

Perceiving Humans

Angjoo Kanazawa

CS280

March 31, 2025

Logistics

- Today: 2D/3D Humans
- HW3 up on keypoint detection
- Wednesday: Jitendra
- Next Monday: Learning to predict correspondences
 - → Released papers for you to read in advance on Ed
- Today after class project proposal

Perceiving Humans

From Recognition to Detection to Reconstruction

Why perceive humans?

- Well they are the most important thing



Learning to act from visual observation



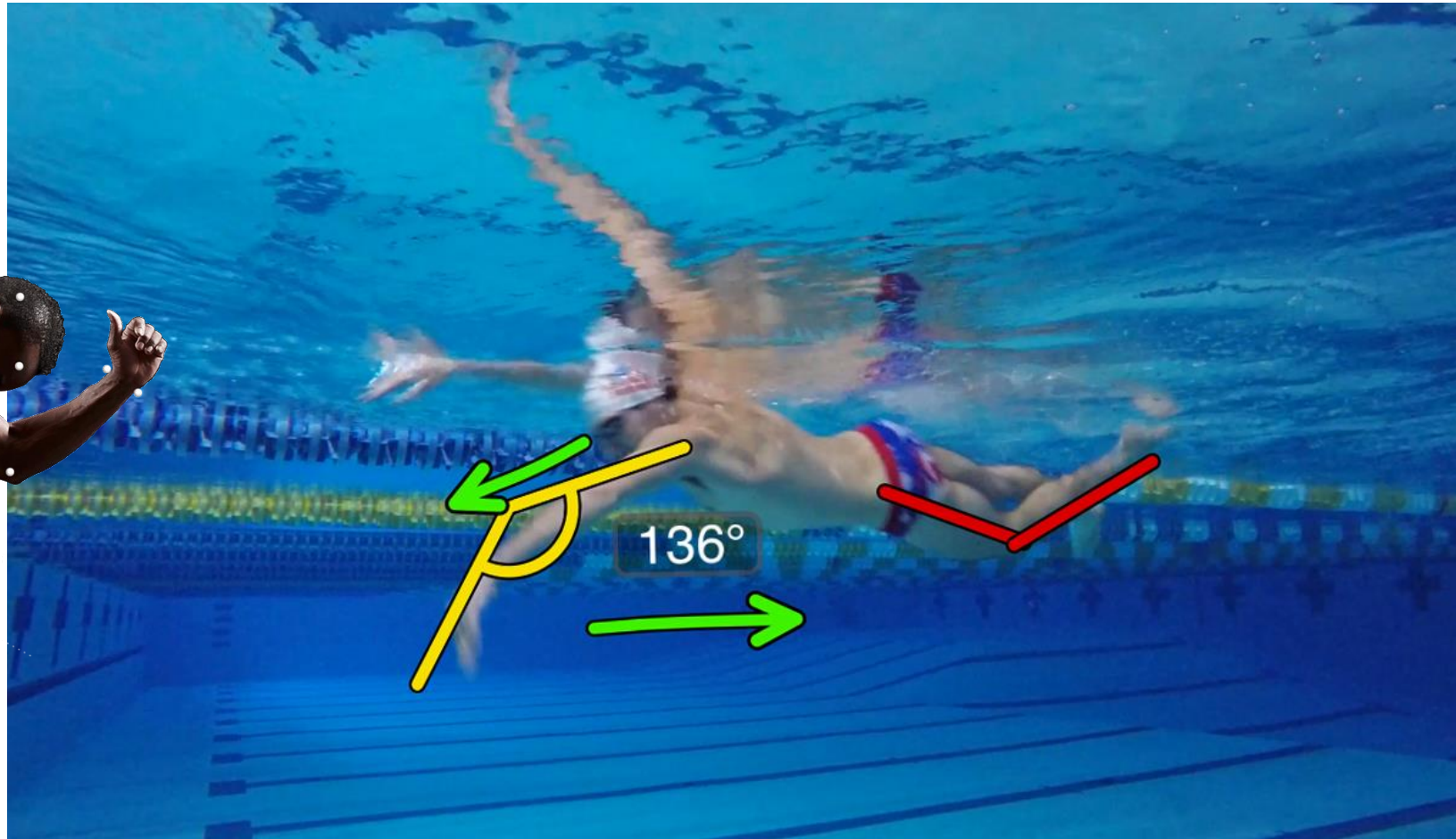
Anticipating human behavior



Sport analysis



OptiTrack



MySwimPro

Medical diagnosis and treatment



Photo Credit: Qualisys

Challenges

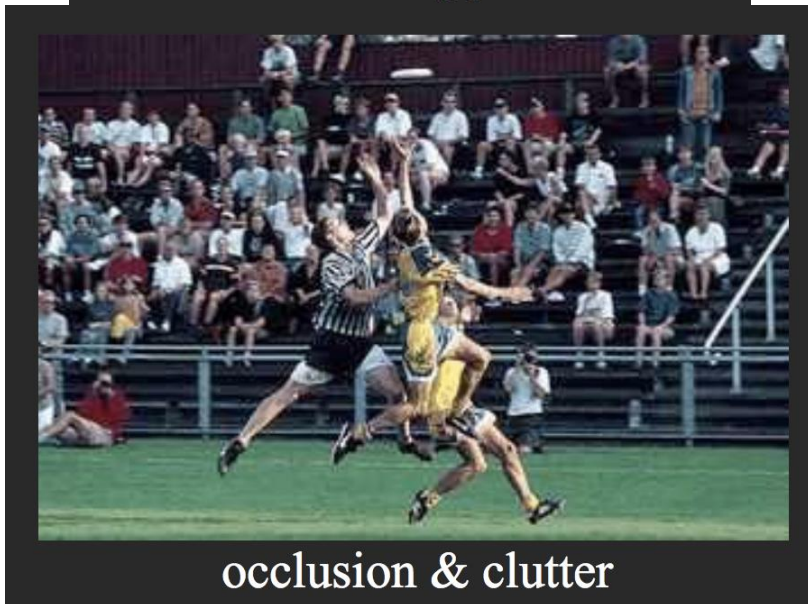
Why is perceiving humans hard?



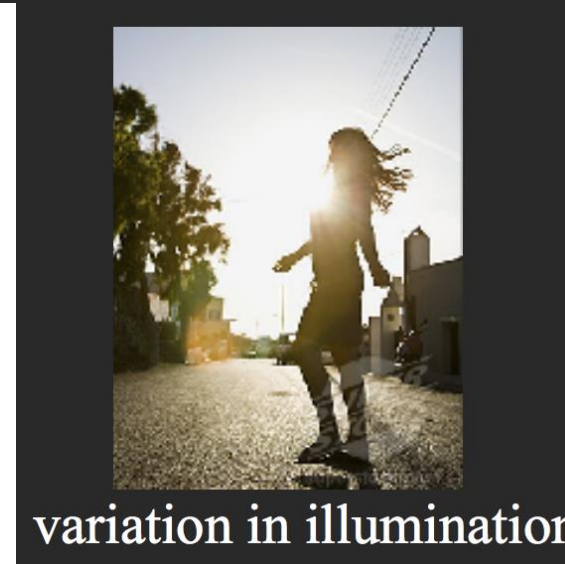
variation in appearance



variation in pose, viewpoint

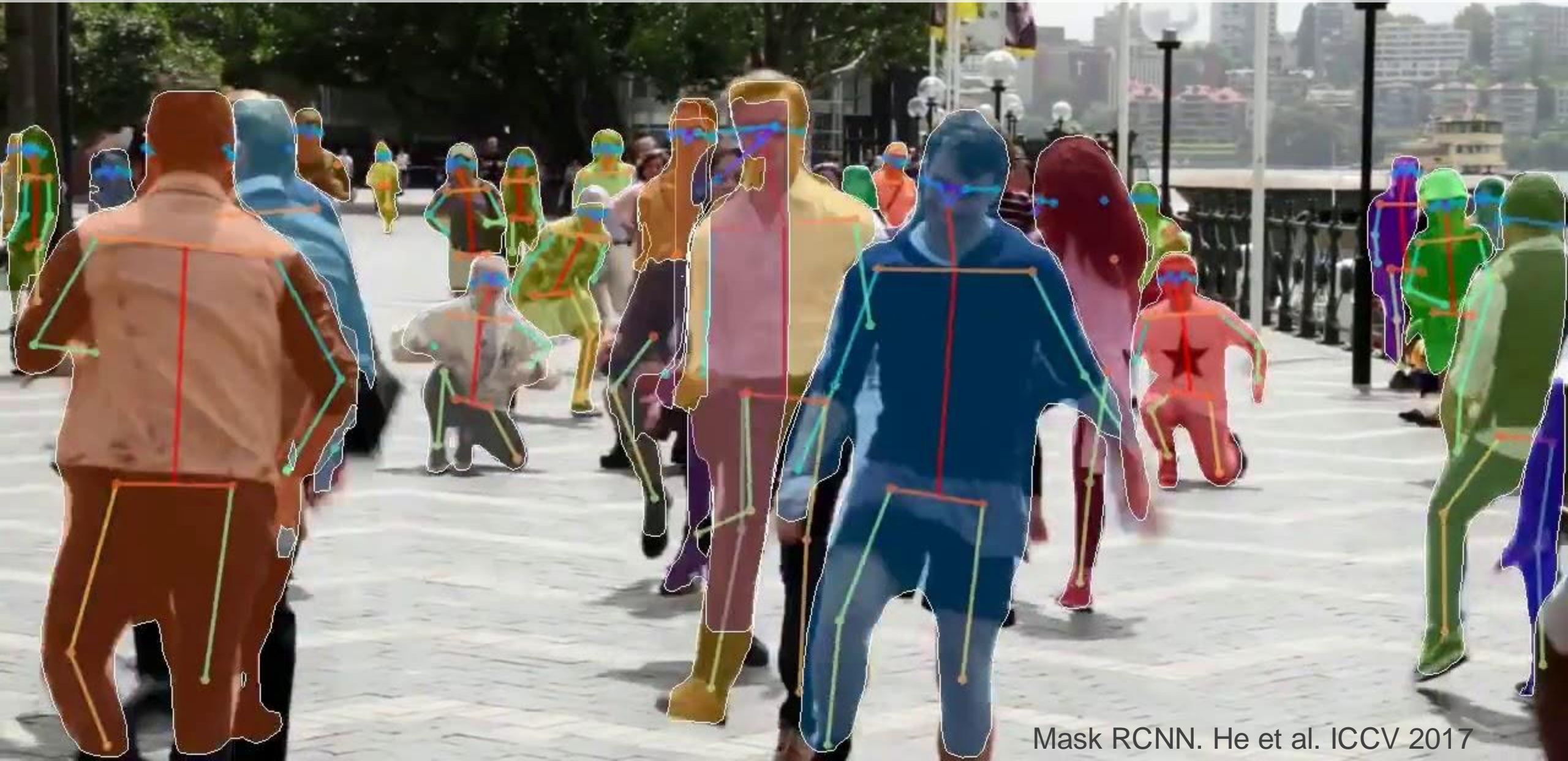


occlusion & clutter

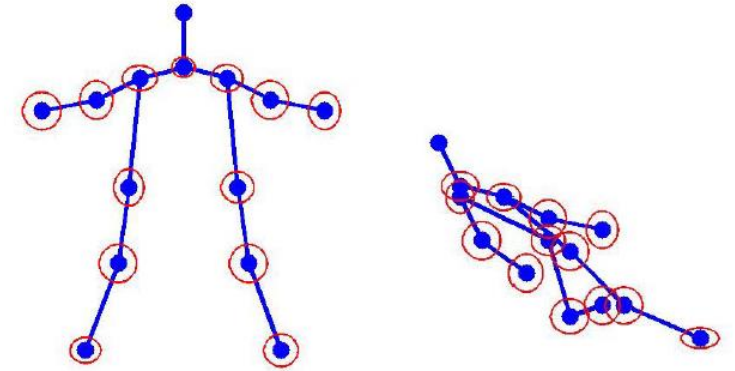
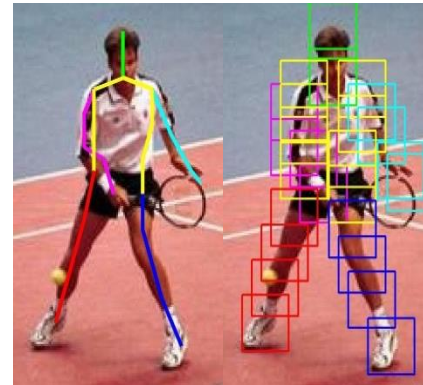
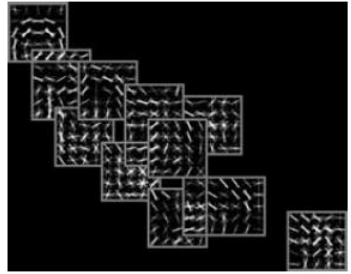
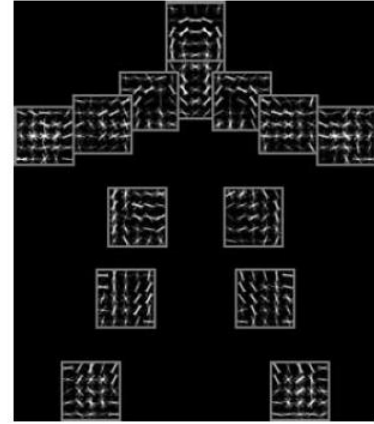
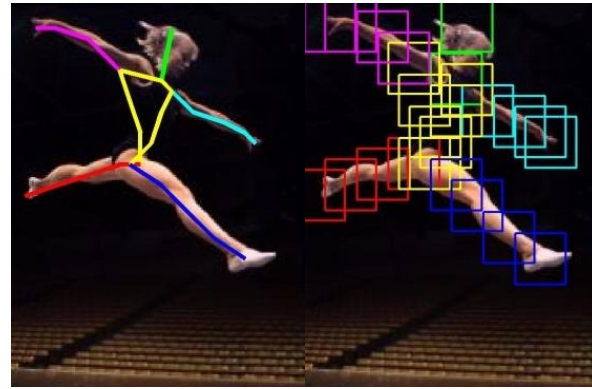
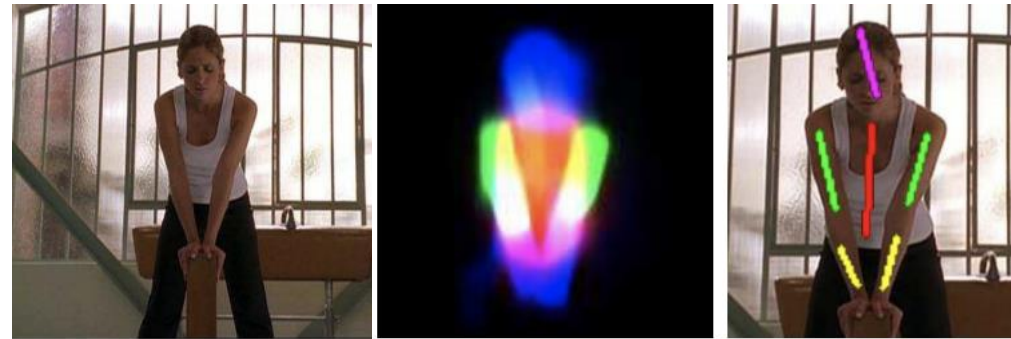


variation in illumination

2D Humans



Parts develop finer into joints & keypoints



[Ferrari, Marín-Jiménez and Zisserman CVPR '08]

Articulated Human Pose Estimation with Flexible Mixtures of Parts
[Yang and Ramanan CVPR '11]

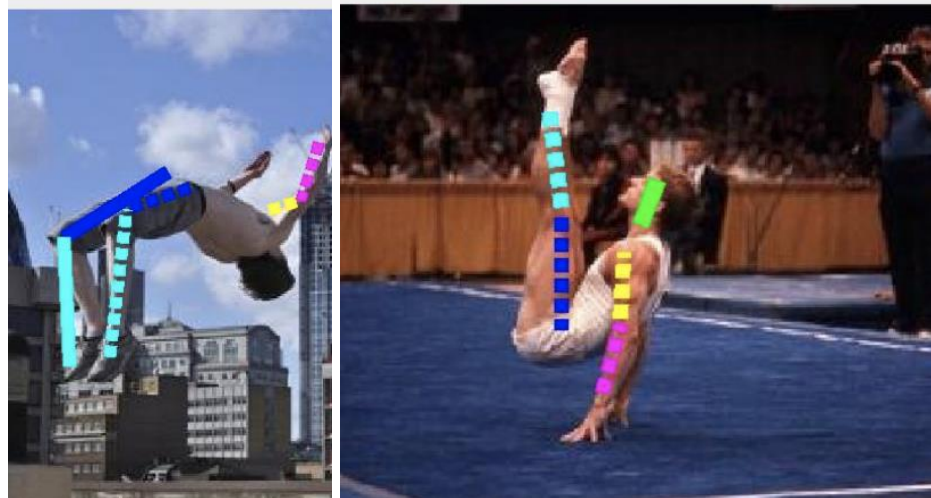
Datasets are introduced

Leeds Sports Pose (**LSP**)

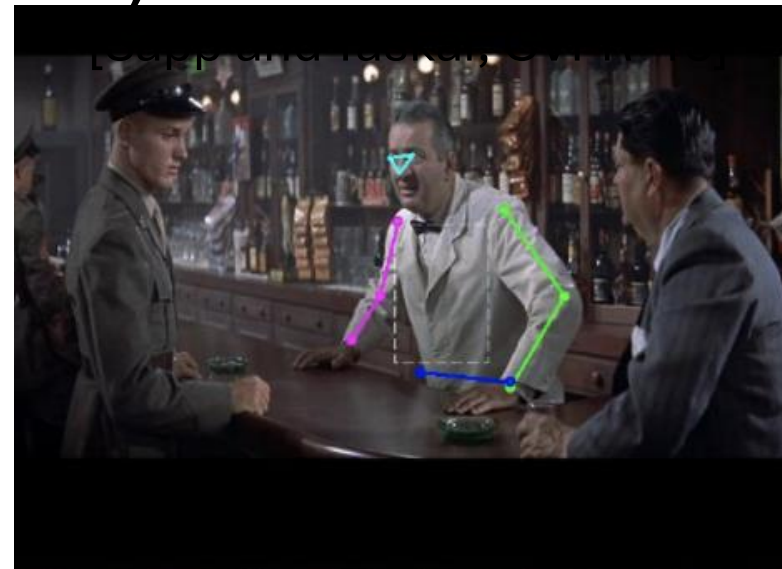
[Johnson and Everingham, CVPR '11]



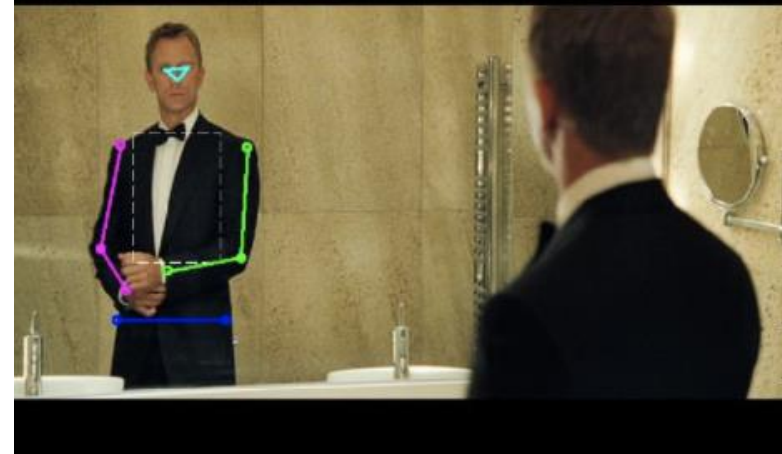
11000 Train
1000 Test



Frames Labeled in Cinema (**FLIC**)



4000 Train
1000 Test

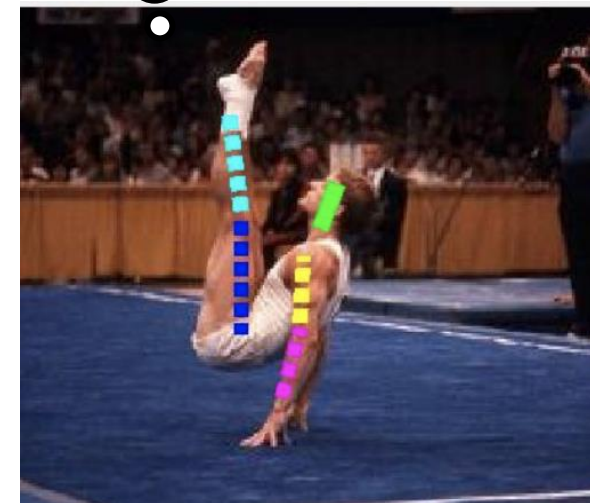
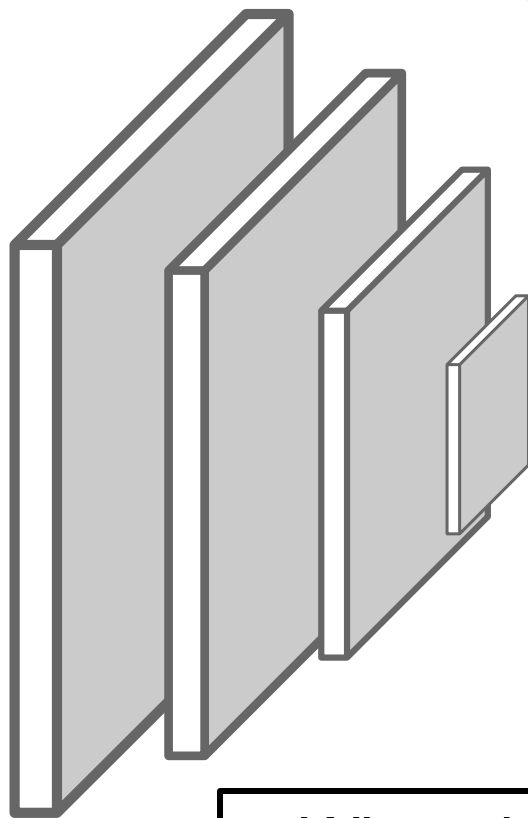


Deep Learning Era

How do we represent/parametrize 2D human pose???



