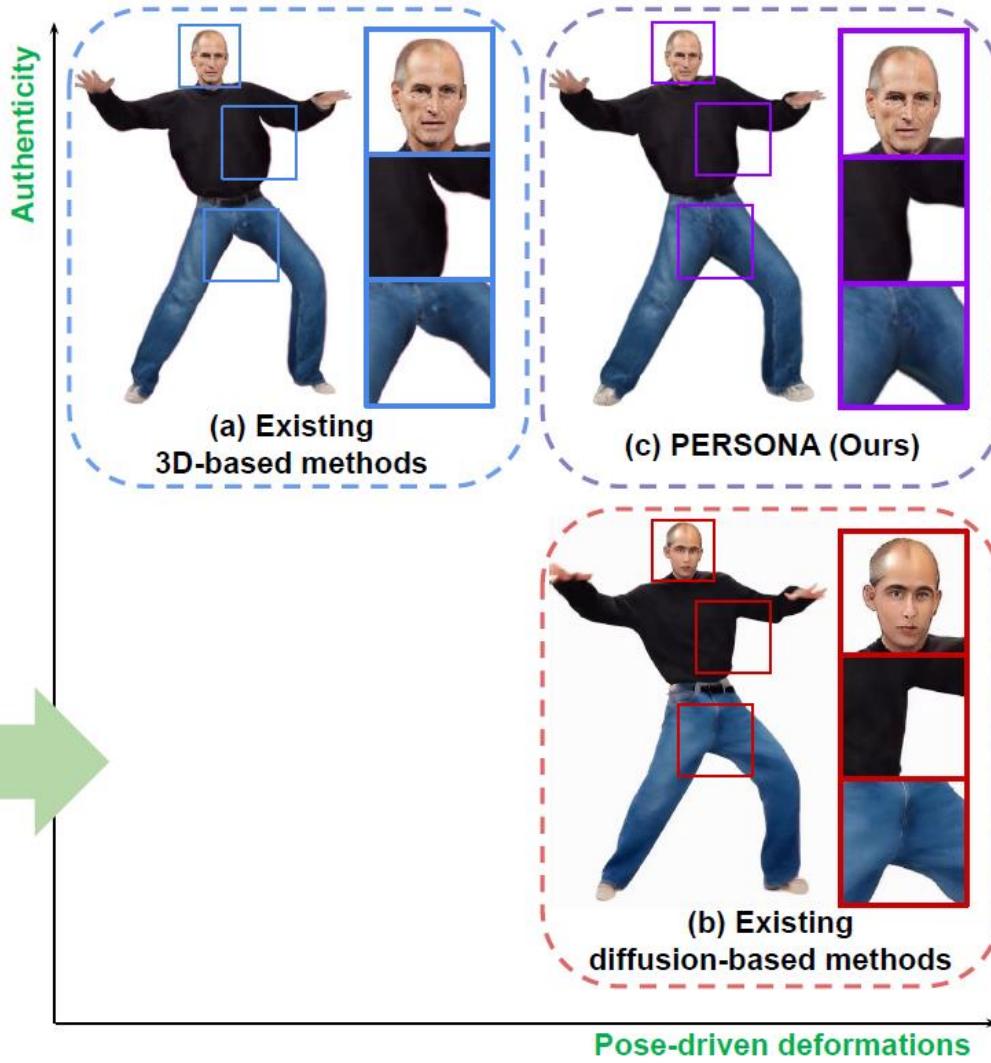
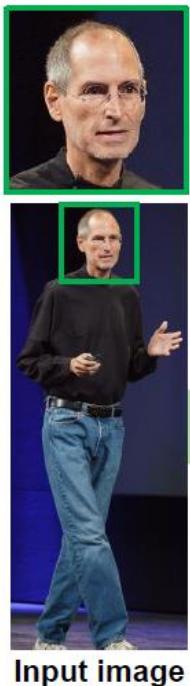
a person in a single reference image + 2D motion -> animated video



[1] Xu, Zhongcong, et al. "Magicanimate: Temporally consistent human image animation using diffusion model." CVPR. 2024.