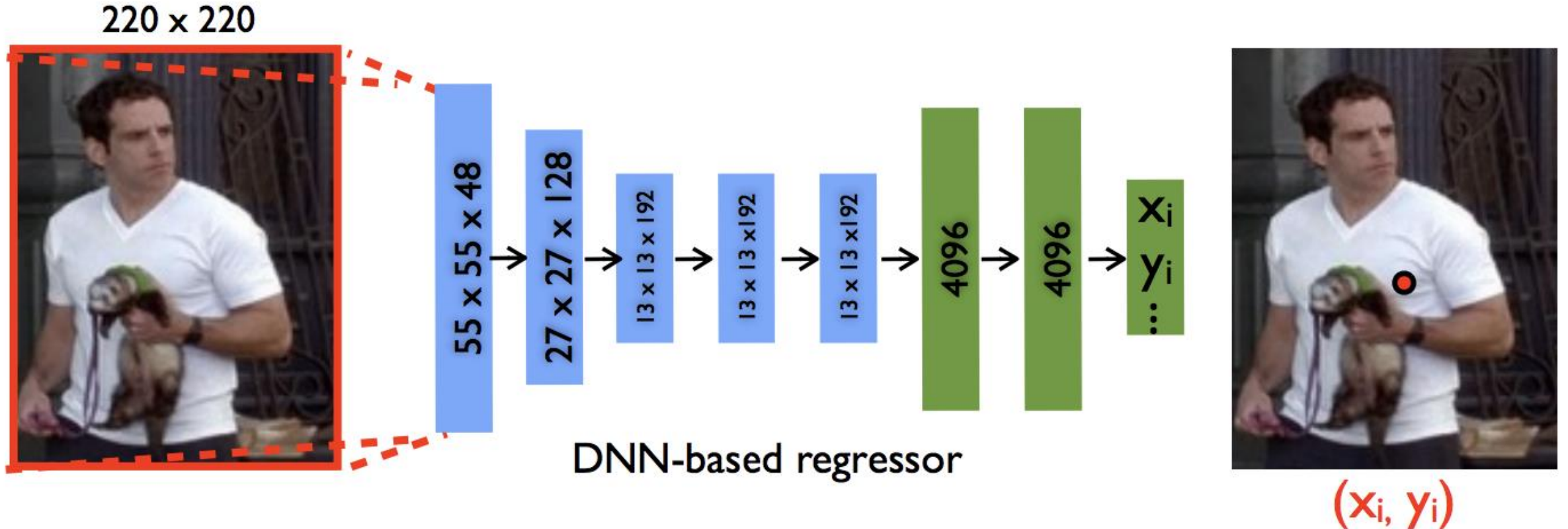
Input



Output

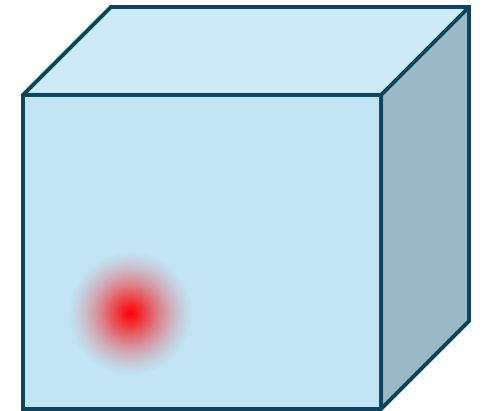
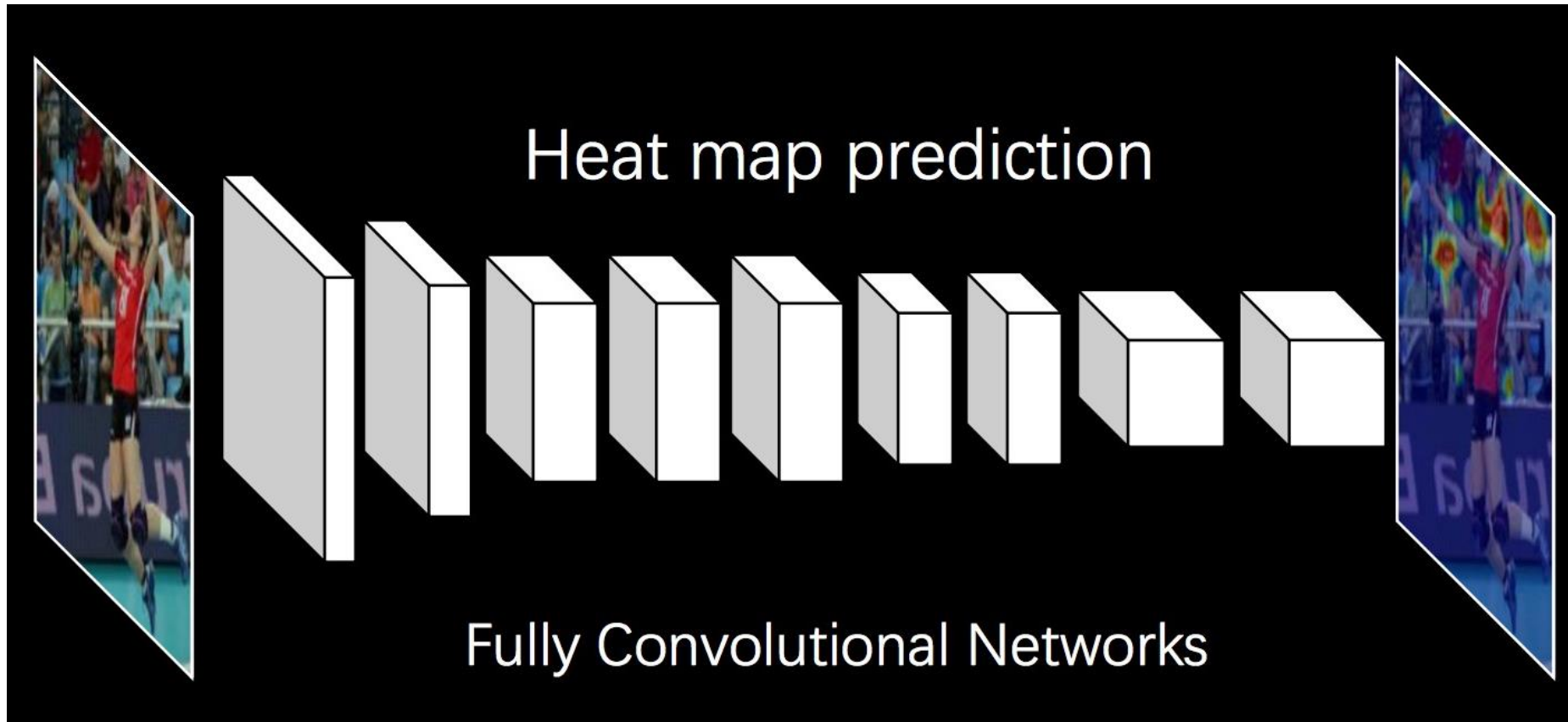
What should the network output?

Predicting keypoints



DeepPose: Human Pose Estimation via Deep Neural Networks
[Toshev and Szegedy 2014]

Predict heat maps



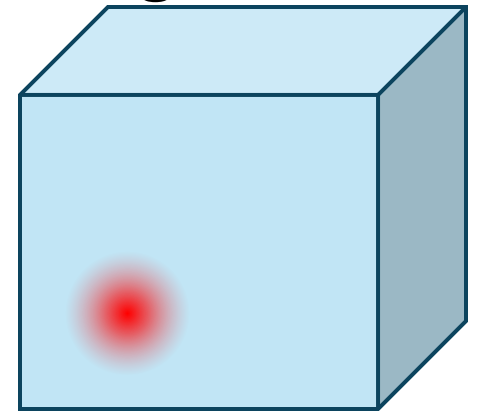
Target: $K+1 \times H \times W$
Gaussian around
 (x,y) for k -th
keypoint in the k -th
channel

$K+1$ for K parts + background

L2 Training Loss

- L2 loss on the target heatmap (peaky gaussian around the gt keypoint)

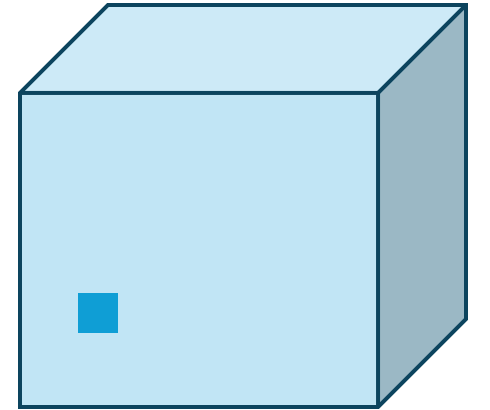
$$L = \sum_{k=1}^{K+1} \sum_{(x,y)} ||b^k(x,y) - b_*^k(x,y)||$$



Target “belief map” :
K+1 x H x W
Gaussian around
(x,y) for k-th
keypoint in the k-th
channel

Log Loss Training Loss

- Log loss (or cross entropy loss) on the target heatmap probabilities
- The target must also sum to 1
- Mask RCNN just uses 1 at the target, 0 everywhere else.
- Experiment



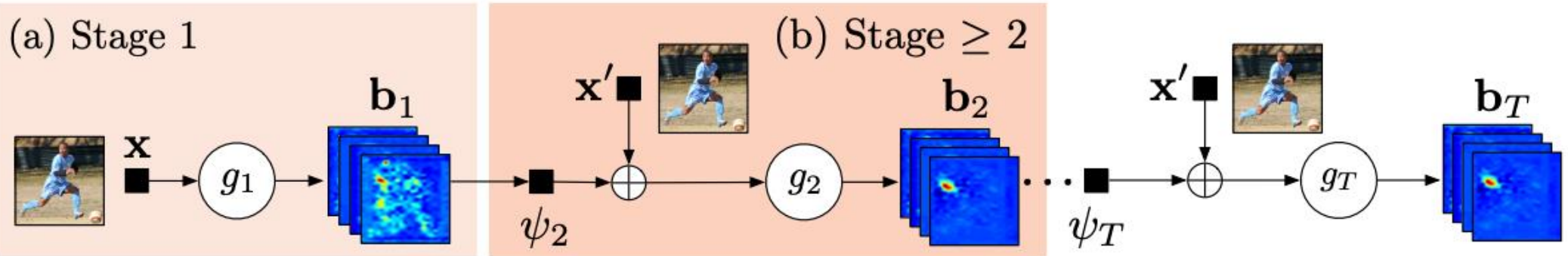
Target “belief map” :
 $K+1 \times H \times W$
1 at Ground truth
location (x,y) for k -th
keypoint in the k -th
channel

Convolutional Pose Machines

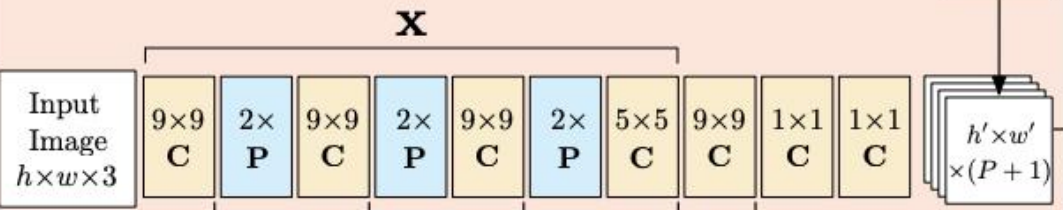
Base architecture for OpenPose

Convolutional Pose Machines (T -stage)

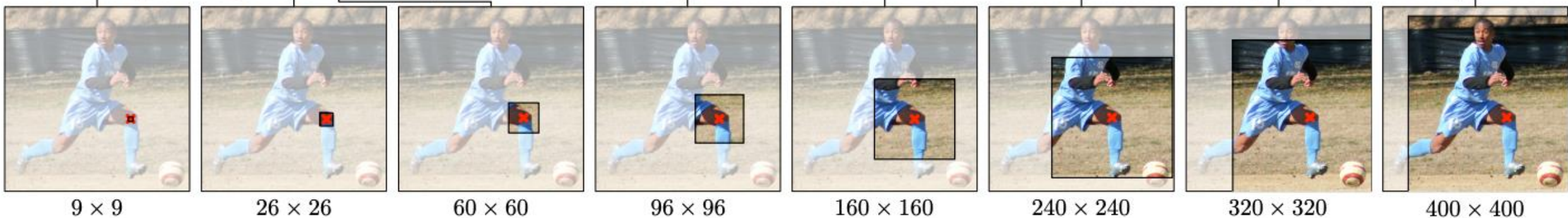
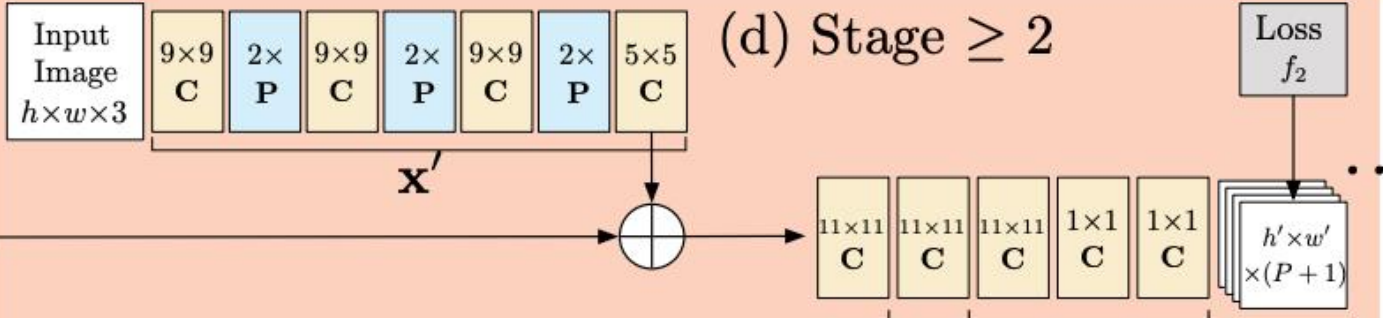
P Pooling
C Convolution



(c) Stage 1

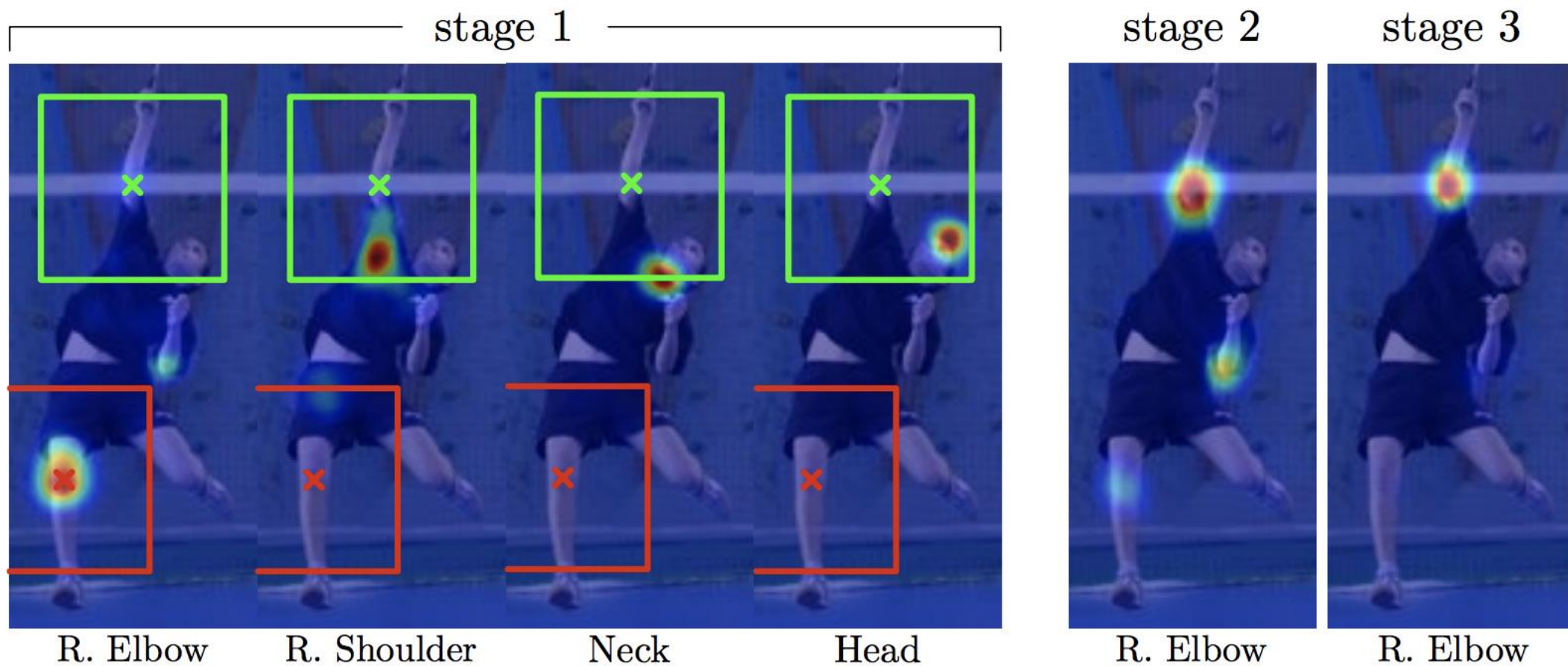


(d) Stage ≥ 2

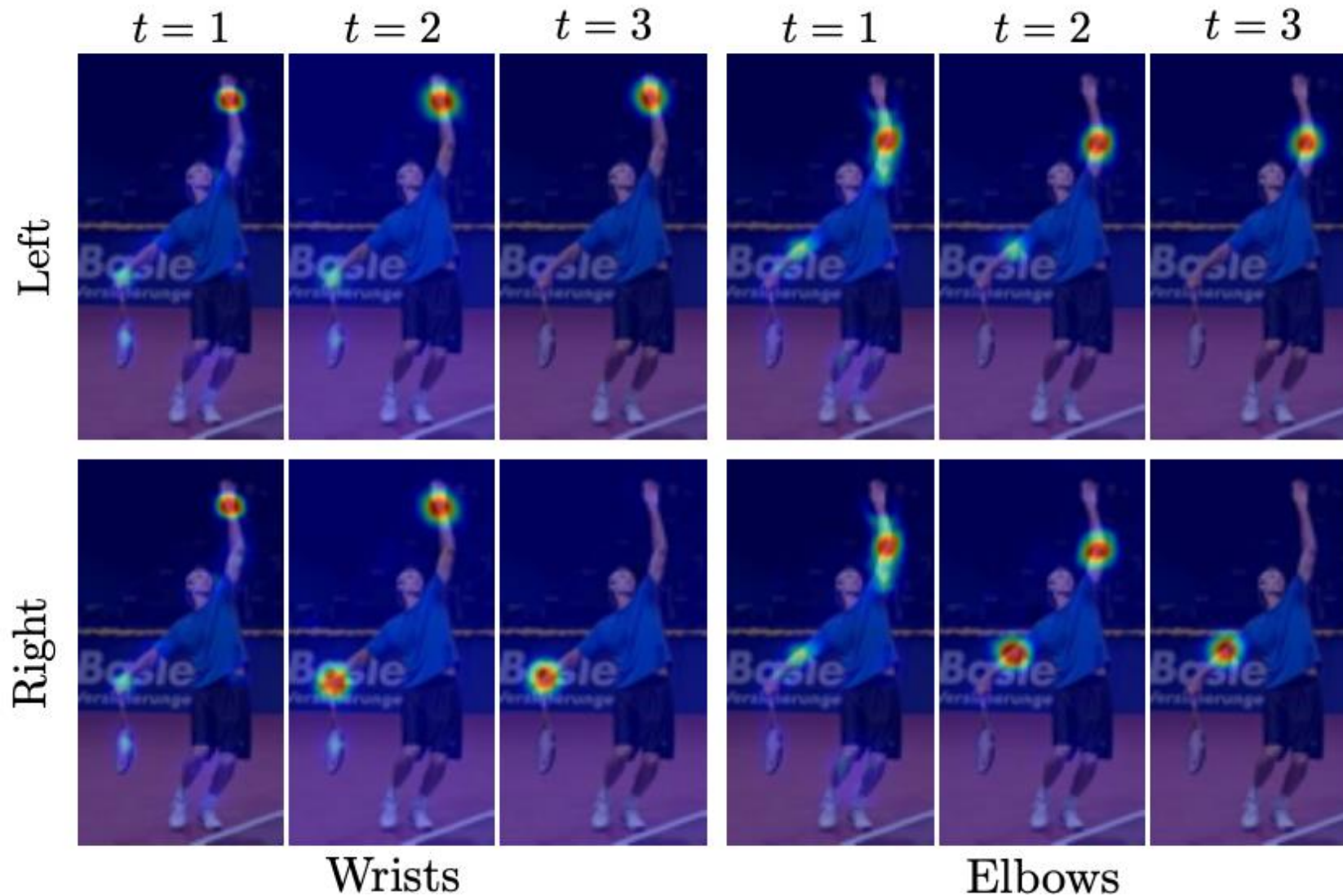


(e) Effective Receptive Field

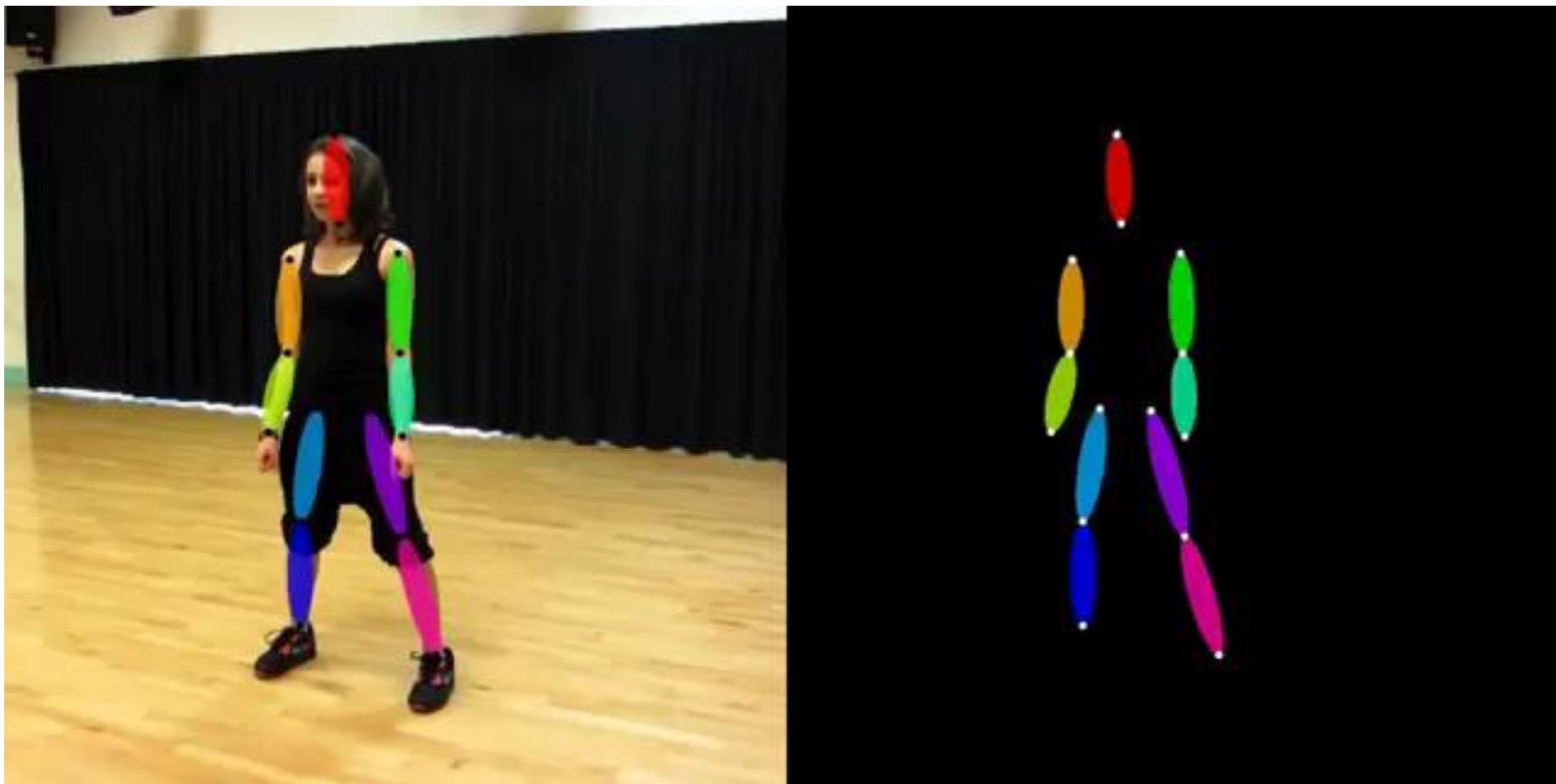
Convolutional Pose Machines



Convolutional Pose Machines



Results



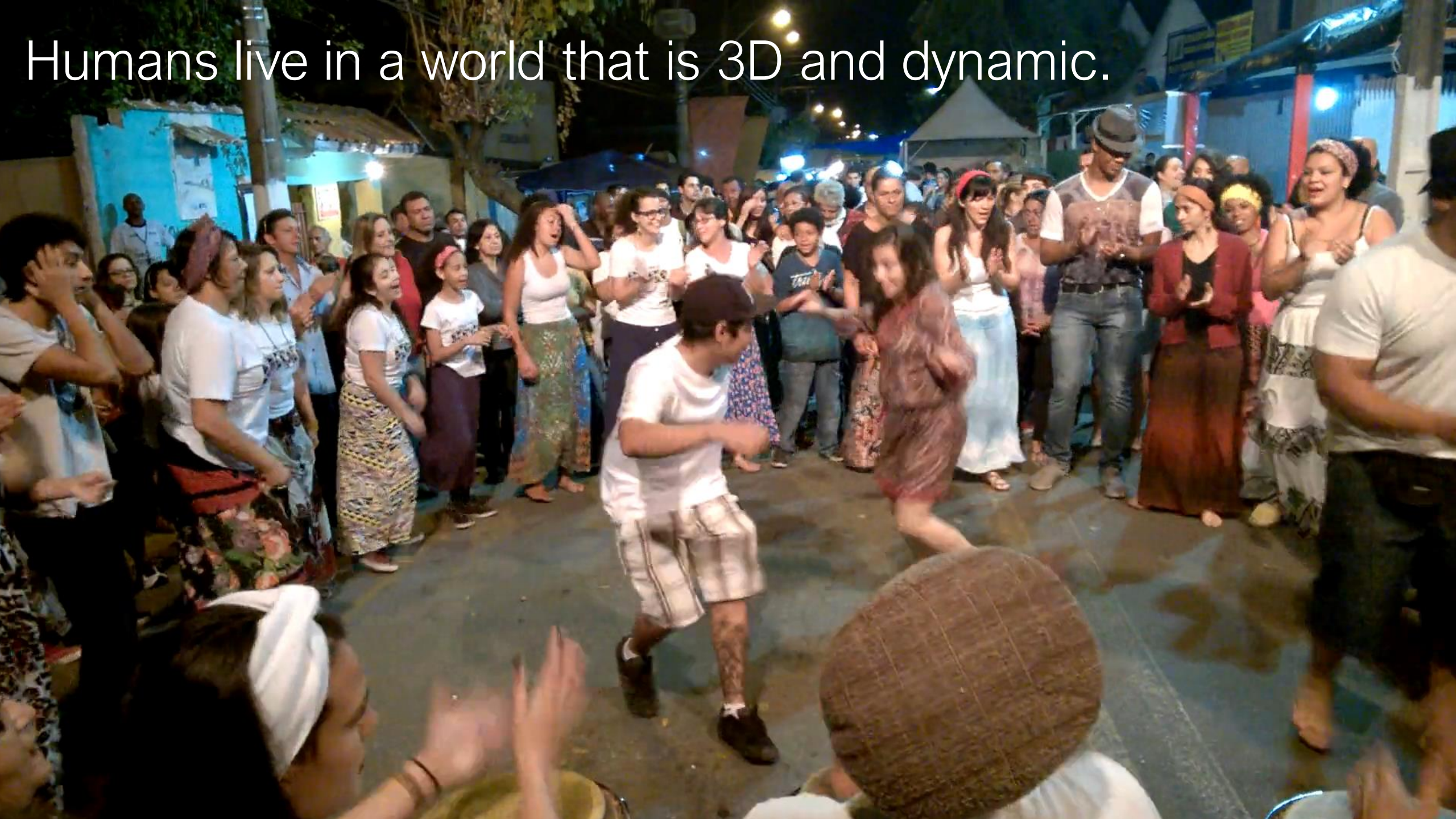
OpenPose

Great opensource tool, builds on convolutional pose machine architecture, adapted to multiple people



Are we done?

Humans live in a world that is 3D and dynamic.

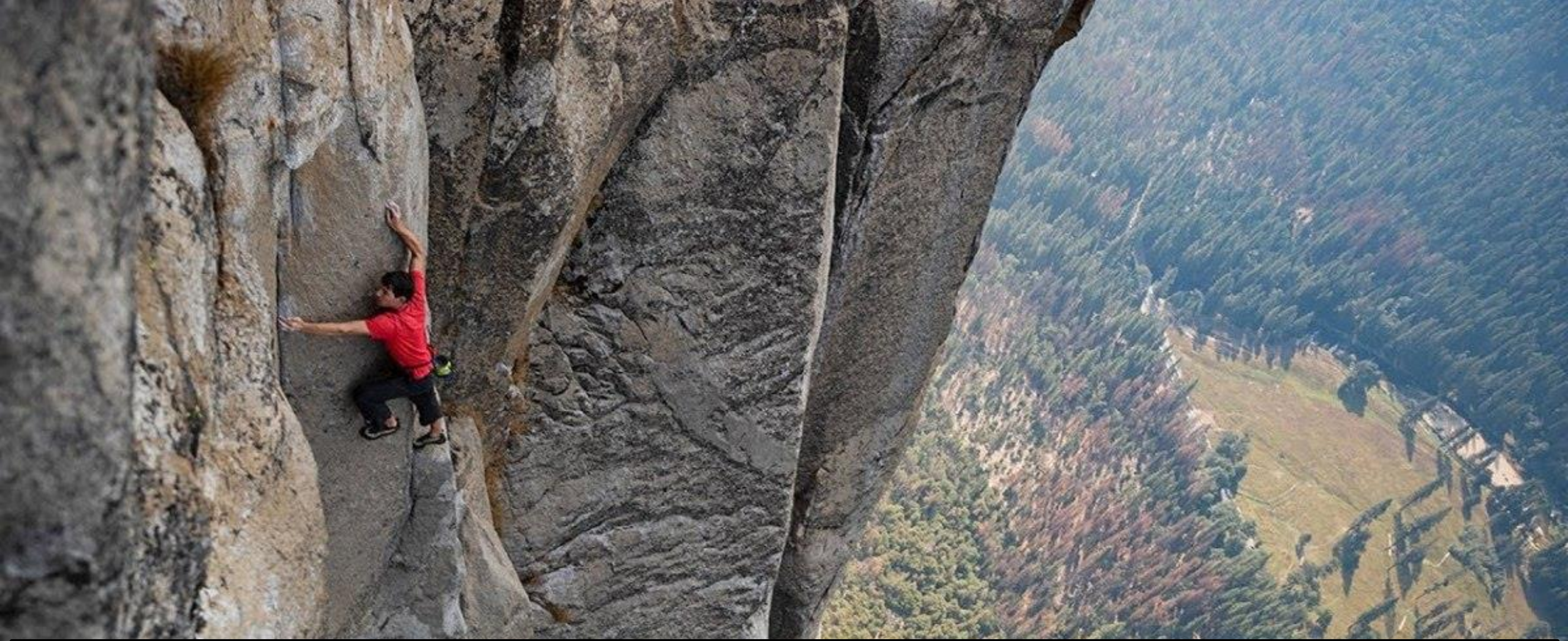


Today's Non-rigid 3D Solution: Motion Capture



The world is so much more than greenscreen!



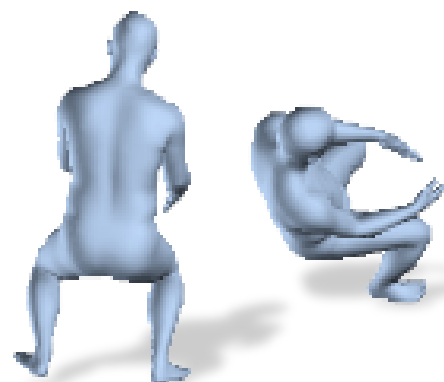


Goal: 3D perception from images “*in-the-wild*”

Single-View 3D Human Mesh Recovery



In everyday photos



Or from Video



Learning to act from visual observation



From video...



Video

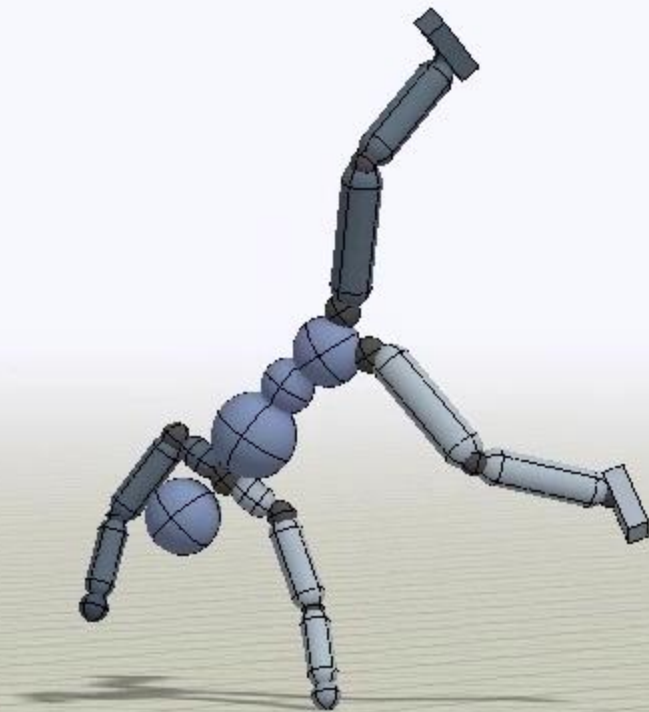


Recovered 3D Body



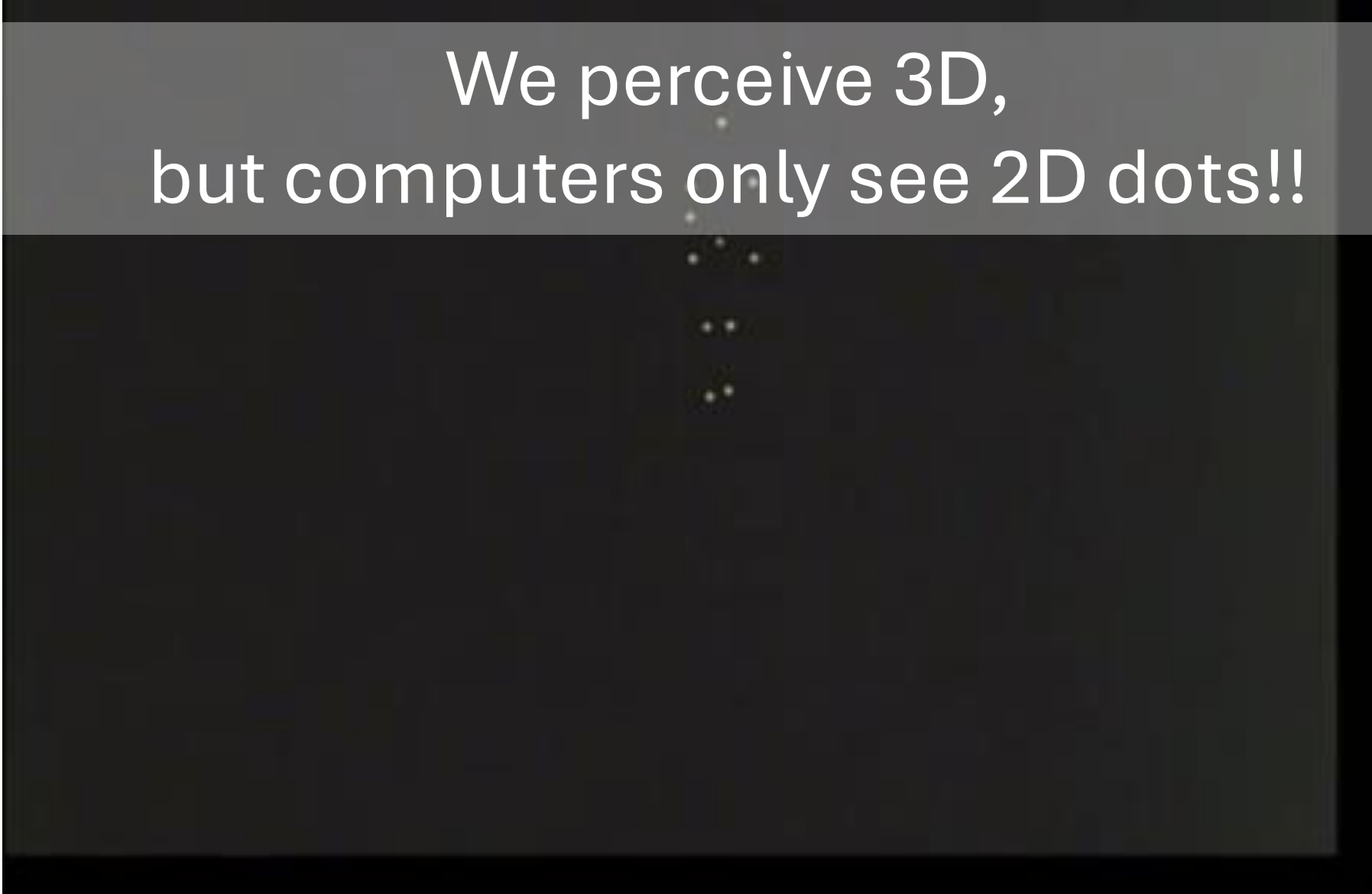
Policy

Animate Virtual Characters



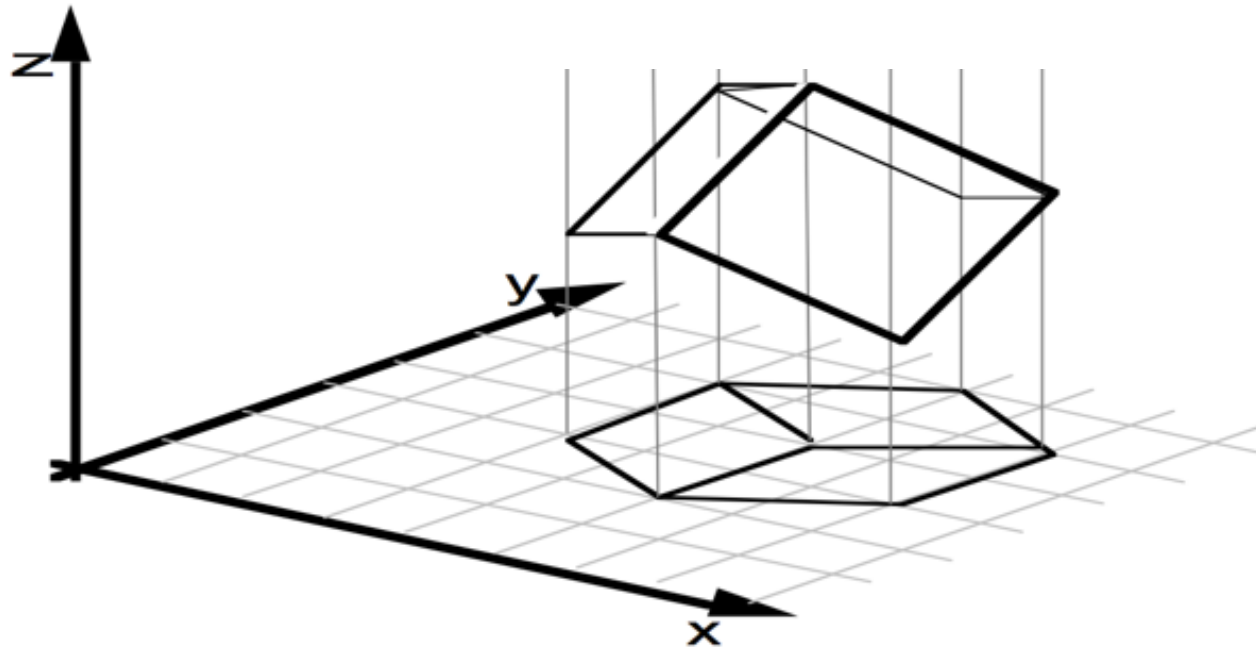
Human 3D perception

We perceive 3D,
but computers only see 2D dots!!

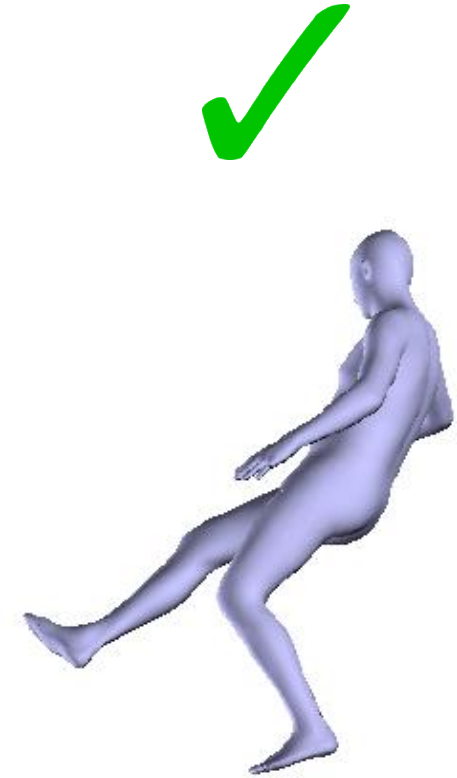
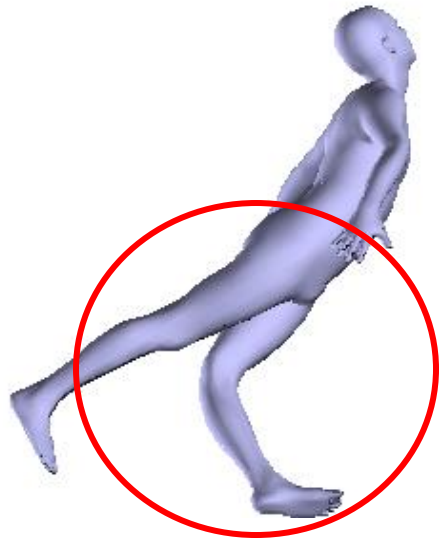


Johansson
experiment,
James Maas, 1971

3D from 2D is inherently under-constrained



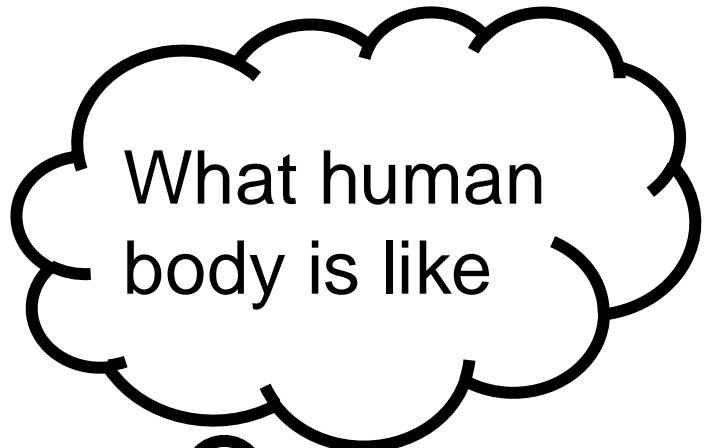
How do we resolve this?



How do we resolve this?

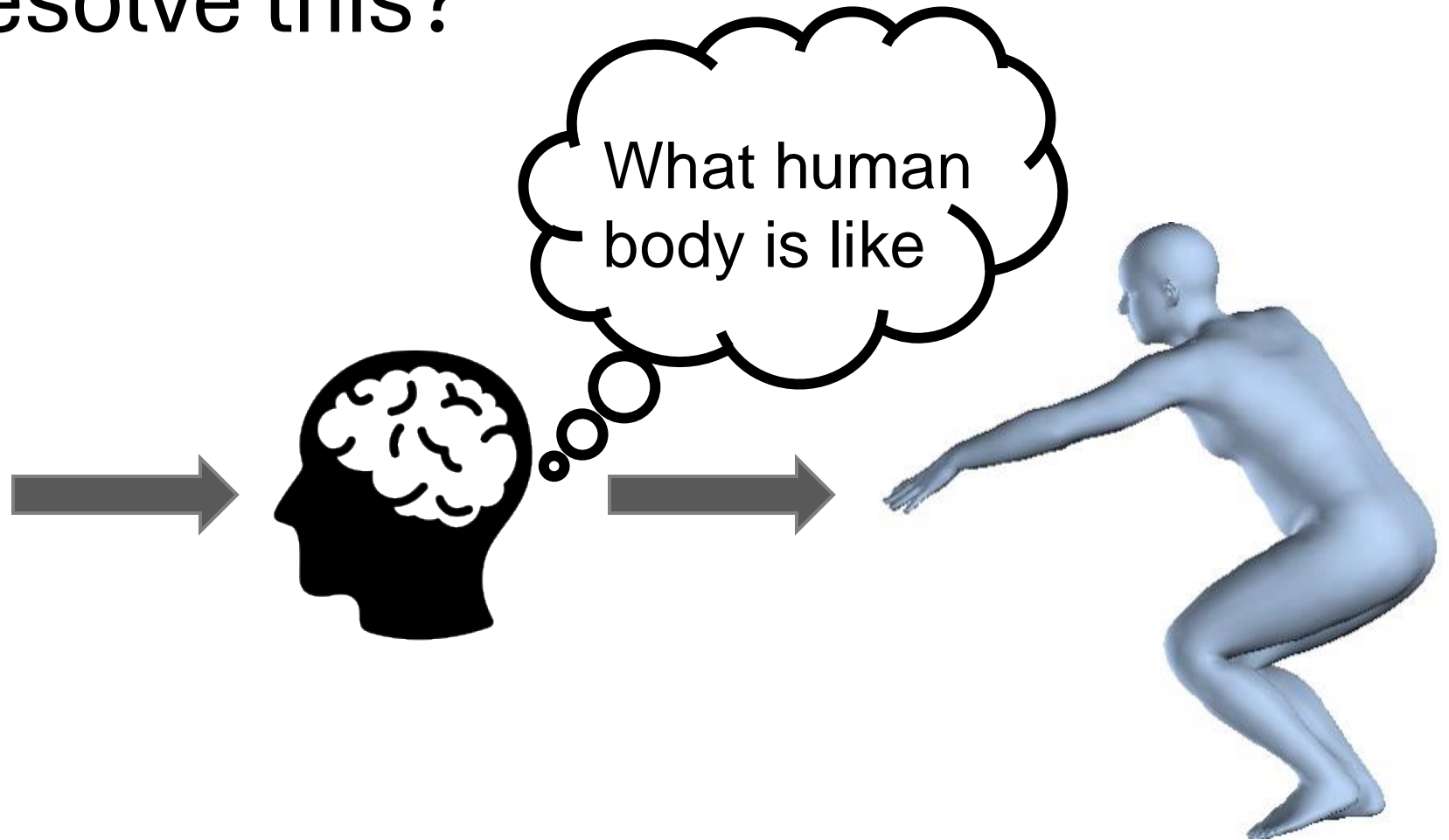


“Remembrance of Objects Past”

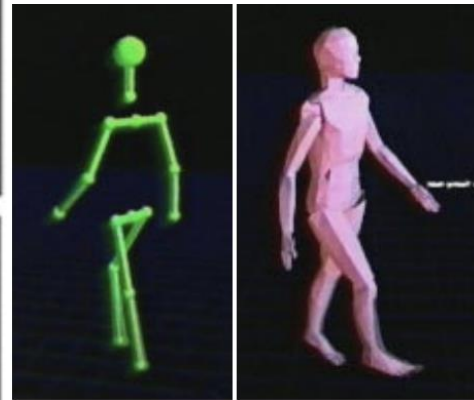


What human
body is like

How do we resolve this?



Bregler and Malik CVPR 1998

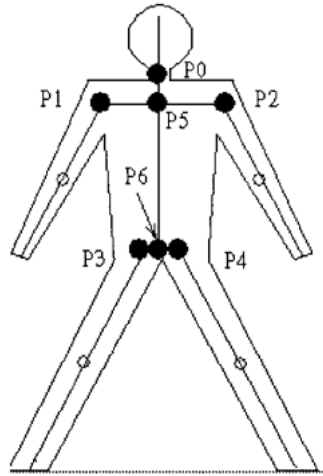


Tracking based:

1. Initialize 3D model in first frame
2. Track parts in next frames via Lukas-Kanade, over joint angles

More stable with 2 views

And many more model-based methods



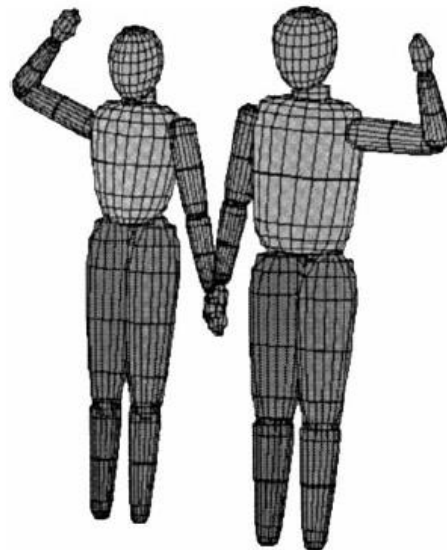
[Leung and Yang '95]



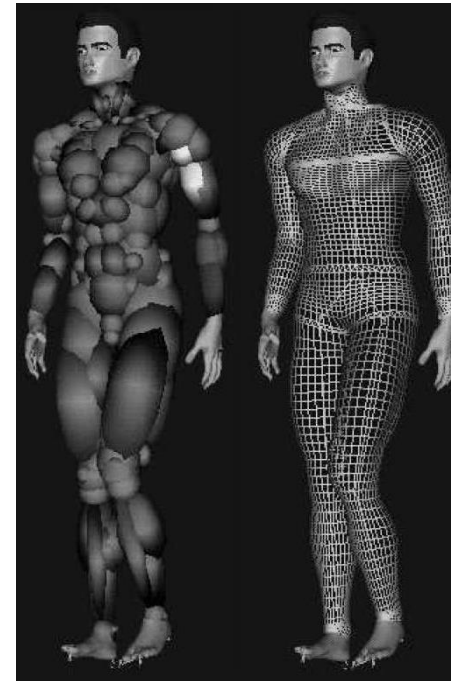
[Terzopoulos and Metaxas '93]



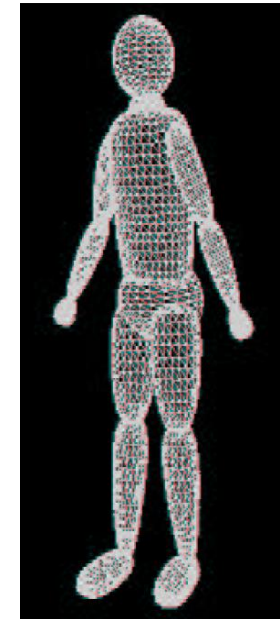
[Kakadiaris and Metaxas '00]



[Gavrilla, '96]



[Plänkers and Fua '01]



[Sminchisescu and Triggs '03]

3D Humans from known 2D joints



Reconstruction of articulated objects from point correspondences
in single uncalibrated image
[CJ Taylor CVIU 2000]