[2] Hu, Li. "Animate anyone: Consistent and controllable image-to-video synthesis for character animation." CVPR. 2024.

[3] Zhang, Yuang, et al. "Mimicmotion: High-quality human motion video generation with confidence-aware pose guidance." *arXiv preprint arXiv:2406.19680* (2024).



1375

Reference image



Chomp



MimicMotion



StableAnimator



AniGS



PERSONA (Ours)







(a) Input image

(b) Target pose

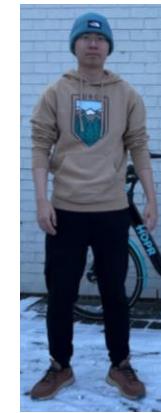
(c) Champ

(d) MimicMotion

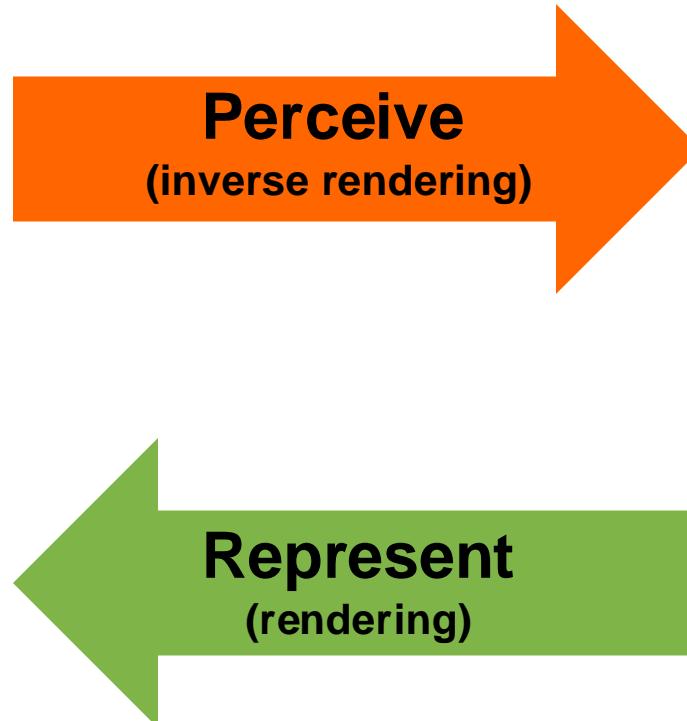
(e) StableAnimator

(f) PERSONA (Ours)