Same issue as single-view 3D reconstruction

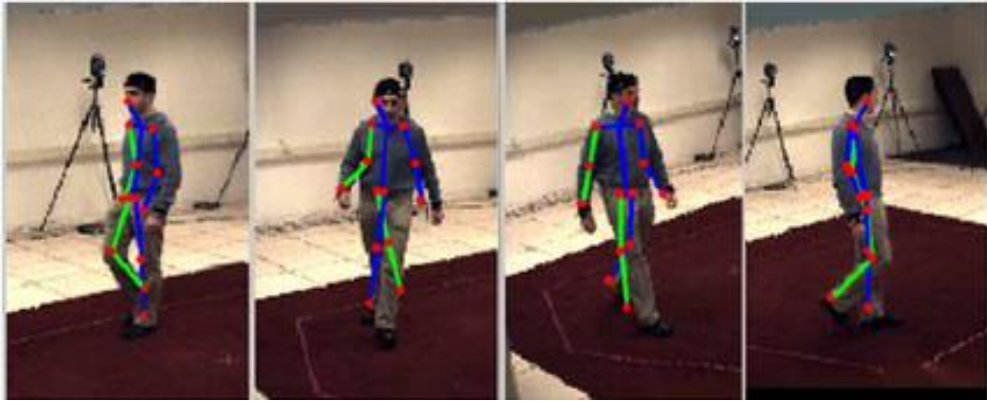
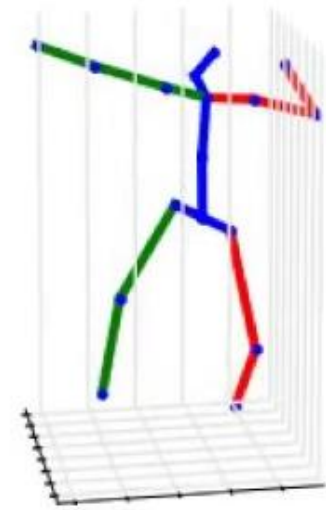
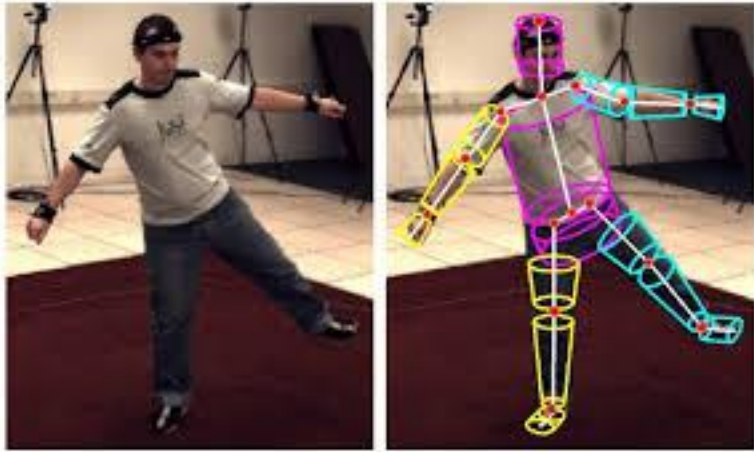


(a)



You need priors!!
Here: Known ratio of
limb length

Datasets

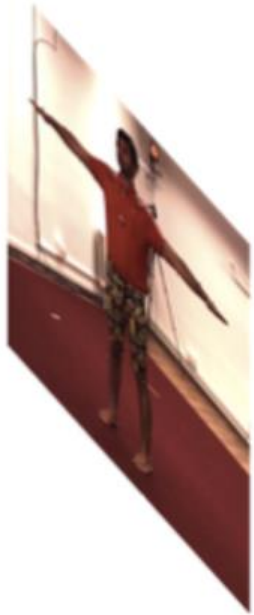


HumanEva [Sigal et al. IJCV 2010]

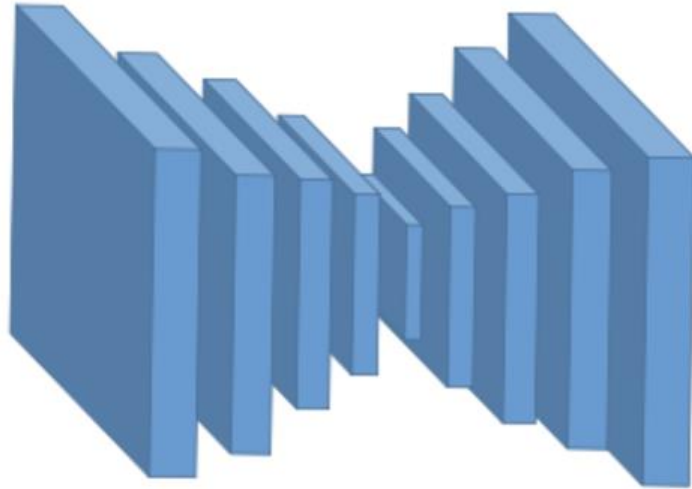
Human3.6M [Ionescu et al. 2014]

Deep Learning based approaches

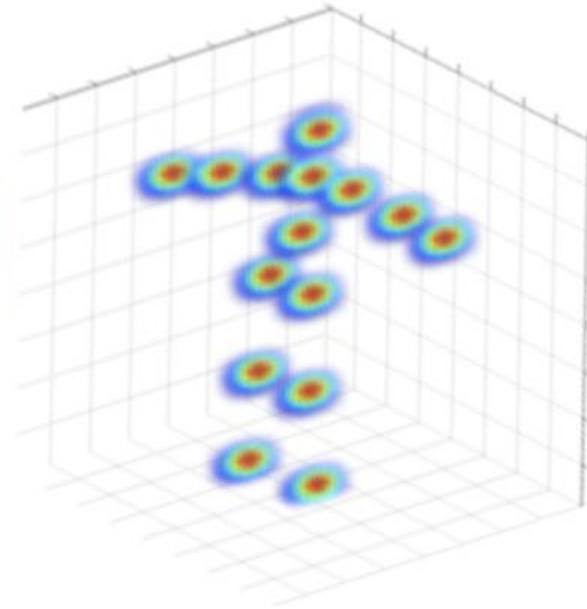
Lots of activities + progress made in this area after datasets



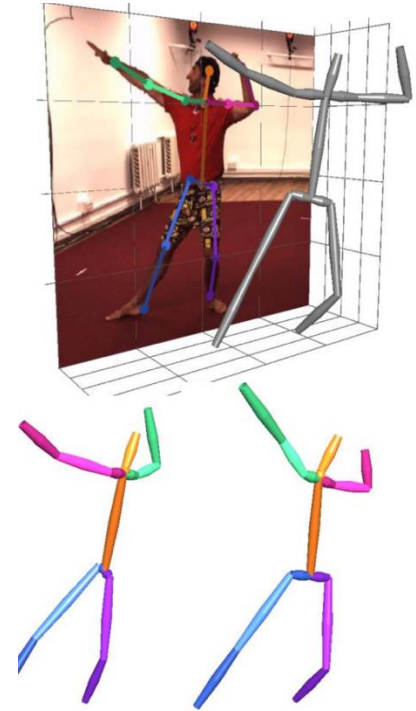
Image



ConvNet



Volumetric Output

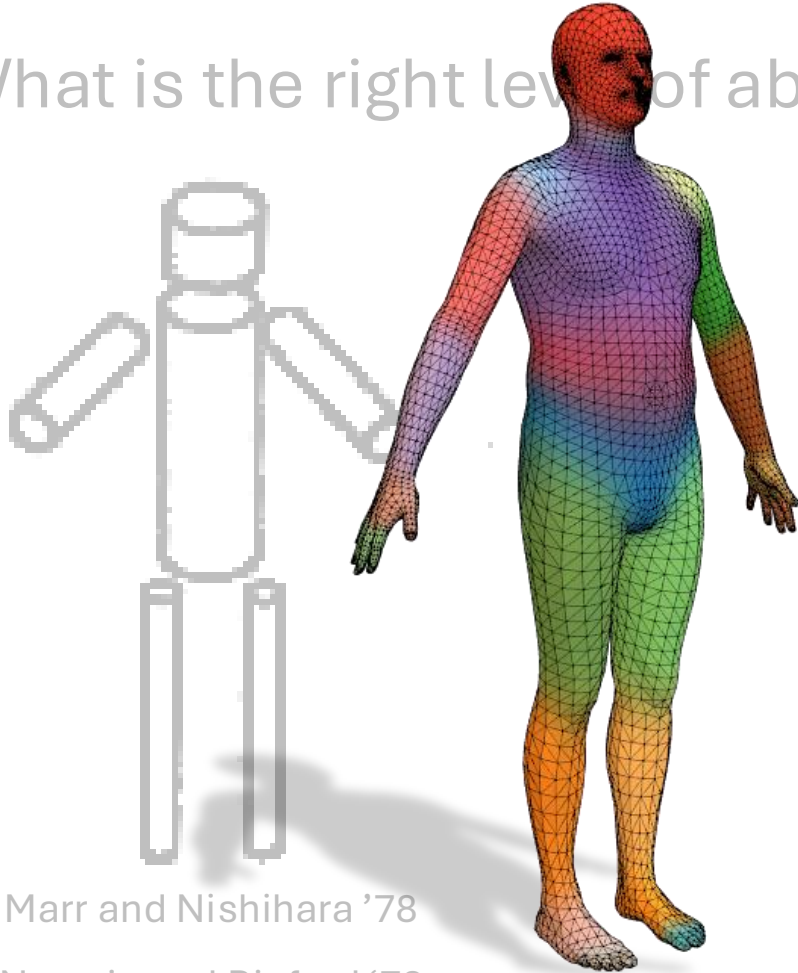


[Pavlakos et al. CVPR'17, Sun et al ICCV'17, Vnect Mehta et al SIGGRAPH '17 ...]

**What about the 3D
representation??
Are we all just stick figures?**

To do more than joints, we need to discuss how to model human bodies

What is the right level of abstraction?



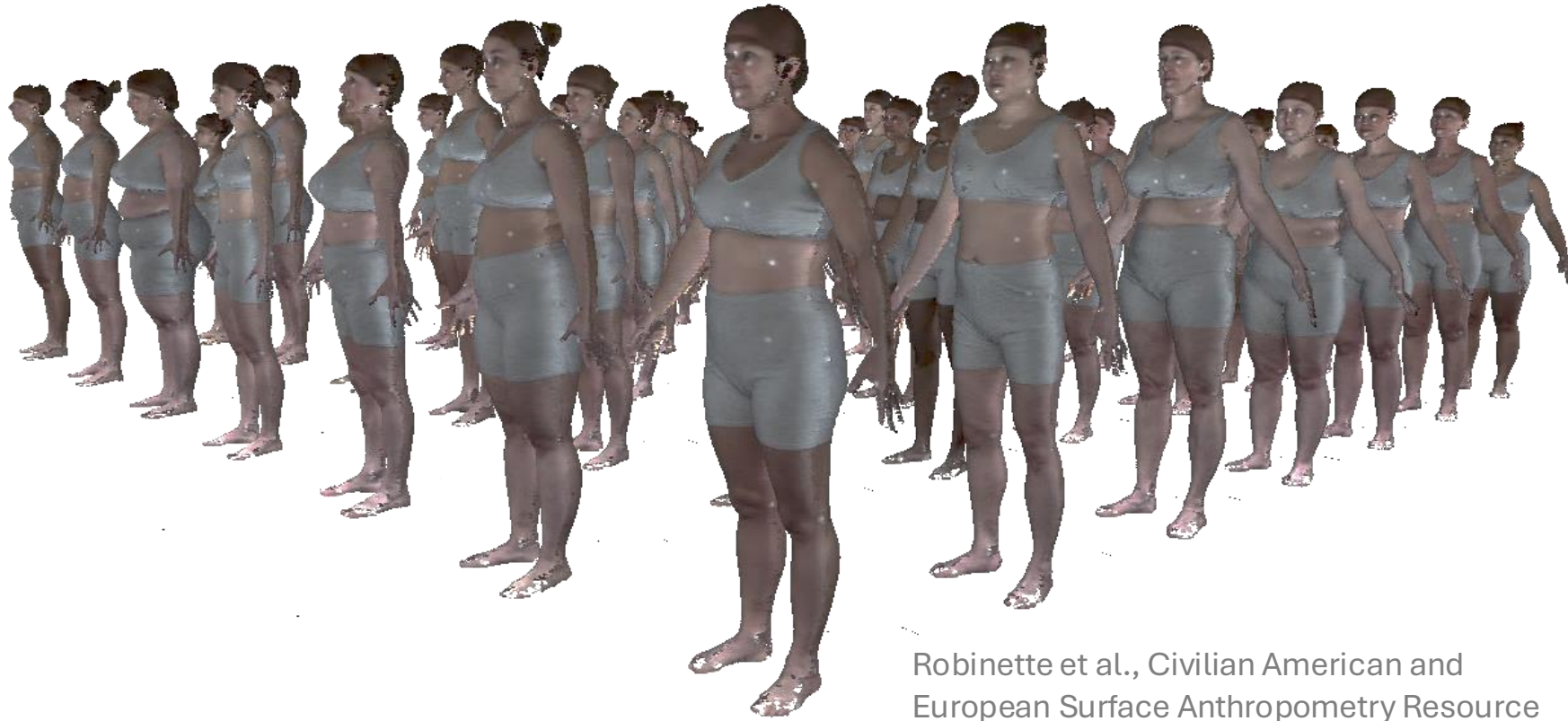
Marr and Nishihara '78

Nevatia and Binford '73

Practical and popular answer has been to model what you can see → the surface

Lee, Sifakis, Terzopoulos, ToG'09

Humans are special



Robinette et al., Civilian American and European Surface Anthropometry Resource (CAESAR) 2002.

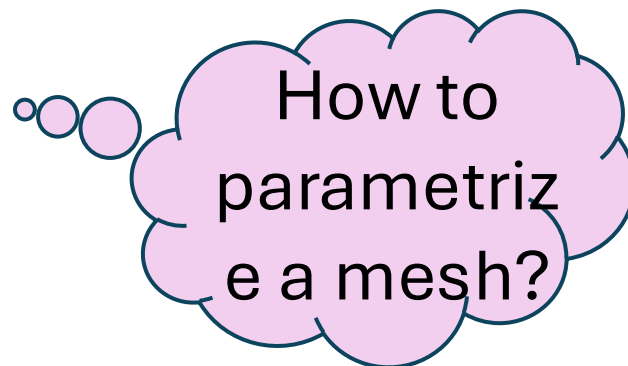
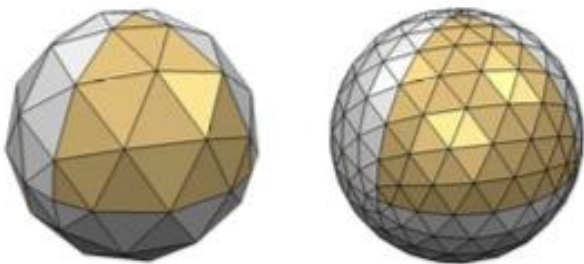


Morphable Model of Human Bodies

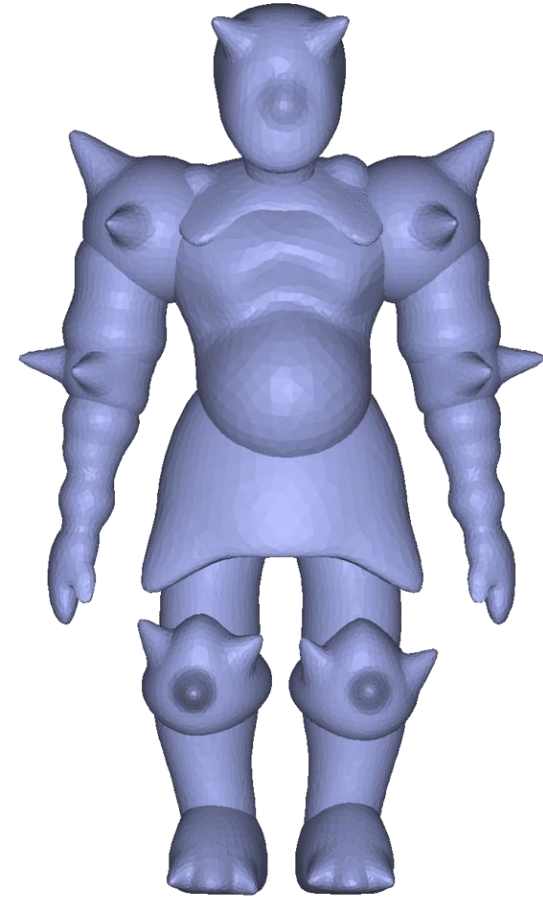


How to represent surfaces?

- Meshes are a popular, practical choice for surfaces
- Mesh = $\{V, F\}$
 - Vertices: $N \times 3$
 - Faces: $|F| \times \{3, 4, \dots\}$ polygons, “triangles”

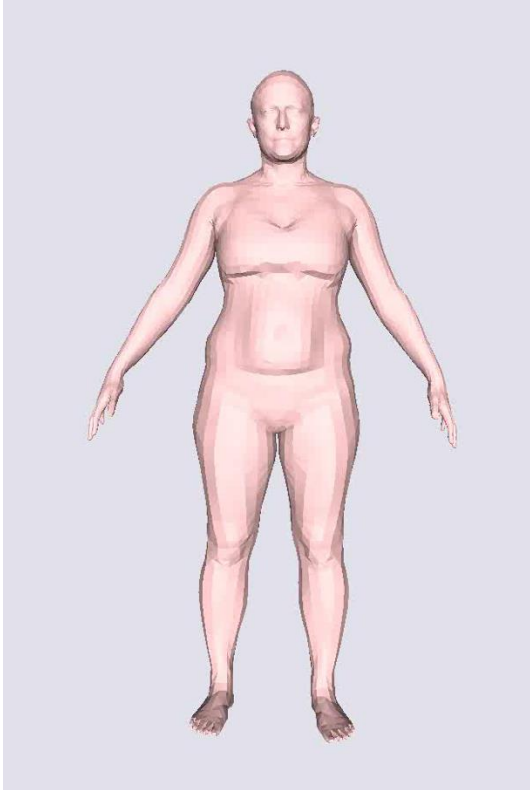


We need a low-dimensional parametrization!!

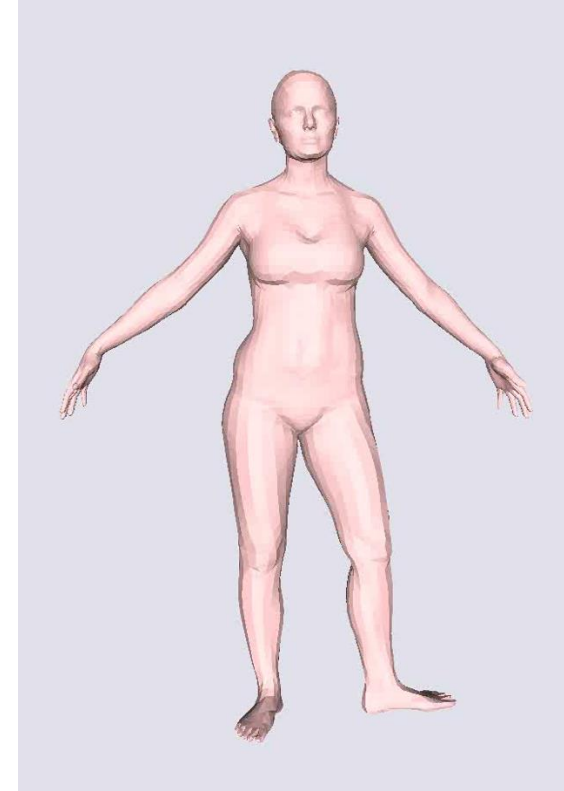


Key in modelling 3D Human Surfaces: Factorization into Shape and Pose

“Identity”



Individual Shape Variation



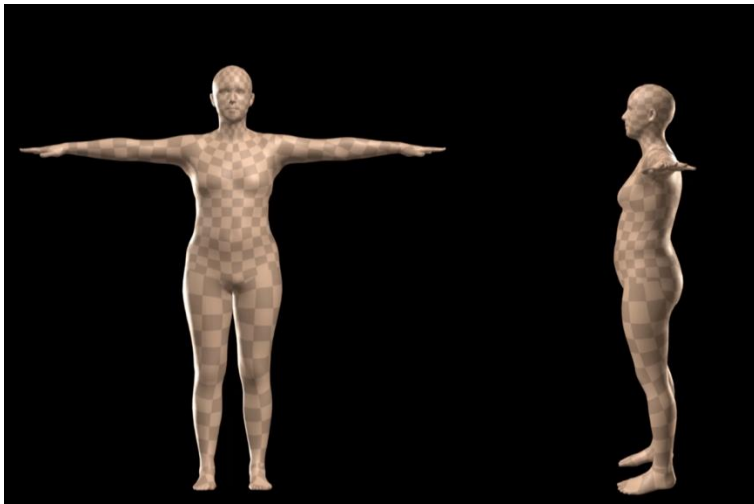
Pose changes (Articulation)

Skinned Multi-Person Linear Model (SMPL)

Shape: PCA coefficients

Pose: Rotation of joints

Mesh



$\vec{\beta}$

10
dimensions

+



$\vec{\theta}$

$3 \times 23 = 69$
dimensions

=

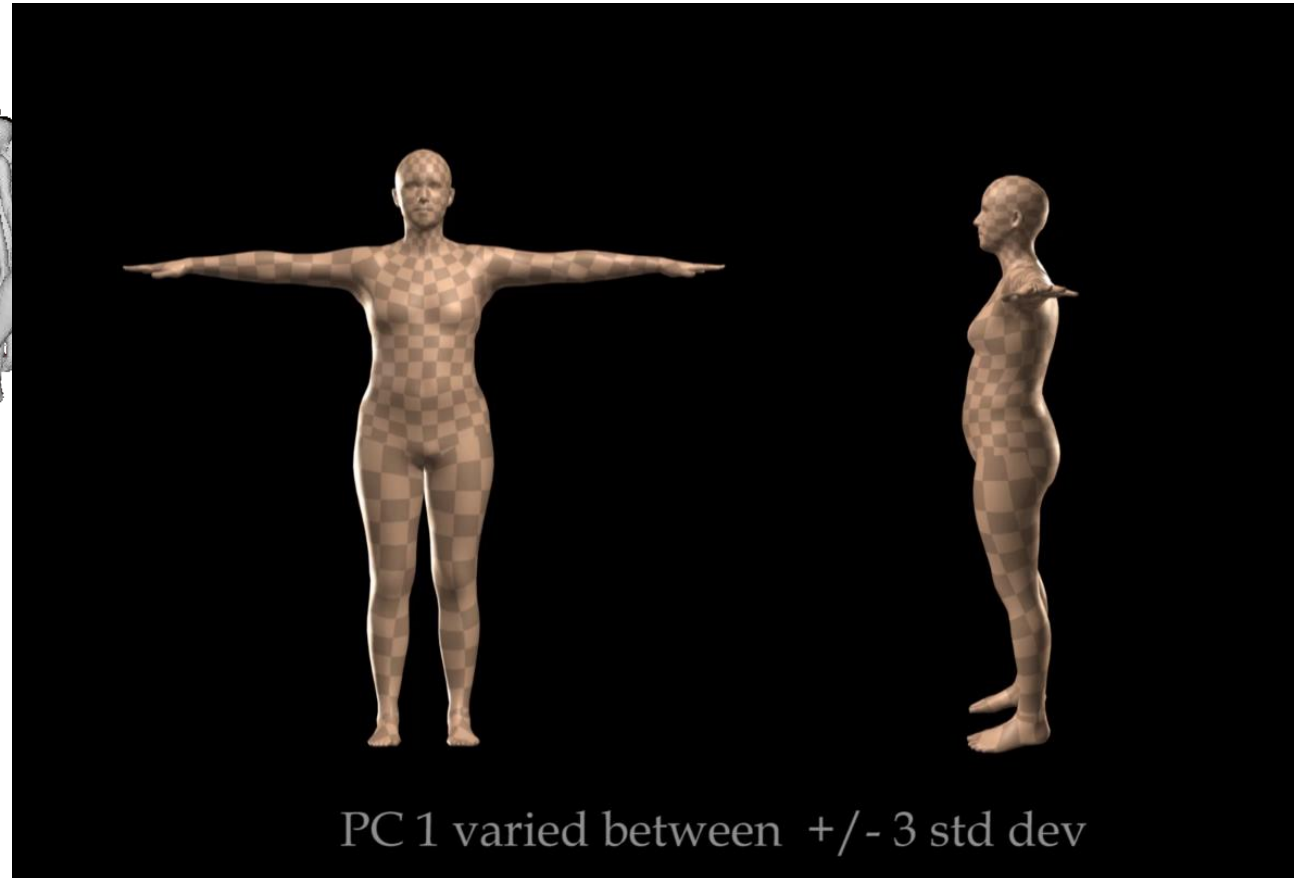


Learning Shape from 3D Scans

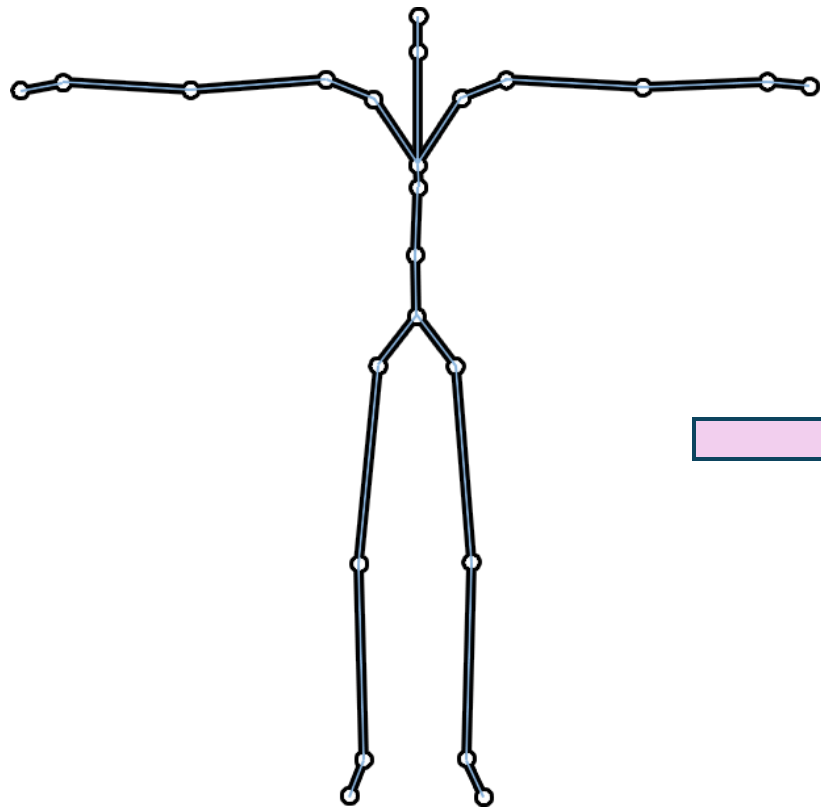
4000 bodies of different shapes in roughly the same pose.



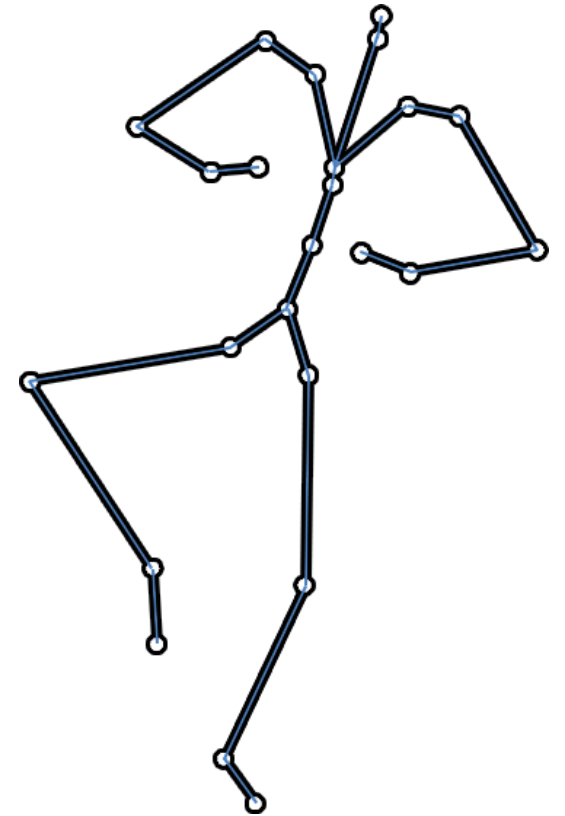
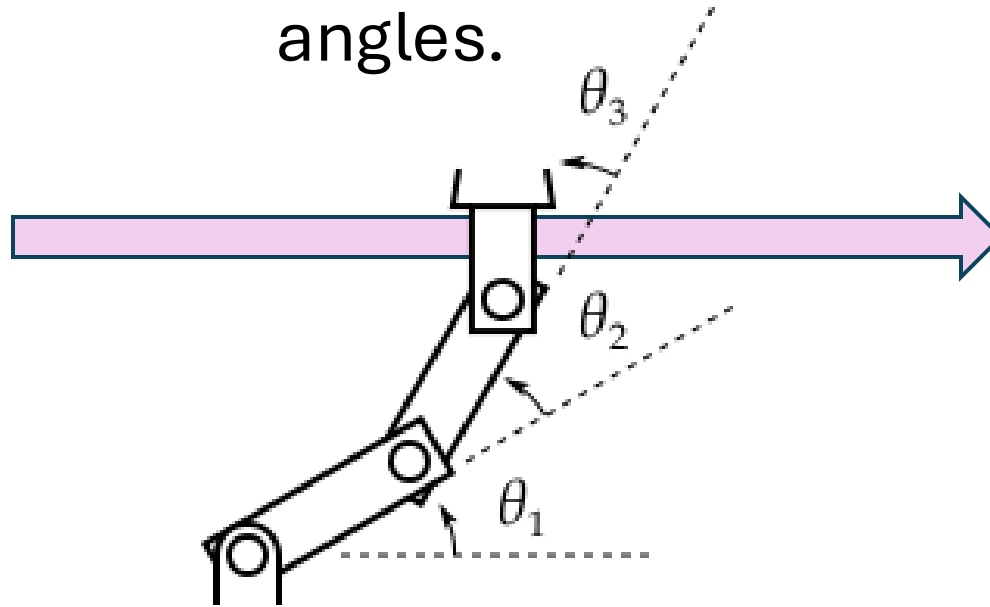
Run PCA on this:
Shape = linear
combination of basis
shapes



Pose: Forward kinematics on the skeleton tree



Defined by
relative joint
angles.

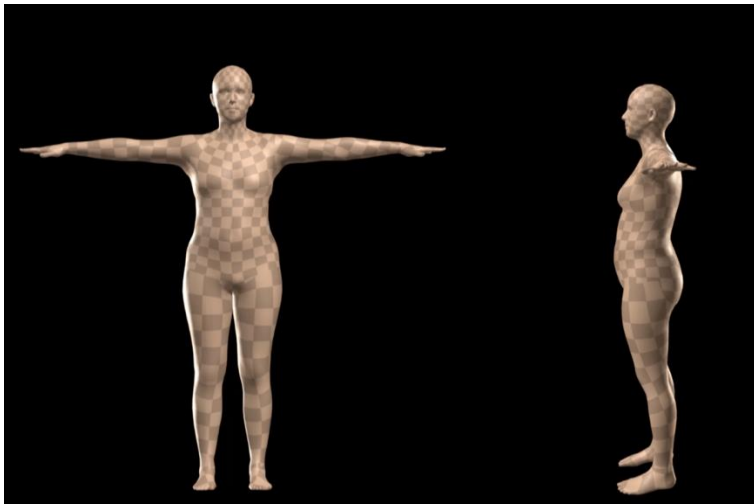


Skinned Multi-Person Linear Model (SMPL)

Shape: PCA coefficients

Pose: Rotation of joints

Mesh



$\vec{\beta}$

10
dimensions

+



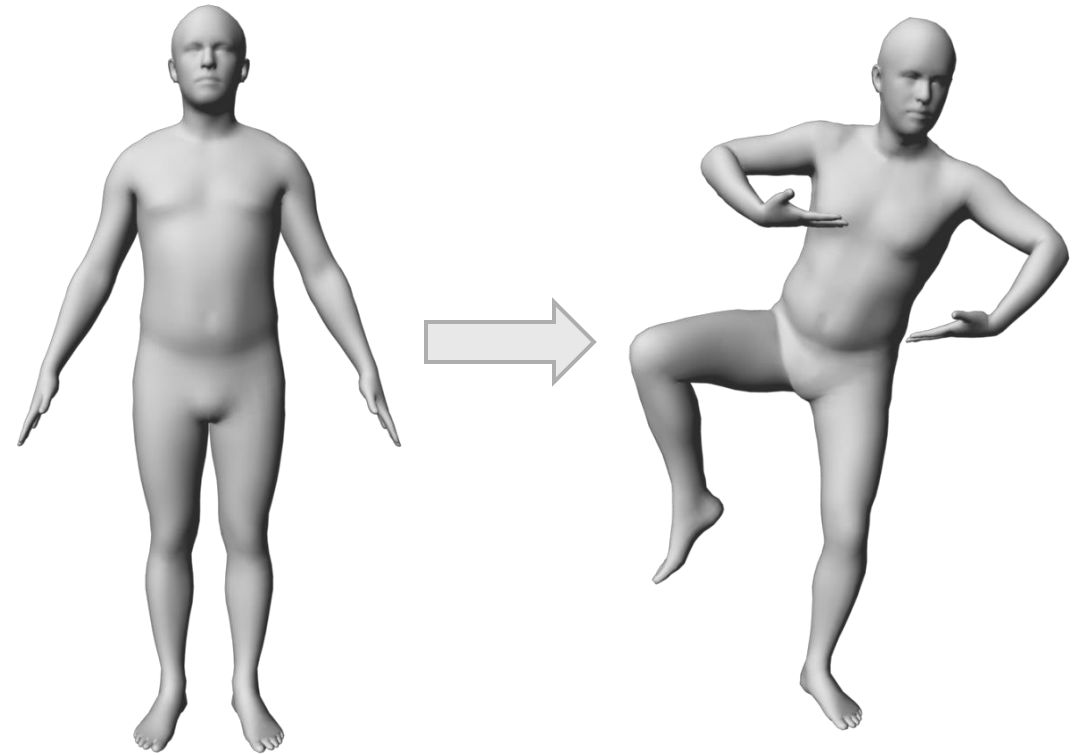
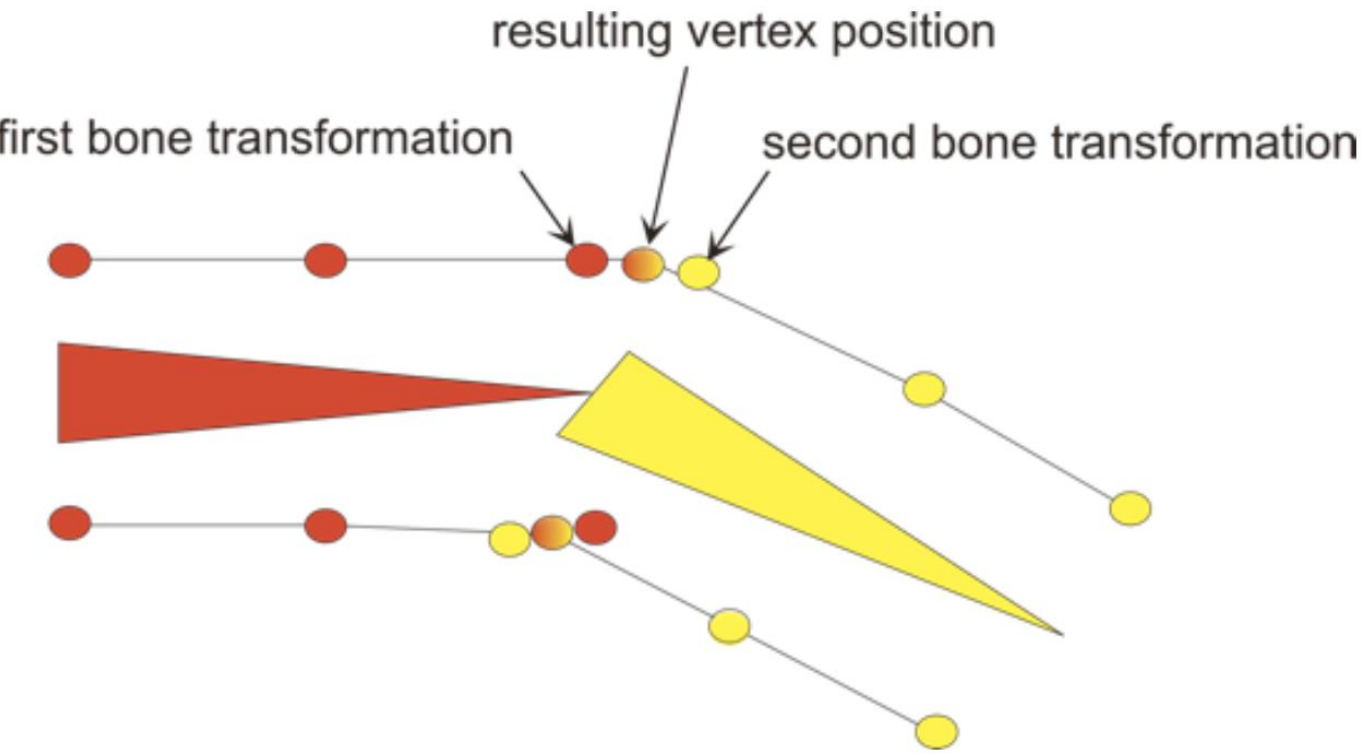
$\vec{\theta}$

$3 \cdot 23 = 69$
dimensions

=



Pose: Forward kinematics on the skeleton tree



Morphable Model of Human Bodies

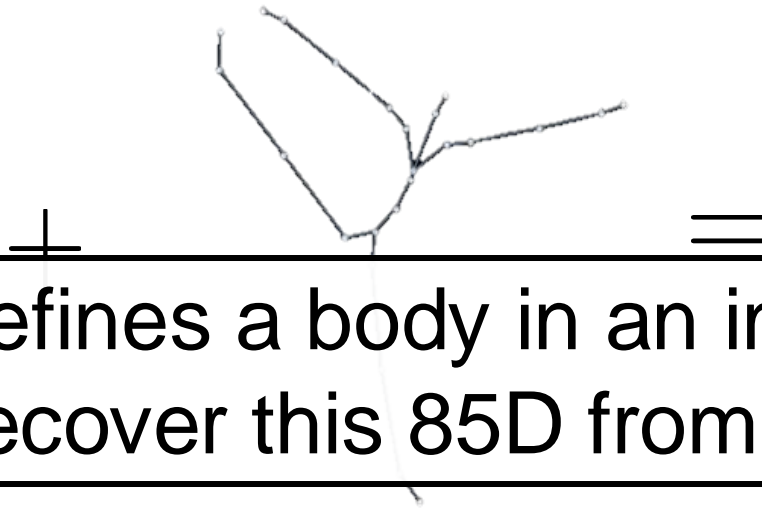
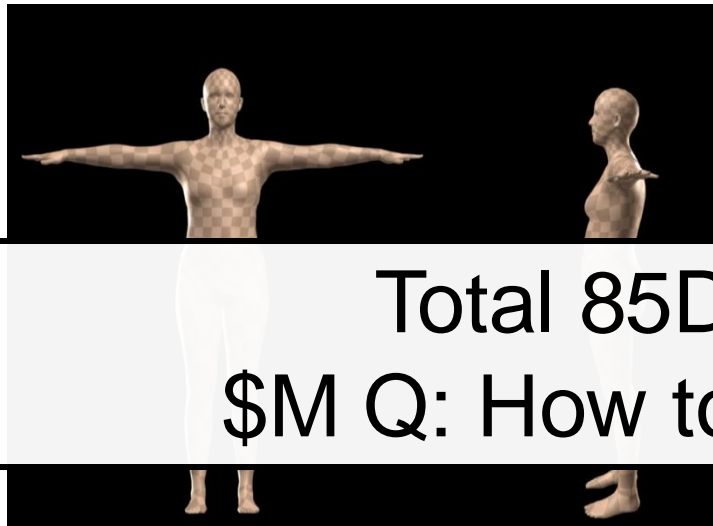


Morphable model for humans

Shape: low-D subspace

Pose: 23 Joint Rotations

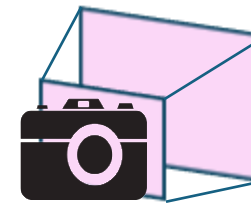
Mesh



Total 85D defines a body in an image!!!
\$M Q: How to recover this 85D from an image?

$$\vec{\beta} \in \mathbb{R}^{10}$$

$$\vec{\theta} \in \mathbb{R}^{23 \times 3}$$



$$\Pi \in \mathbb{R}^6$$

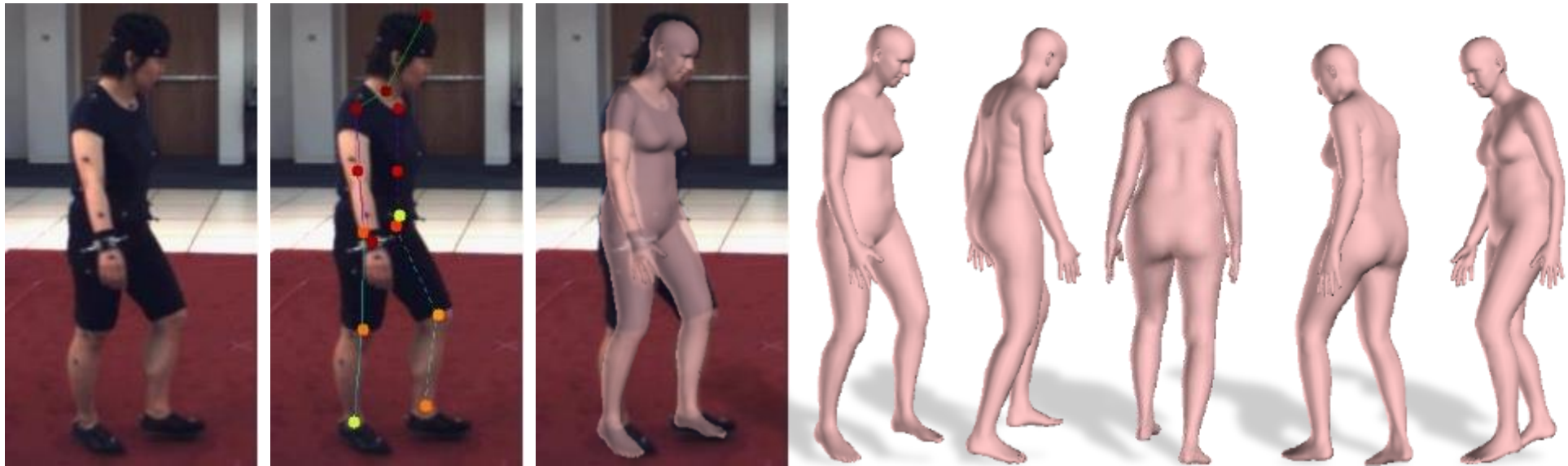
Back to images...

3D Shape and Pose from a Single Image

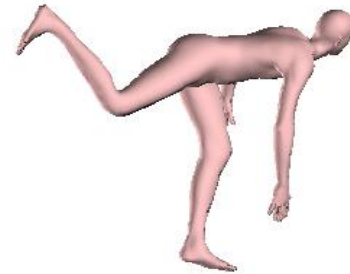


Overview: SMPLify

1. Automatic 2D joint detection via CNNs
2. Fit SMPL pose and shape parameters



SMPLify Objective Function



camera joints

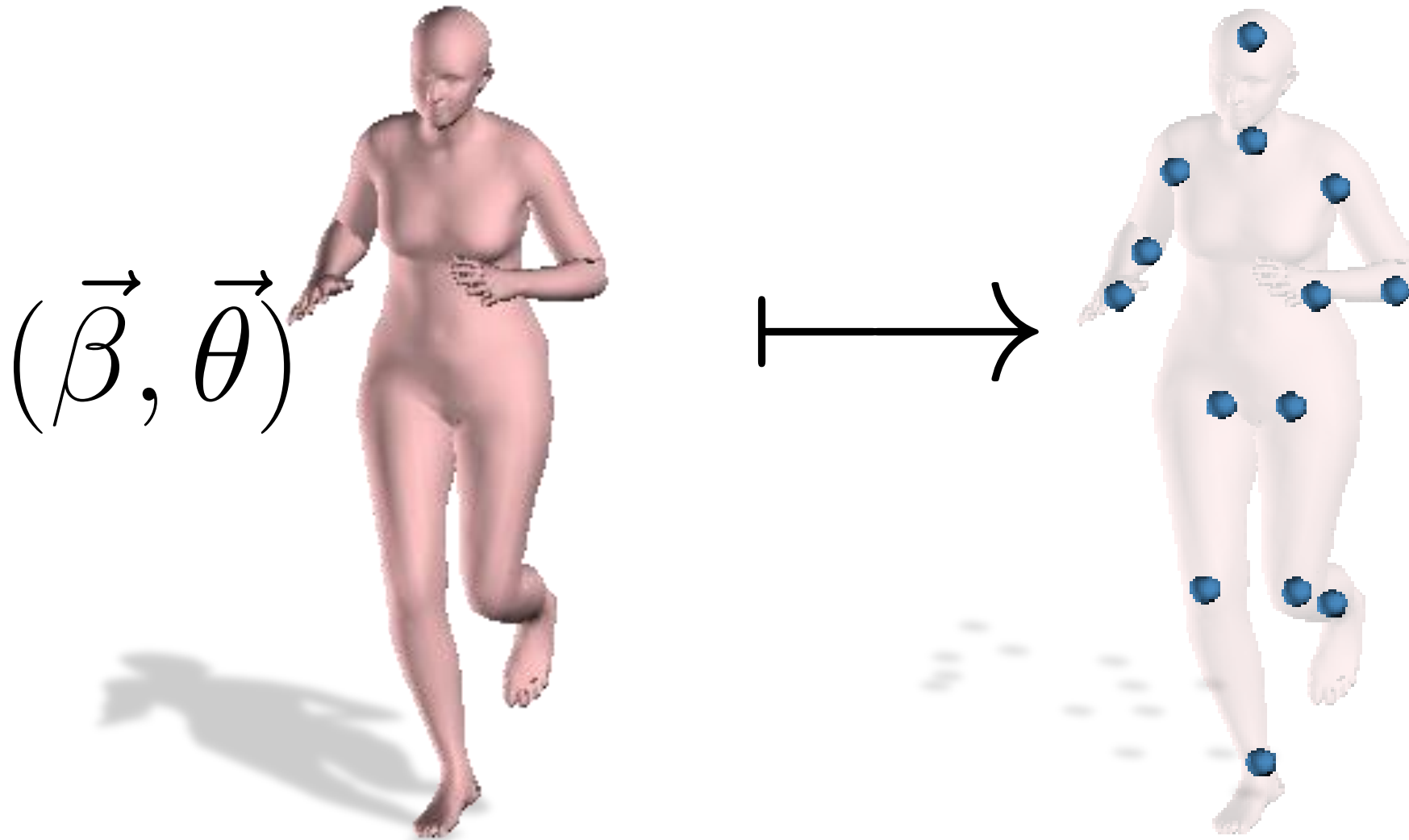
$$E(\vec{\beta}, \vec{\theta}, K; \mathcal{J}_{est}) =$$

$$E_J(\vec{\beta}, \vec{\theta}, K; \mathcal{J}_{est}) + E_a(\vec{\theta}) + E_\theta(\vec{\theta}) + E_{sp}(\vec{\theta}, \vec{\beta}) + E_\beta(\vec{\beta})$$