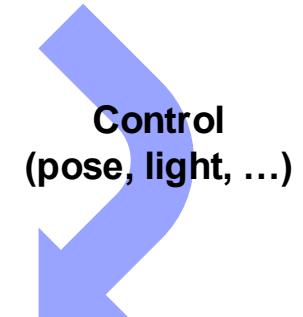
Perceive and Represent Human Images



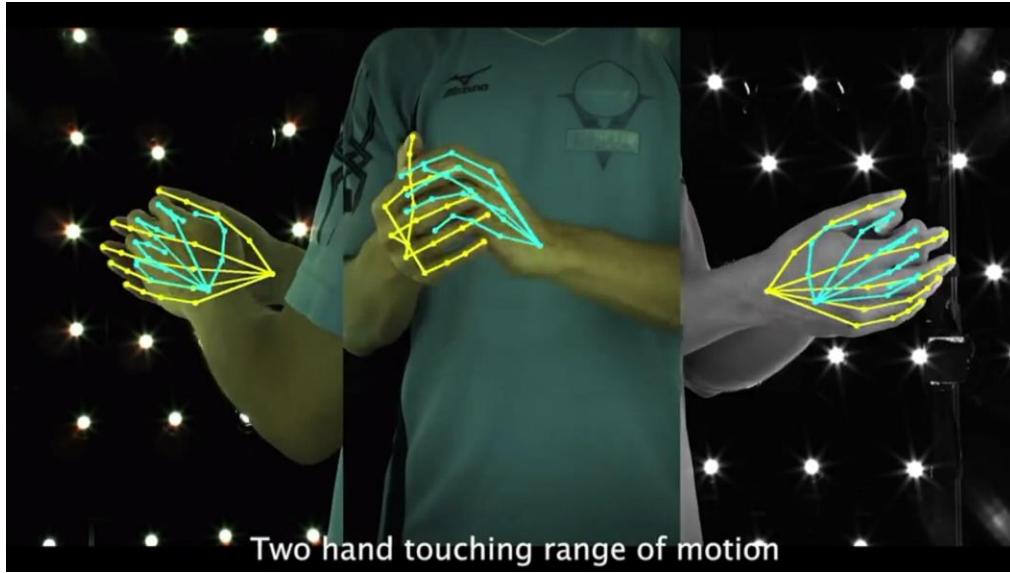
Human images



Human features
(e.g., geometry, texture, light, latent features, ...)



3D Whole-Body **Pose** Estimation from a Single Image



PARTE: Part-Guided Texturing for 3D Human Reconstruction from a Single Image

Work in progress

3D Whole-Body Pose + Surface + Texture Estimation from a Single Image

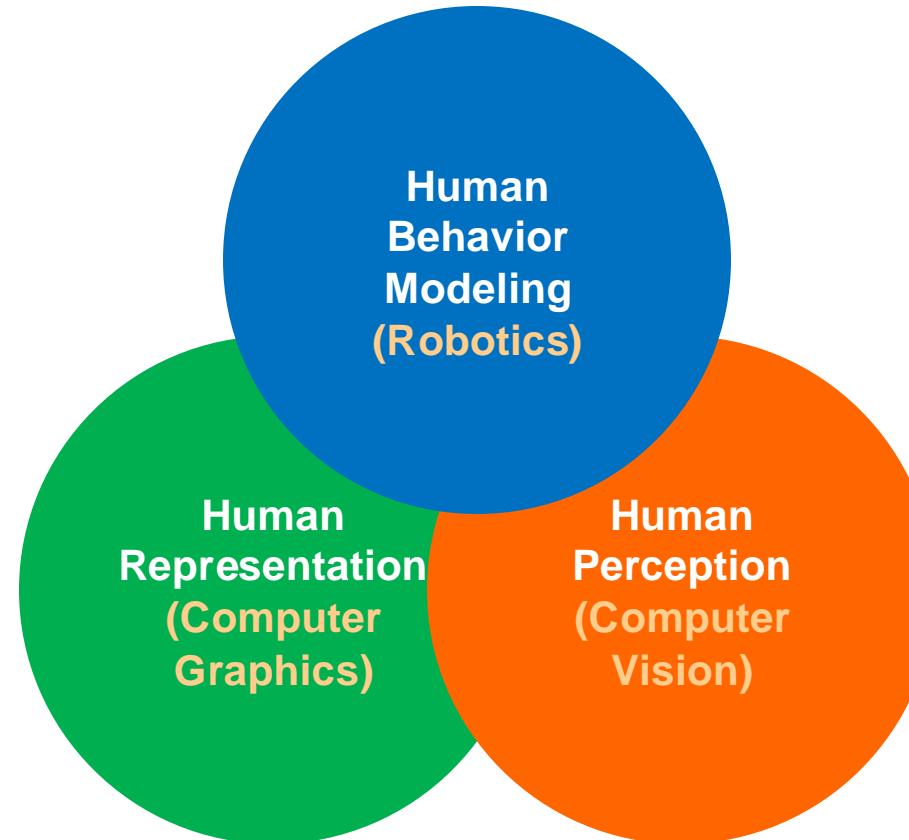


Input image

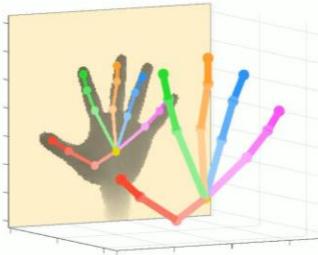
Textured
human surface



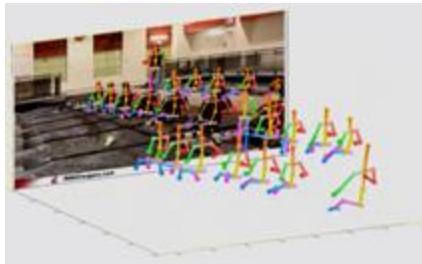
Visual Computing and AI Lab



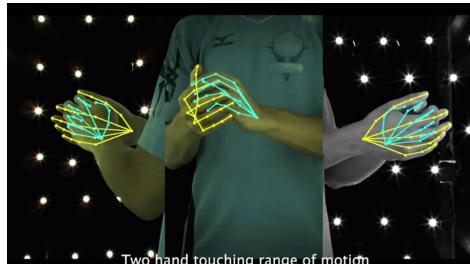
Computer Vision for Human Perception



3D hand pose estimation
(CVPR 2018)



3D multi-person pose estimation
(ICCV 2019)



3D interacting hand pose estimation
(ECCV 2020)



3D body shape estimation
(ECCV 2020)



3D body shape estimation
(CVPR 2021)



Hand-object interaction
(CVPR 2022)



3D whole-body pose estimation
(CVPR 2022)

Computer Vision for Human Representation



(b) DeepHandMesh (ours)



(c) 3D reconstruction

High-fidelity 3D hand geometry model
(ECCV 2020)



Relighted 3D interacting hands
(NeurIPS 2023)



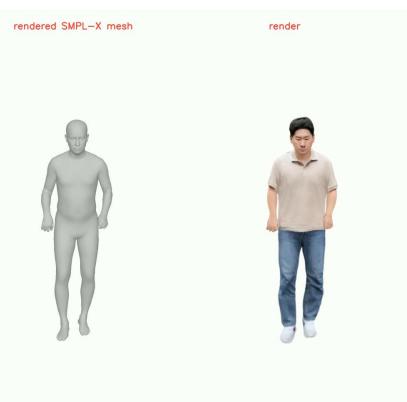
Authentic 3D hand avatar
(CVPR 2024)



Universal relightable hand model
(CVPR 2024)

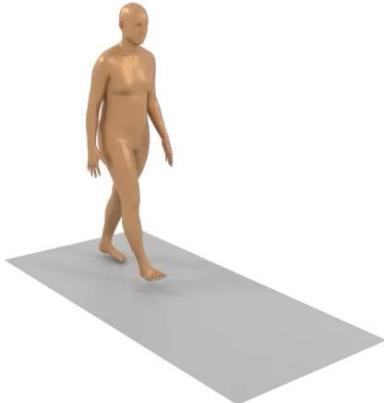


Expressive Whole-Body 3D Gaussian Avatar
(ECCV 2024)