Data term

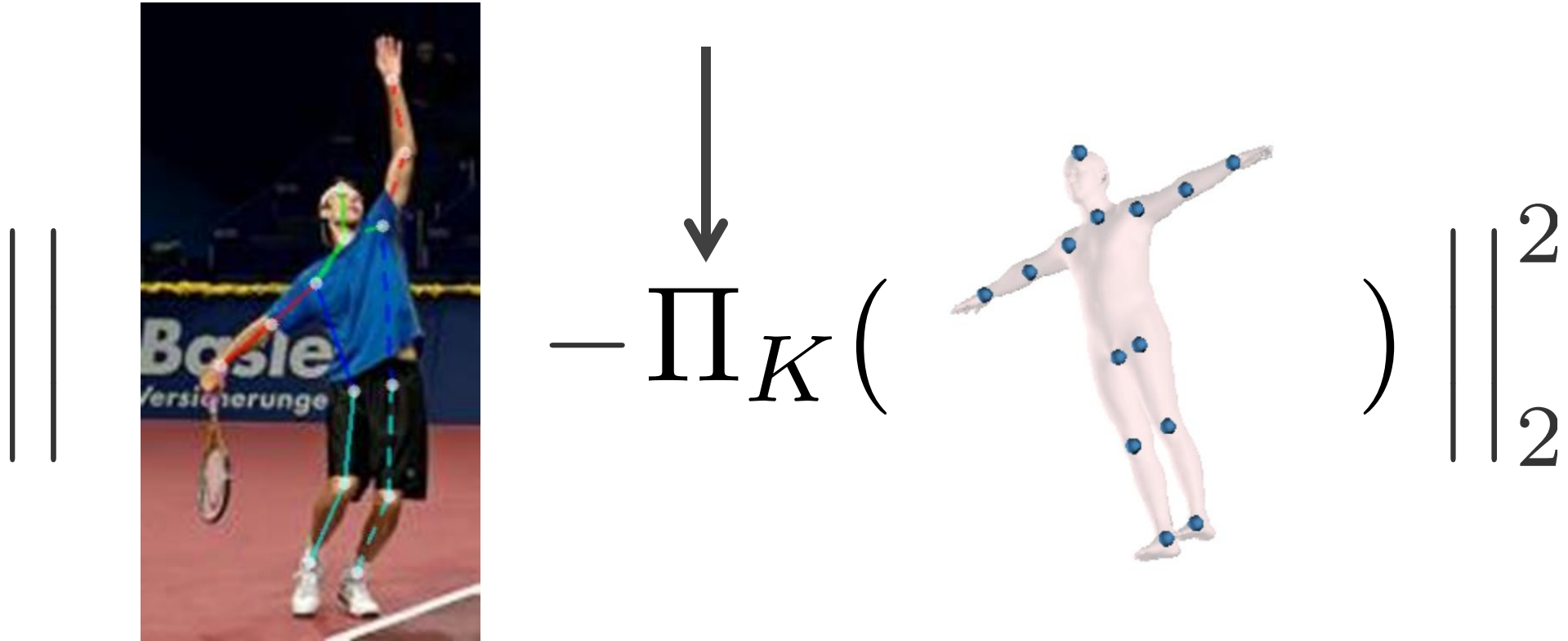
Priors

Data Term: Joint Reprojection Error



Data Term: Joint Reprojection Error

Camera Projection

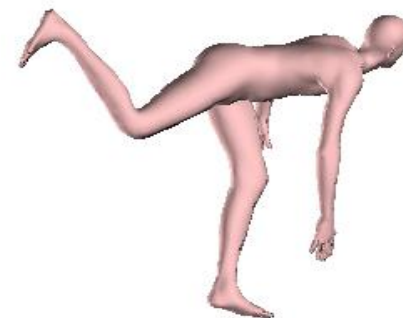


Summary: Fit to 2D joints

1. Automatic 2D joint detection via CNN



2. Solve for pose and shape that explain the 2D joints

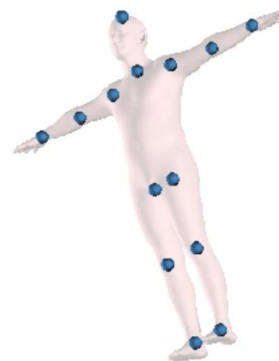


$$\min_{\beta, \theta, \Pi}$$



-

Π



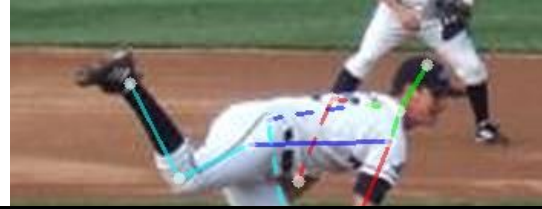
)

$$\left\| \begin{matrix} 2 \\ 2 \end{matrix} \right\|$$

+ lots of priors

Approach: Fit to 2D joints

1. Automatic 2D joint detection via CNN



2. Solve for pose and shape that explain the 2D joints



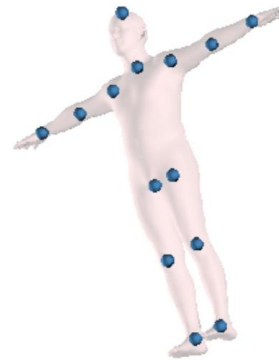
Only looks at 2D joints, not the image
Optimization based inference = too slow for video

$$\min_{\beta, \theta, \Pi}$$



-

Π



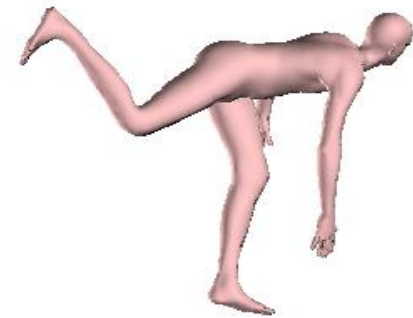
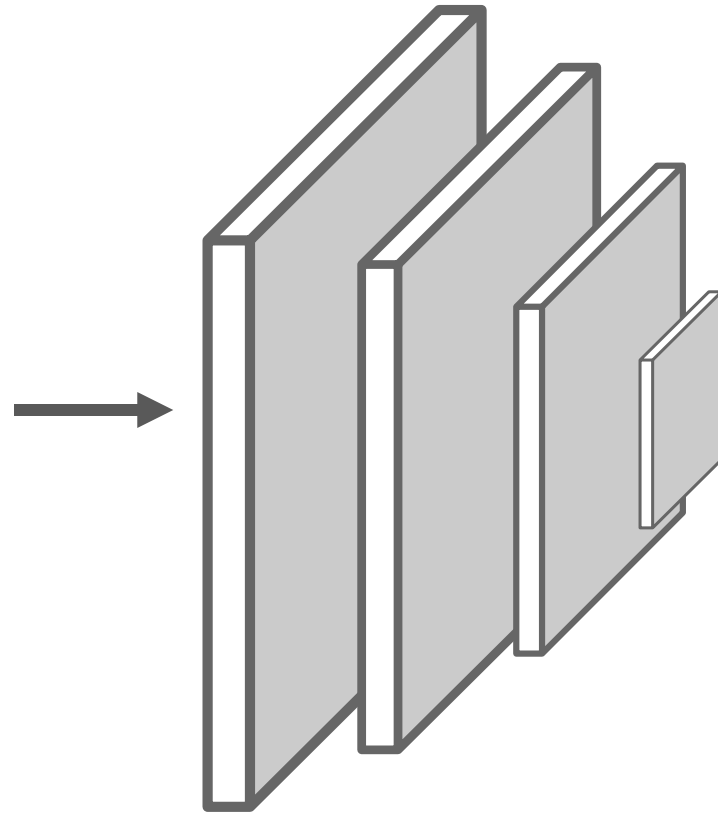
)

$$\left\| \right\|_2^2$$

+ lots of priors

Why not just throw a deep network at it?

- Image in, 85D human parameters out!!



$\{\beta, \theta, \Pi\}$

Challenges

1. Lack of real paired 2D-to-3D labels



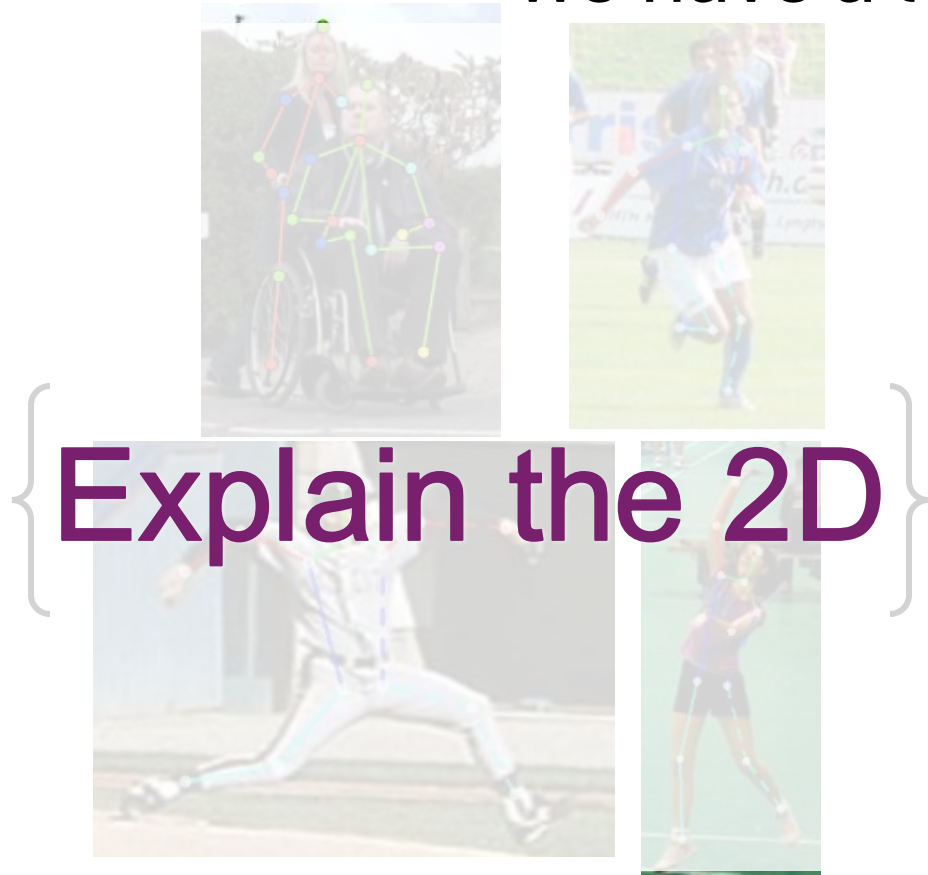
2. Depth ambiguity



[CJ Taylor CVPR 2000]

Solution

Even though we don't have paired 2D-to-3D labels, we have a lot of **unpaired** labels

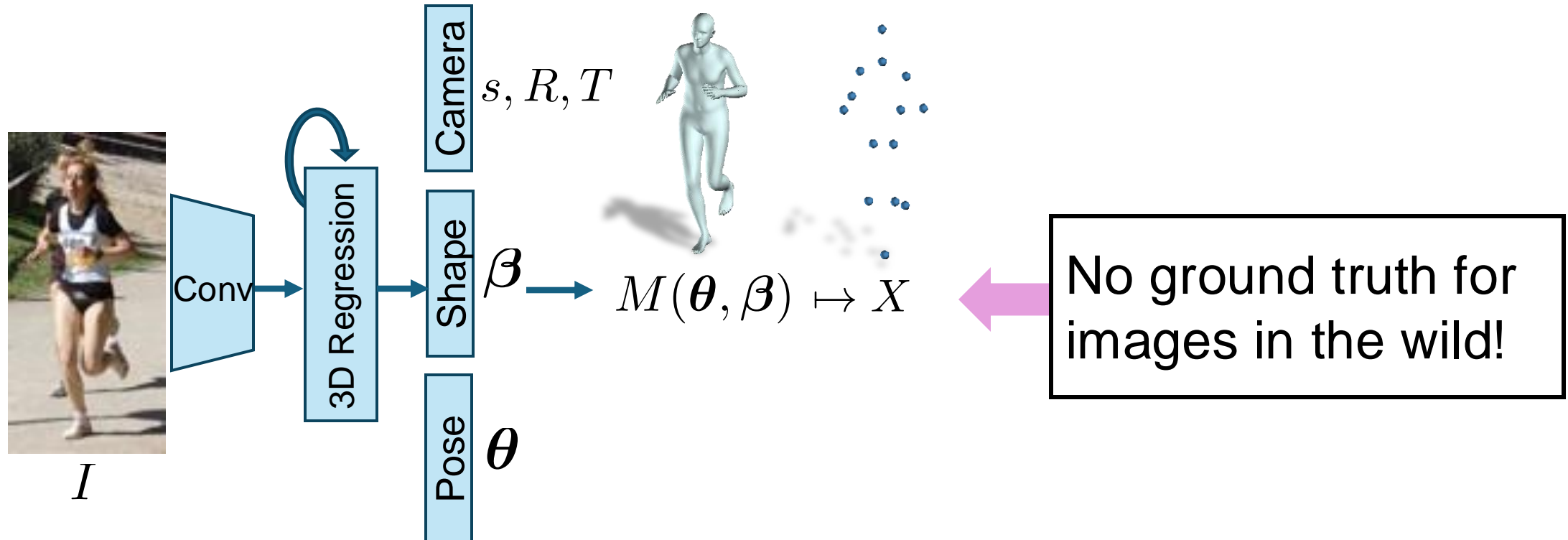


2D Labeled images
[LSP, MPII, MS COCO,...]

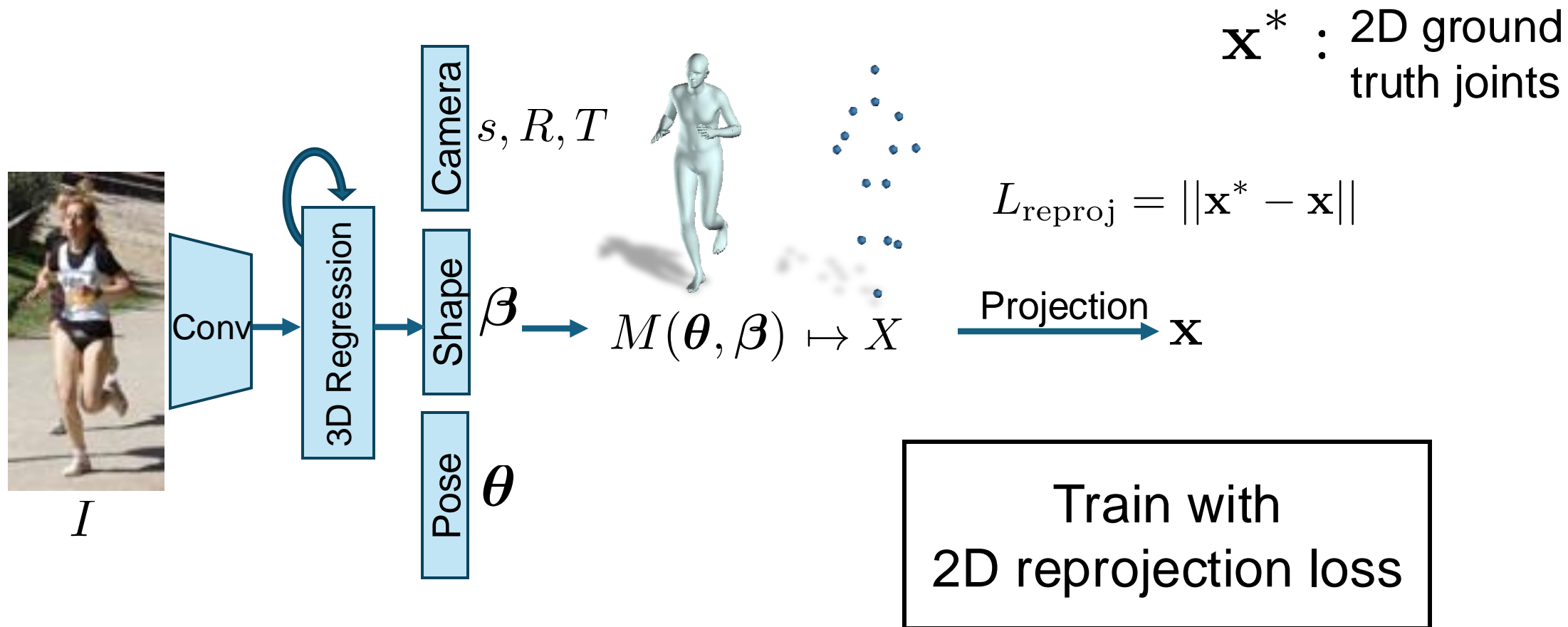


3D Scans/Motion Capture
[CMU Mocap, CAESER, JointLimits..]

Overview: Human Mesh Recovery (HMR)



Overview: Human Mesh Recovery (HMR)



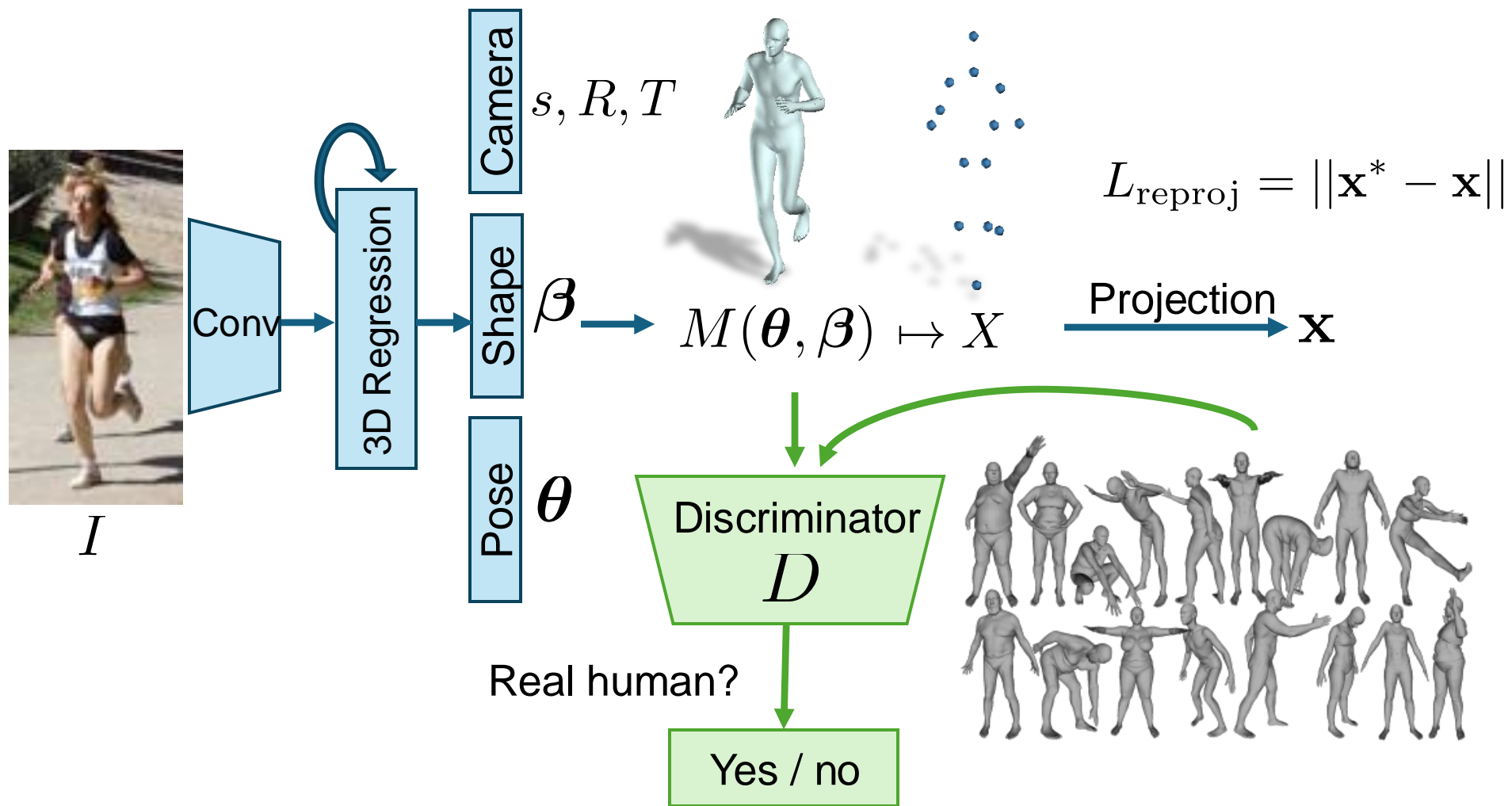
Without any 2D-to-3D supervision...



More monsters from training



Overview: Human Mesh Recovery (HMR)



Training Data

In-the-wild

3D

2D

Human3.6M

[Ionescu et al. PAMI'14]



MS COCO

[Lin et al. ECCV '14]



Overview: HMR



I

Conv

3D Regression

Camera

s, R, T

Shape

β

Pose

θ

$M(\theta, \beta) \mapsto X$



Projection \mathbf{x}

$L_{\text{reproj}} = \|\mathbf{x}^* - \mathbf{x}\|$

~~$L_{3D} = \|X^* - X\|$~~

Can be trained in a fully weakly-supervised mode

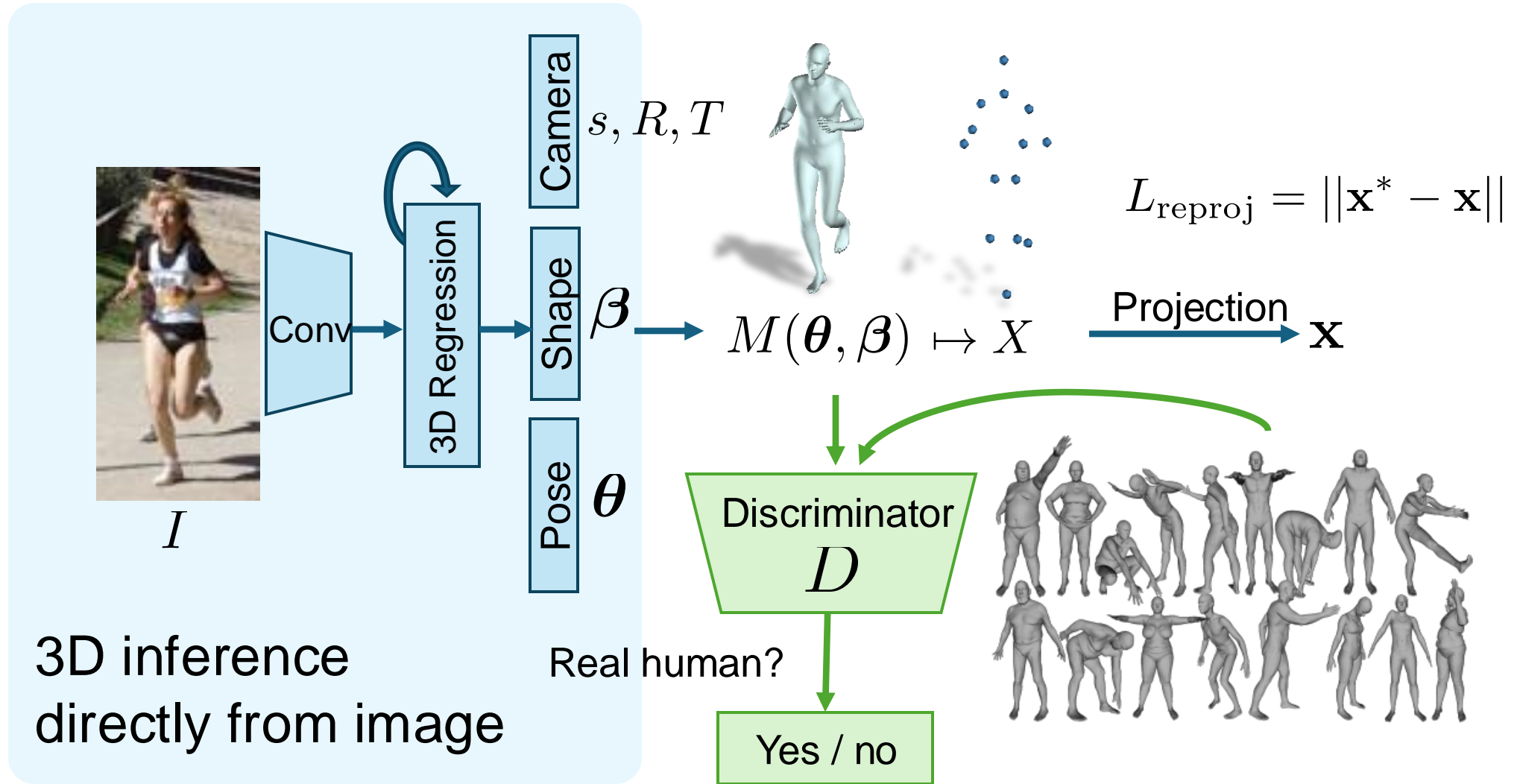
Discriminator D

Real human?

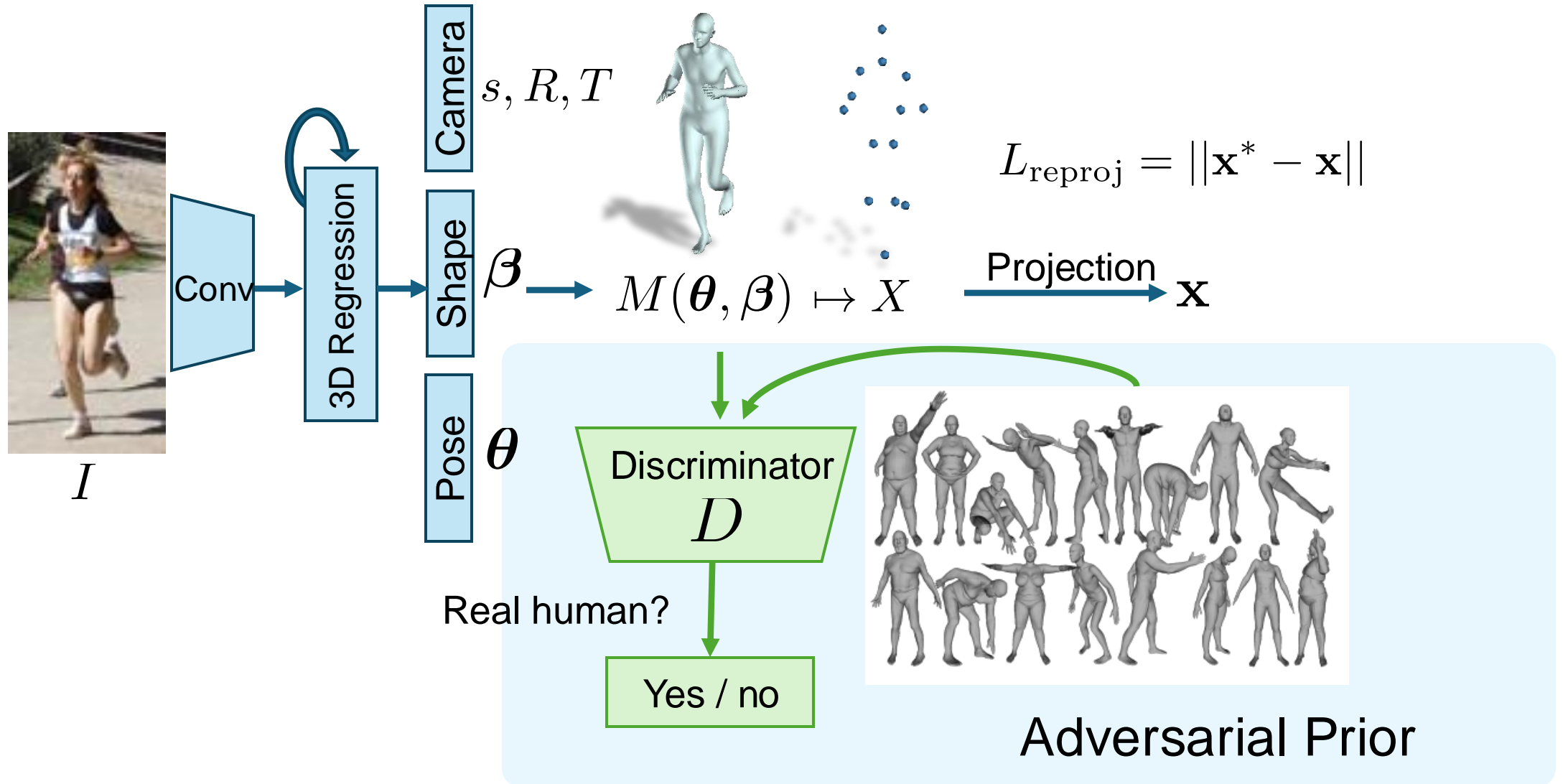
Yes / no



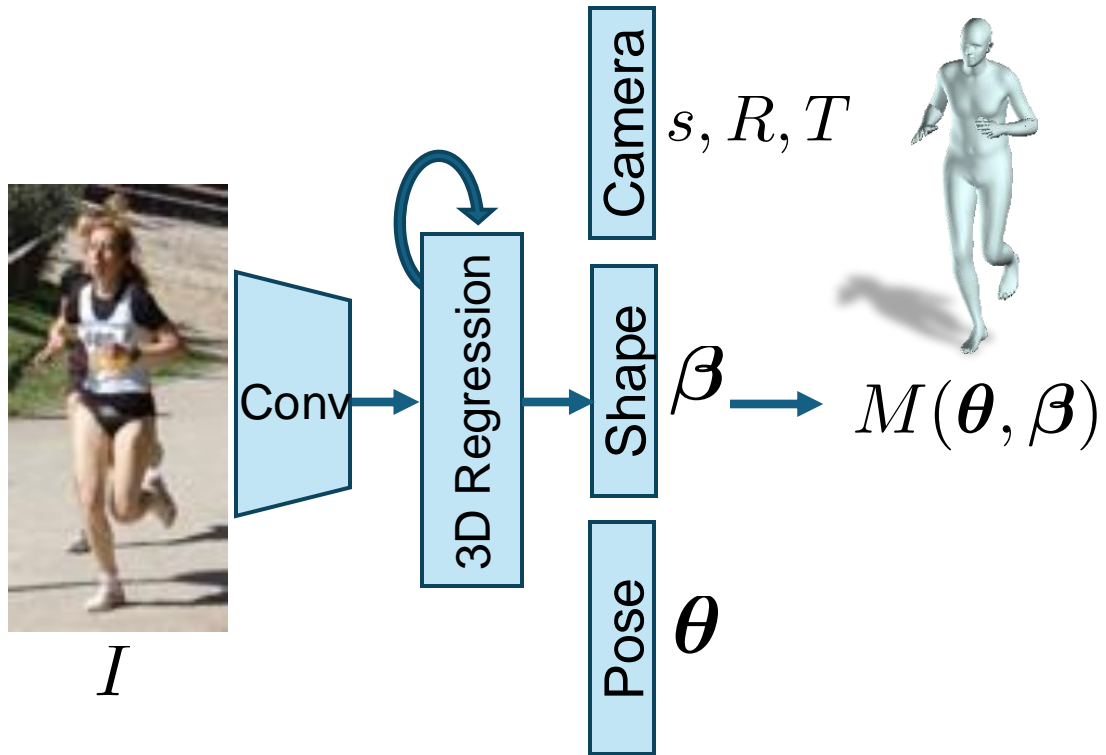
Overview: Human Mesh Recovery (HMR)



Overview: Human Mesh Recovery (HMR)



Test time: just feed forward

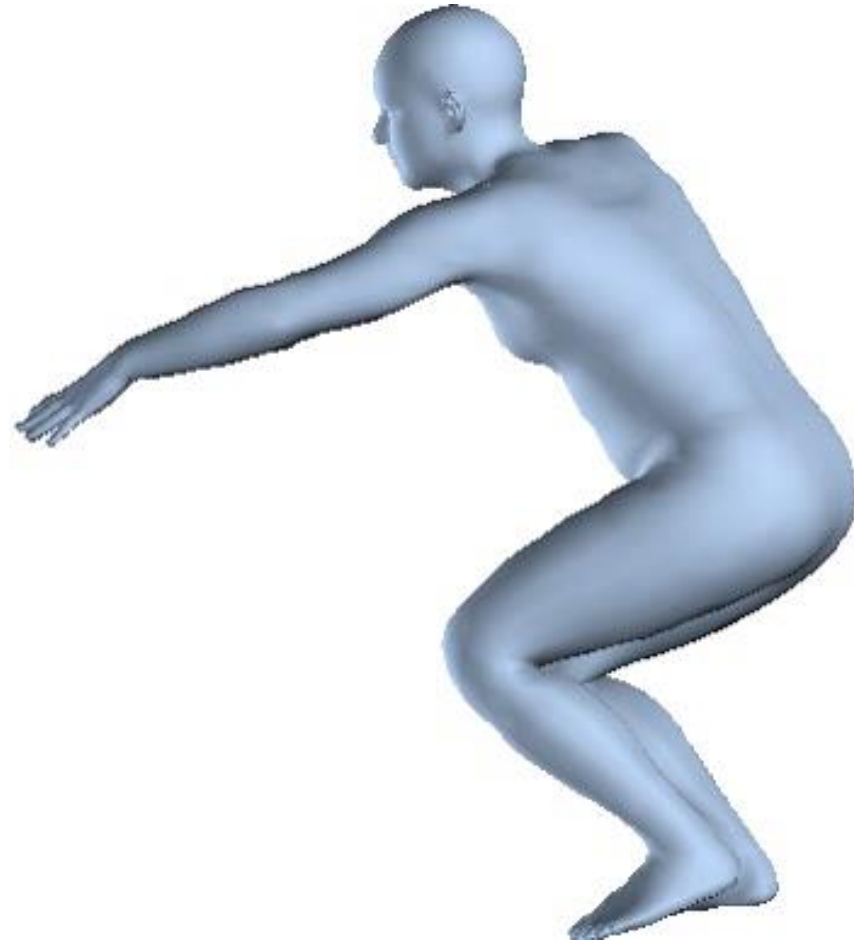


Benefits of recovering a deformable model

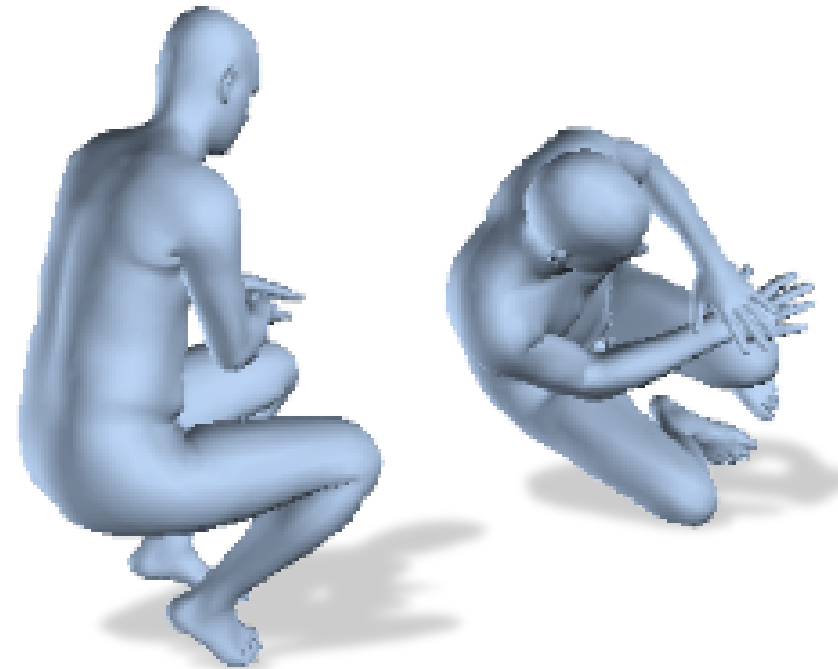
Correspondences across recovered bodies (part segmentations)



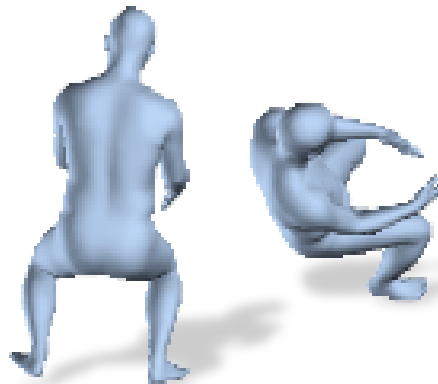
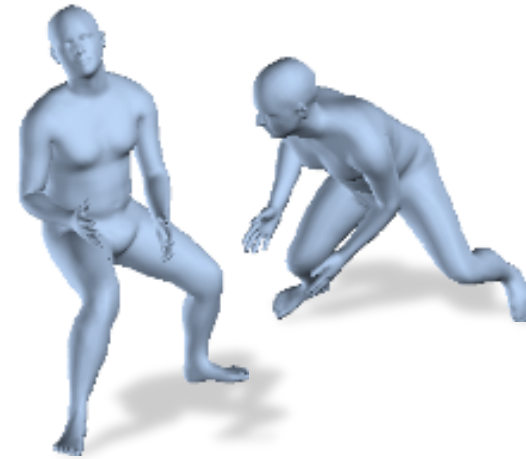
Amodal/holistic prediction

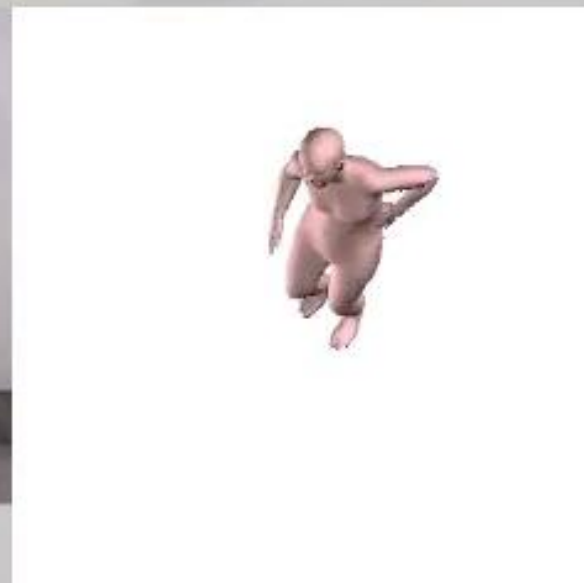
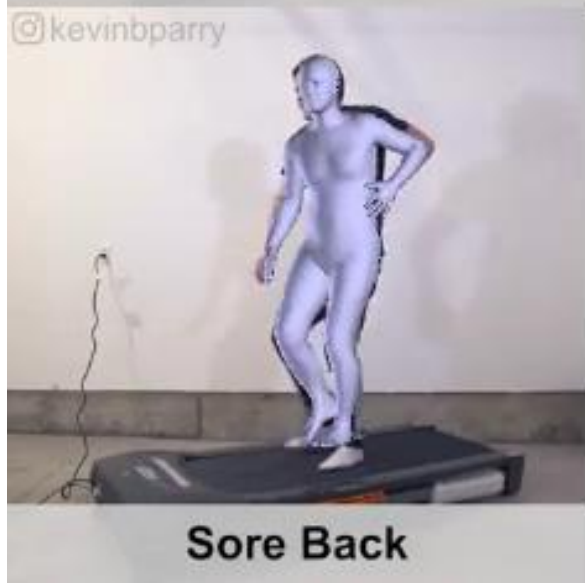
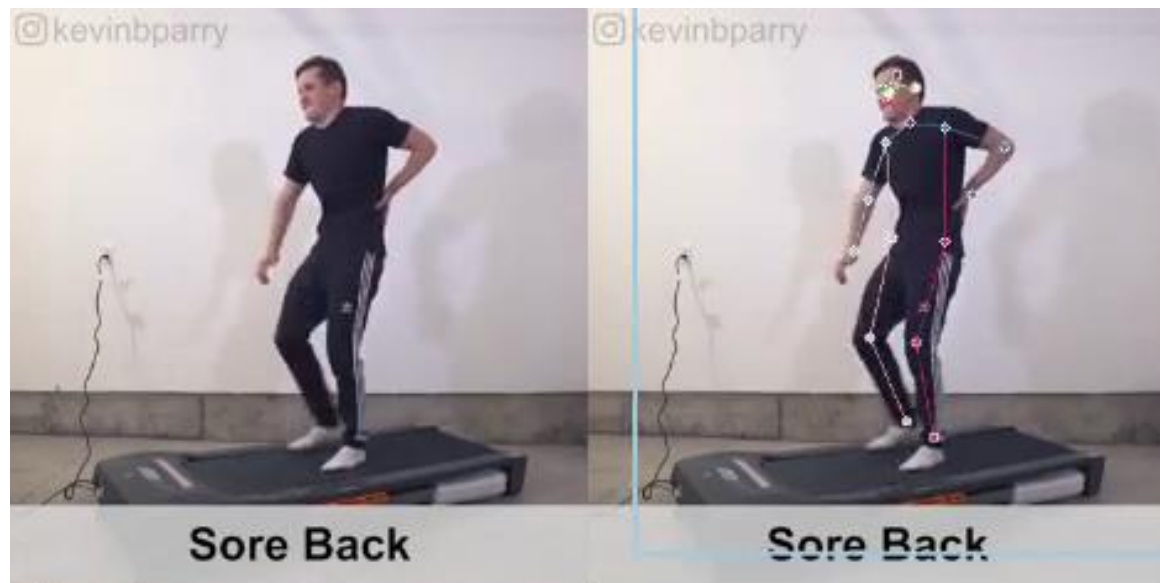


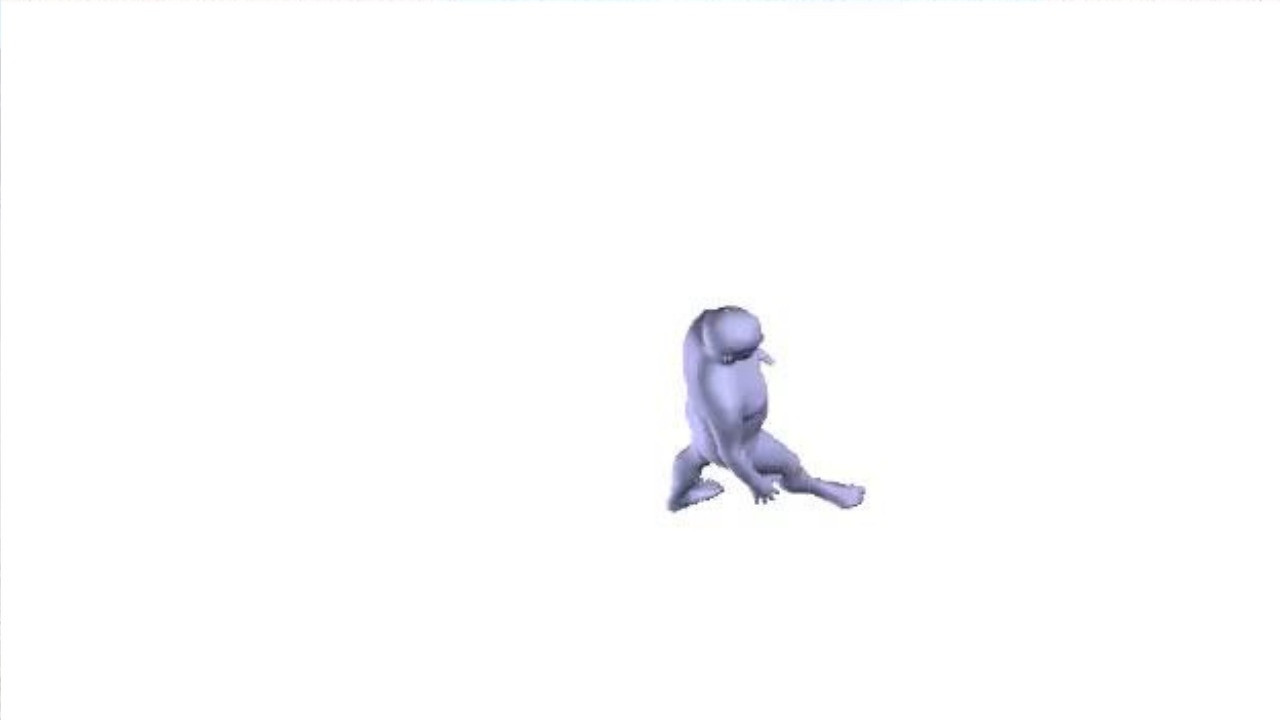
Prediction on occluded body parts



Qualitative results on COCO w/occlusion + clutter



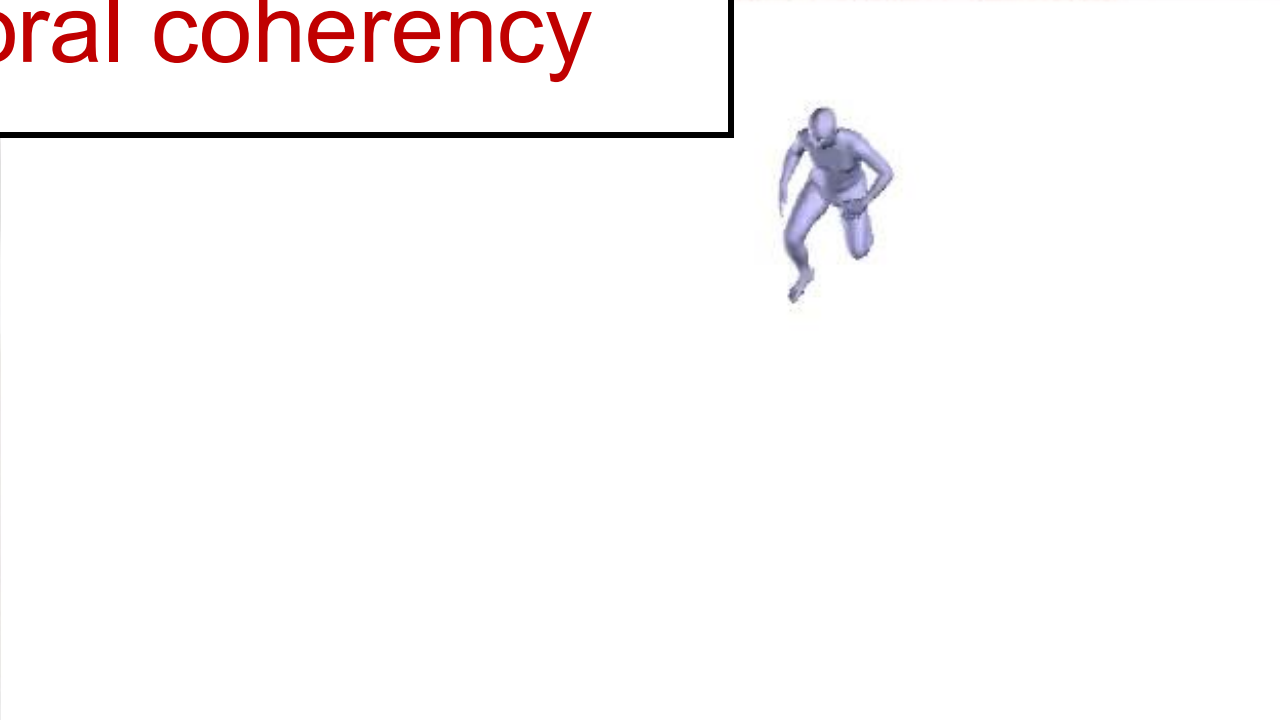








+ Good per frame performance
- Lacks temporal coherency



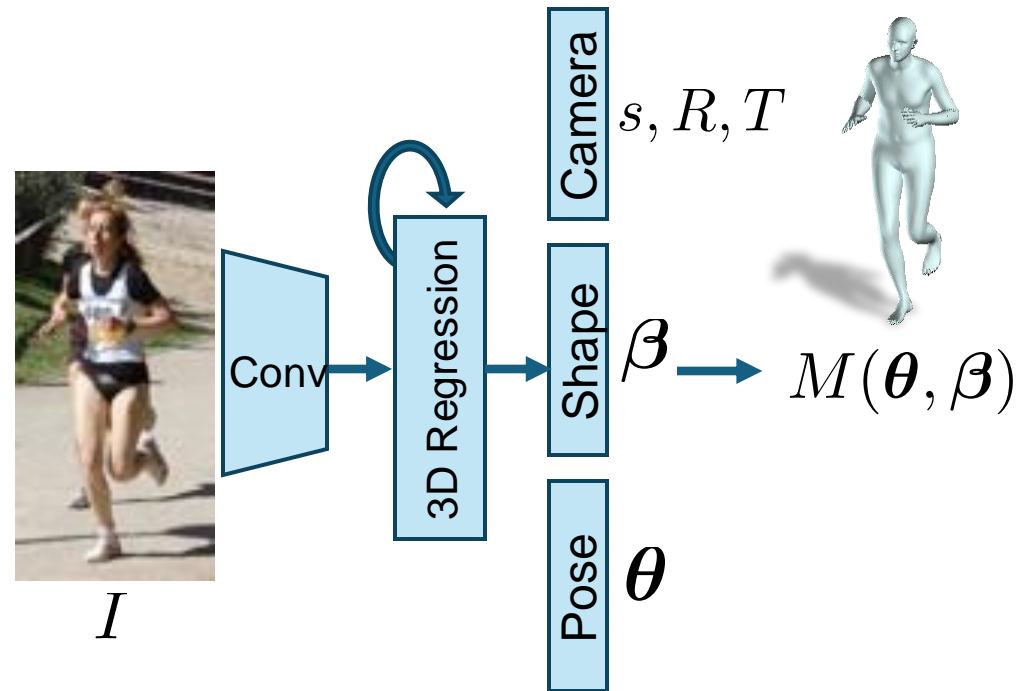
Recap: model based 3D human perception

Iterative optimization in 2016

$$\min_{\beta, \theta, \Pi} \left\| \begin{array}{c} \text{Image} \\ \text{Baseline} \\ \text{Versicherung} \end{array} - \Pi \left(\begin{array}{c} \text{3D Model} \\ \text{Joints} \end{array} \right) \right\|_2^2 + \text{lots of priors}$$

One-shot inference in 2018

Complementary!
Discuss Pros and Cons

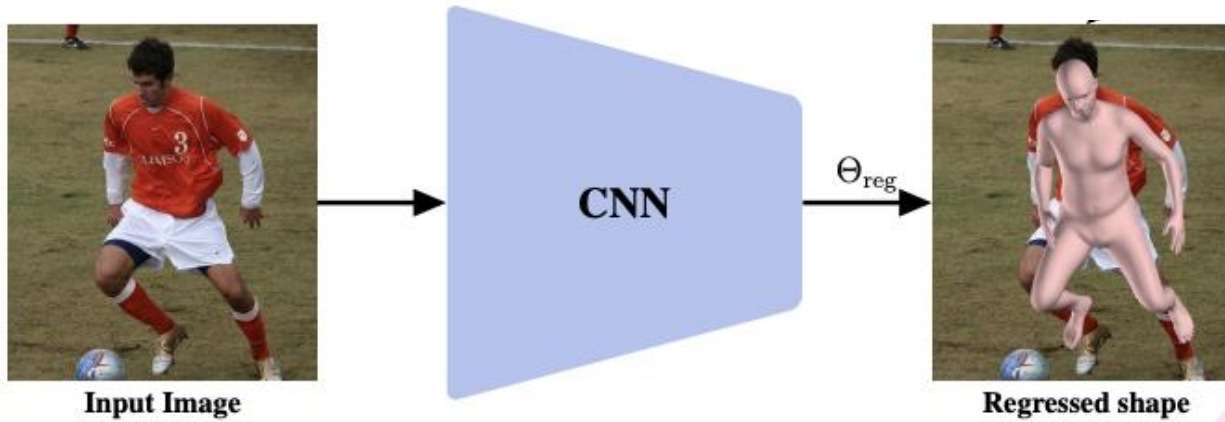


SPIN (SMPL oPtimization IN the loop)

[Kolotouros and Pavlakos et al. ICCV 2019]

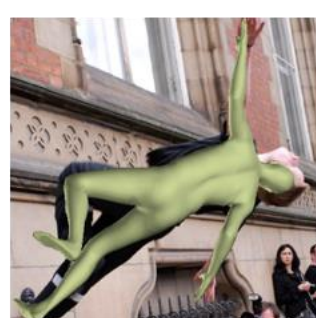


SPIN [Kolotouros and Pavlakos et al. ICCV 2019]



SPIN is self-improving

Starting from an initial set of fits, our method can **improve** them.