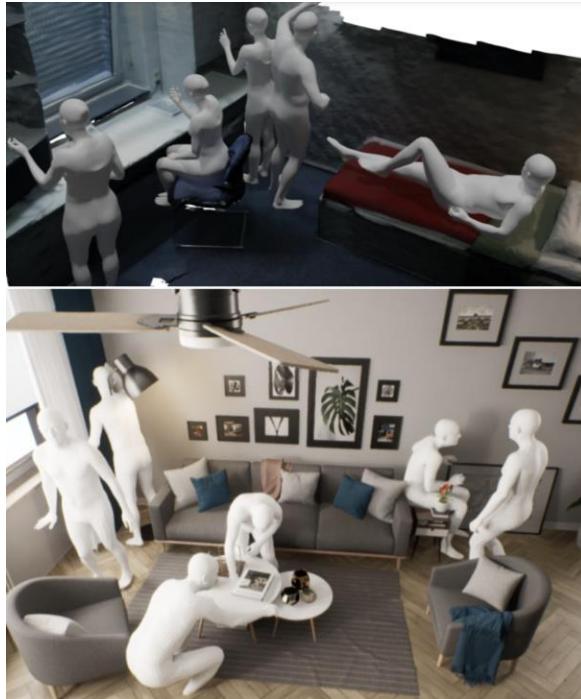


Robotics for Human Behavior Modeling

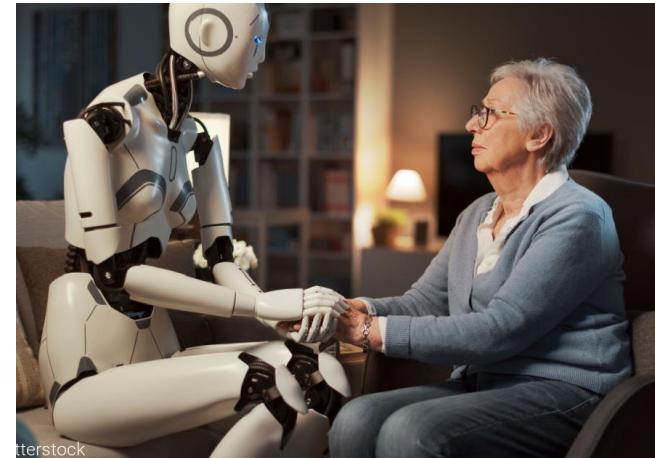
“A person walks forward, bends down to pick something up off the ground.”



Human behavior modeling



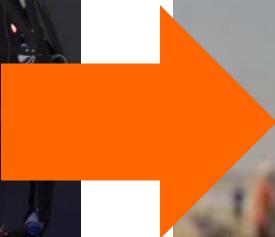
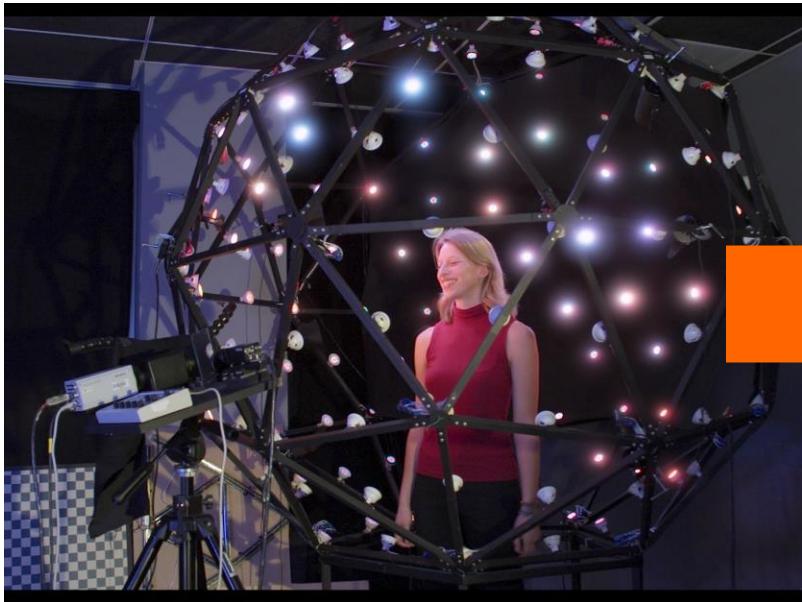
Human-world interaction



Human-robot interaction

Understanding Humans in Casual Environments

- Casual environments have much less available observations
- Need to use **priors to complement limited observations**



Big Wave: Generative AIs

- They are so powerful without any human-specific priors (even without explicit concepts of 3D space)
- **Use priors from them**

OpenAI - SORA