Image

Initial fit

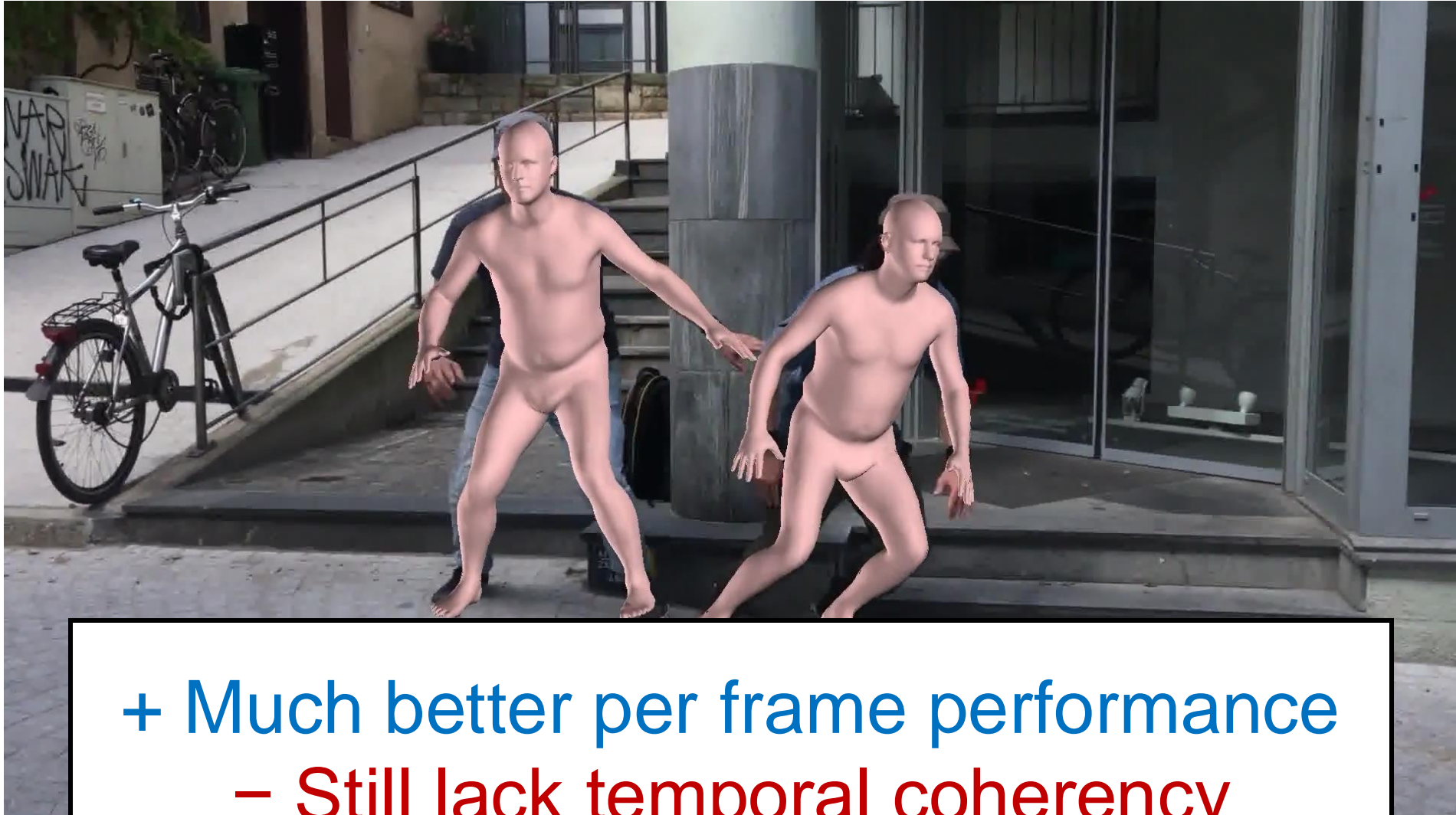
Final fit

Image

Initial fit

Final fit

SPIN results



- + Much better per frame performance
- Still lack temporal coherency

SMPL-X model



SMPL-X estimated independently on each frame

Model fitting

Pavlakos et al.
CVPR 2019

Objective function

$$E(\beta, \theta, \psi) =$$

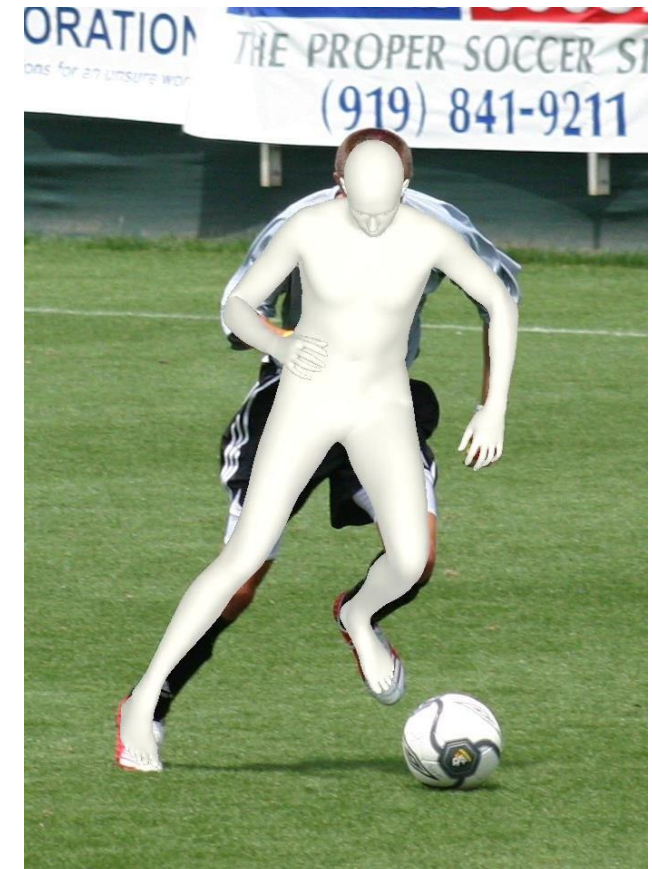
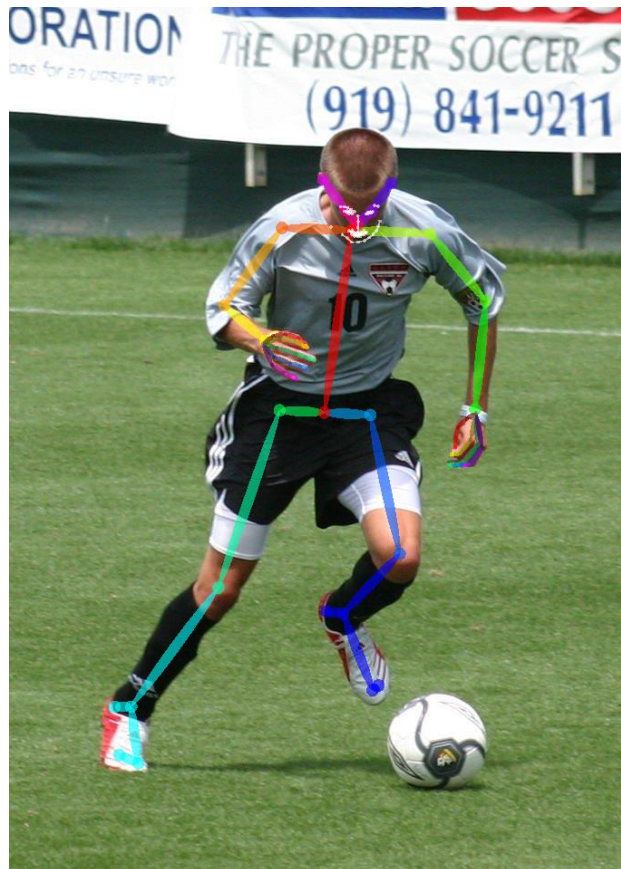
Data term

joints reprojection

+

Priors

pose, shape, expression, interpenetration



TODO replace with 4D Humans, HaMeR, SLAHMR

Progress on Human Mesh Recovery — from 2018 to 2023

Human Mesh Recovery (HMR)

CVPR 2018

Kanazawa, Black, Jacobs, Malik



Human Mesh Recovery 2.0

ICCV 2023

Goel, Pavlakos, Rajasegaran, Kanazawa*, Malik*, ICCV 2023



Per-frame estimation — no smoothness applied
Color = Identity

Human Mesh Recovery 2.0

ICCV 2023

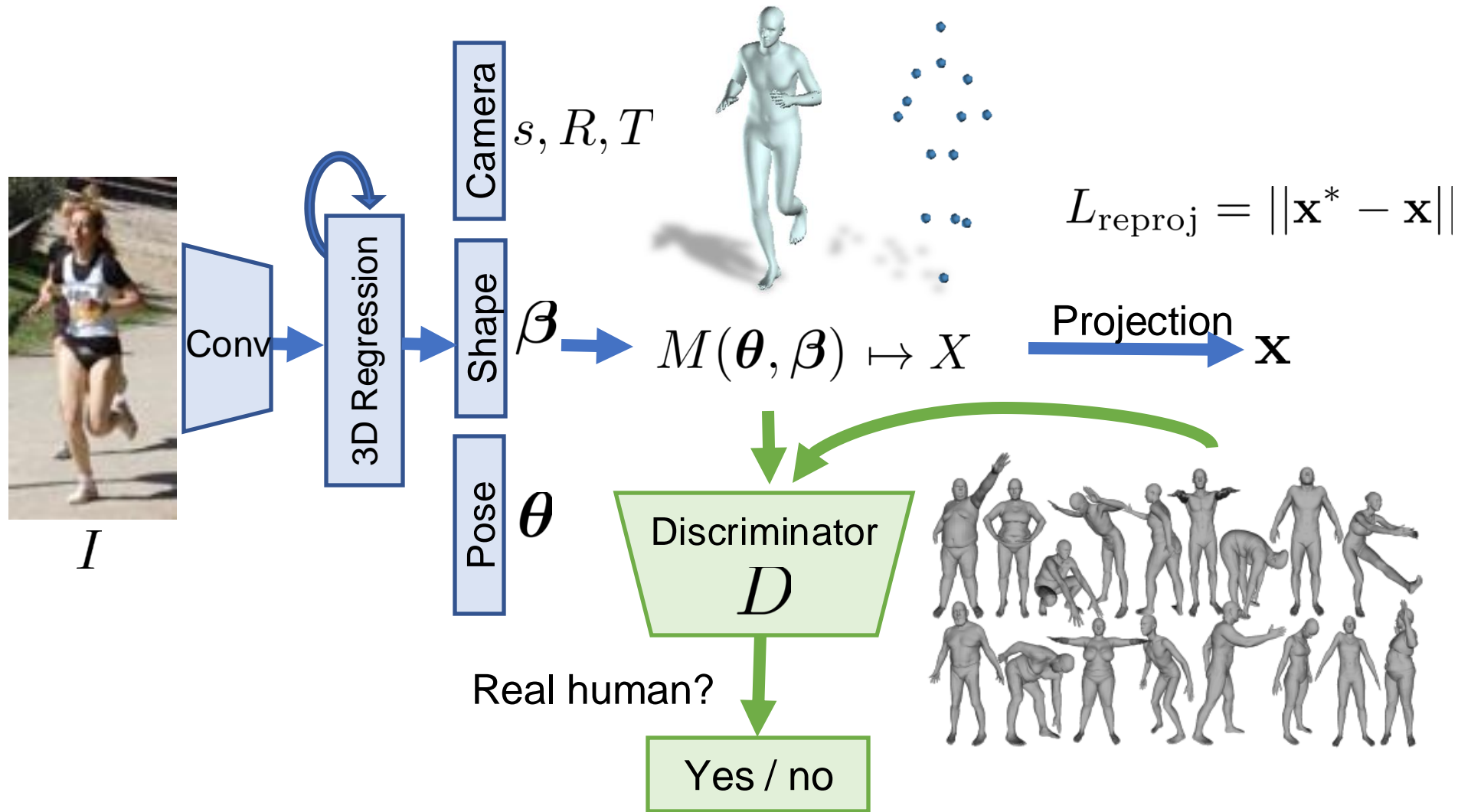
Goel, Pavlakos, Rajasegaran, Kanazawa*, Malik*, ICCV 2023



Per-frame estimation — no smoothness applied

Color = Identity

Human Mesh Recovery (HMR) 2018



Recipe: Big Model and Big Data

Before



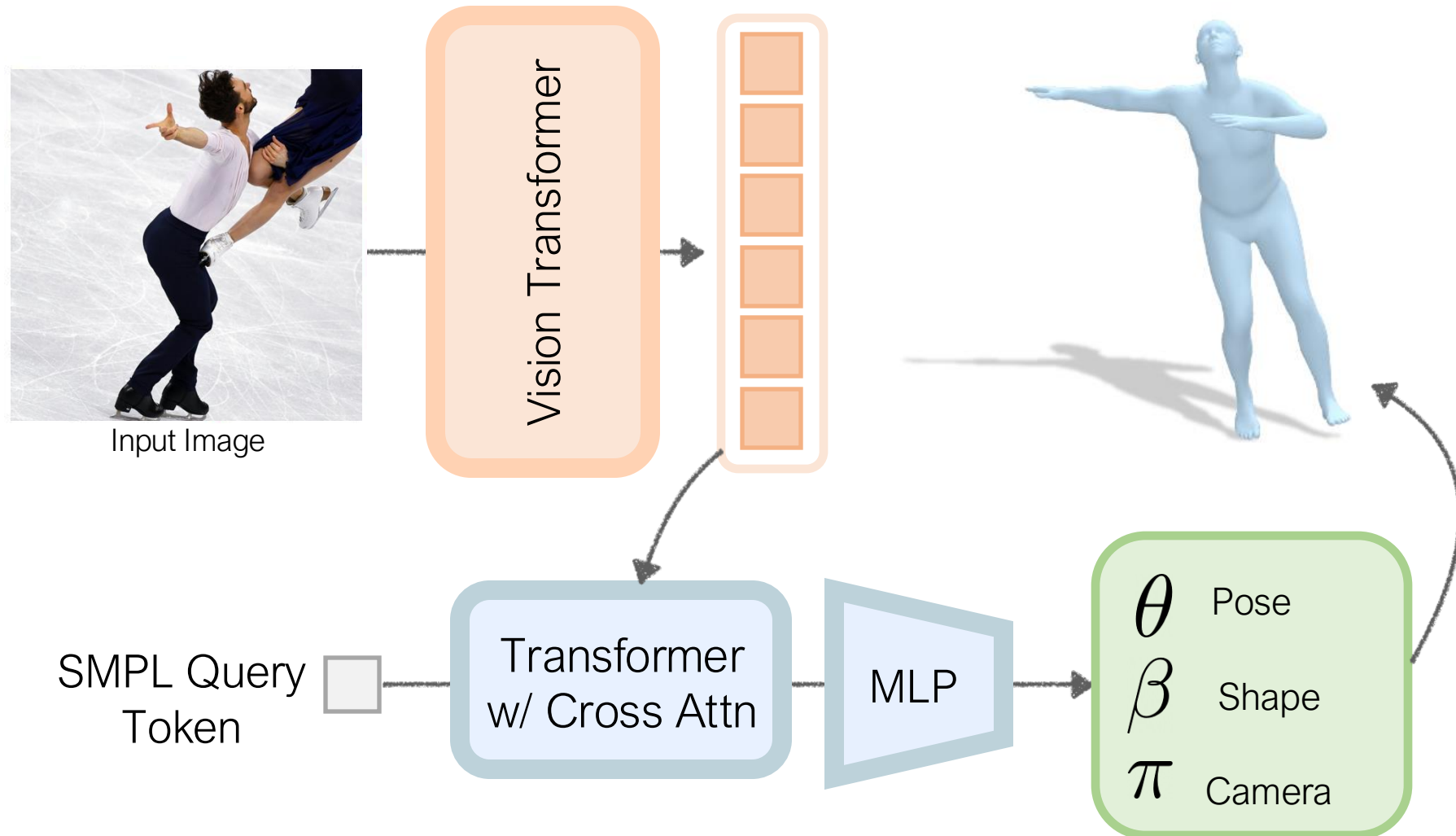
Ours



per-frame estimation - no smoothness applied

Recipe: Big Model and Big Data

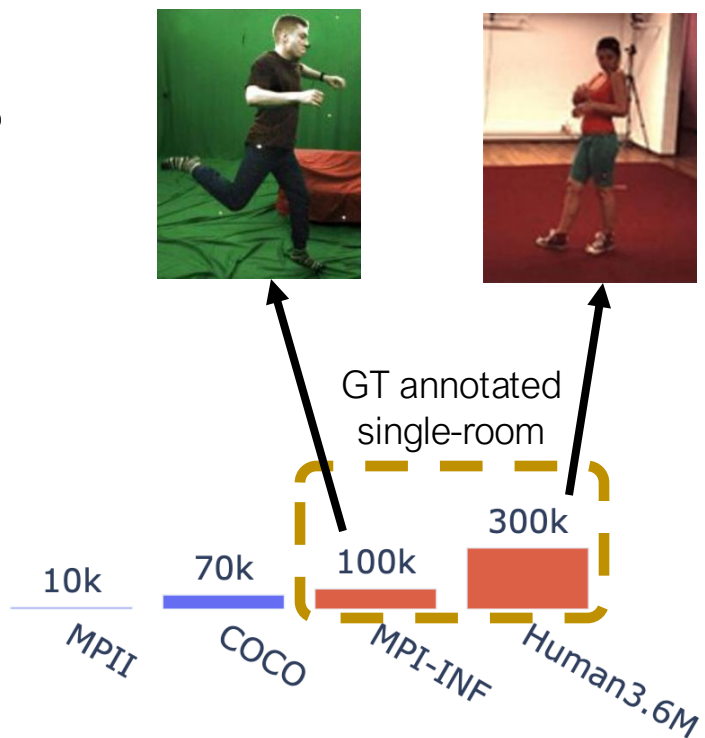
HMR 2.0



Recipe: Big Model and Big Data

Automatic dataset labelling

Number of images



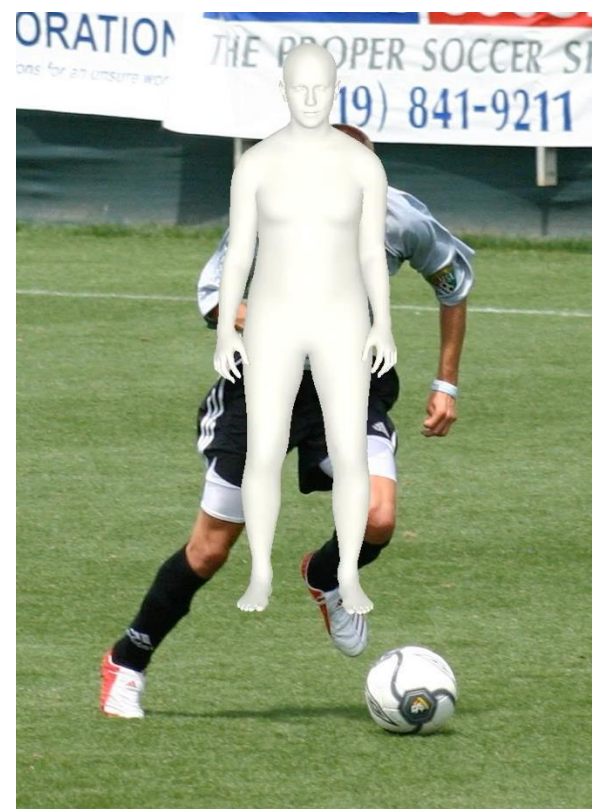
2D Keypoints

Optimize

Priors:
Pose + Shape

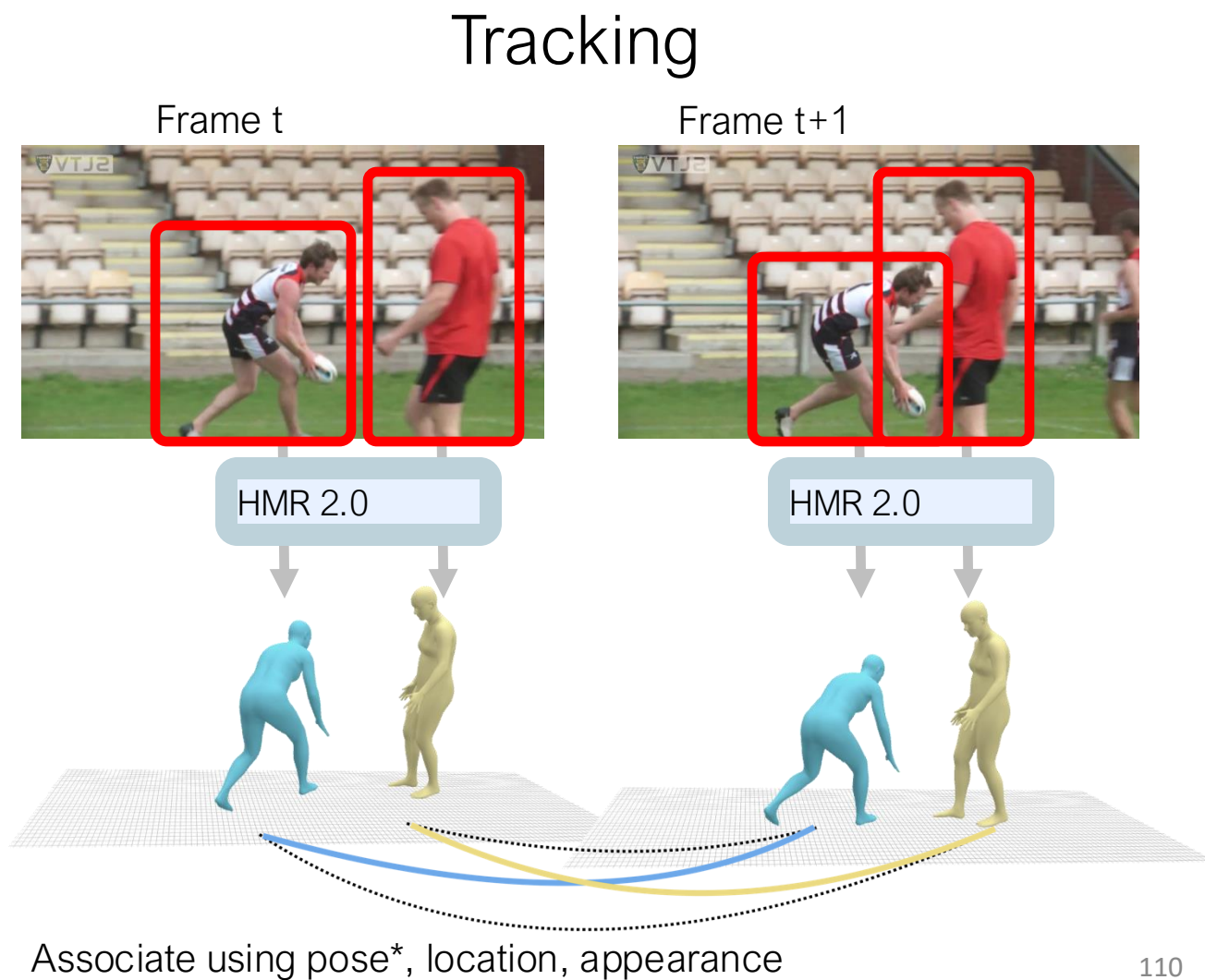
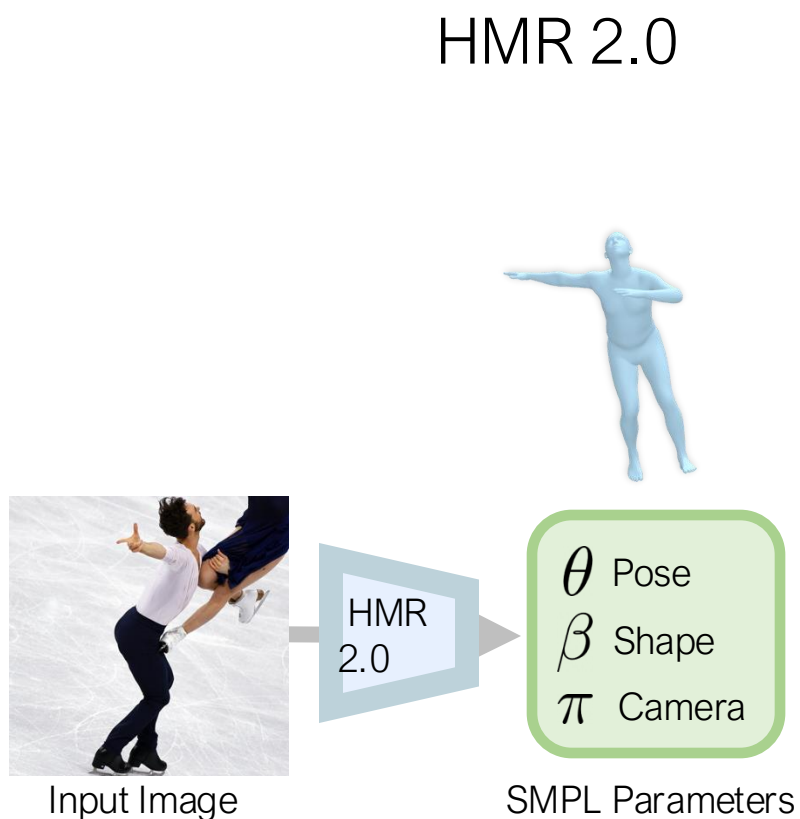
Joint Reprojection
onto keypoints

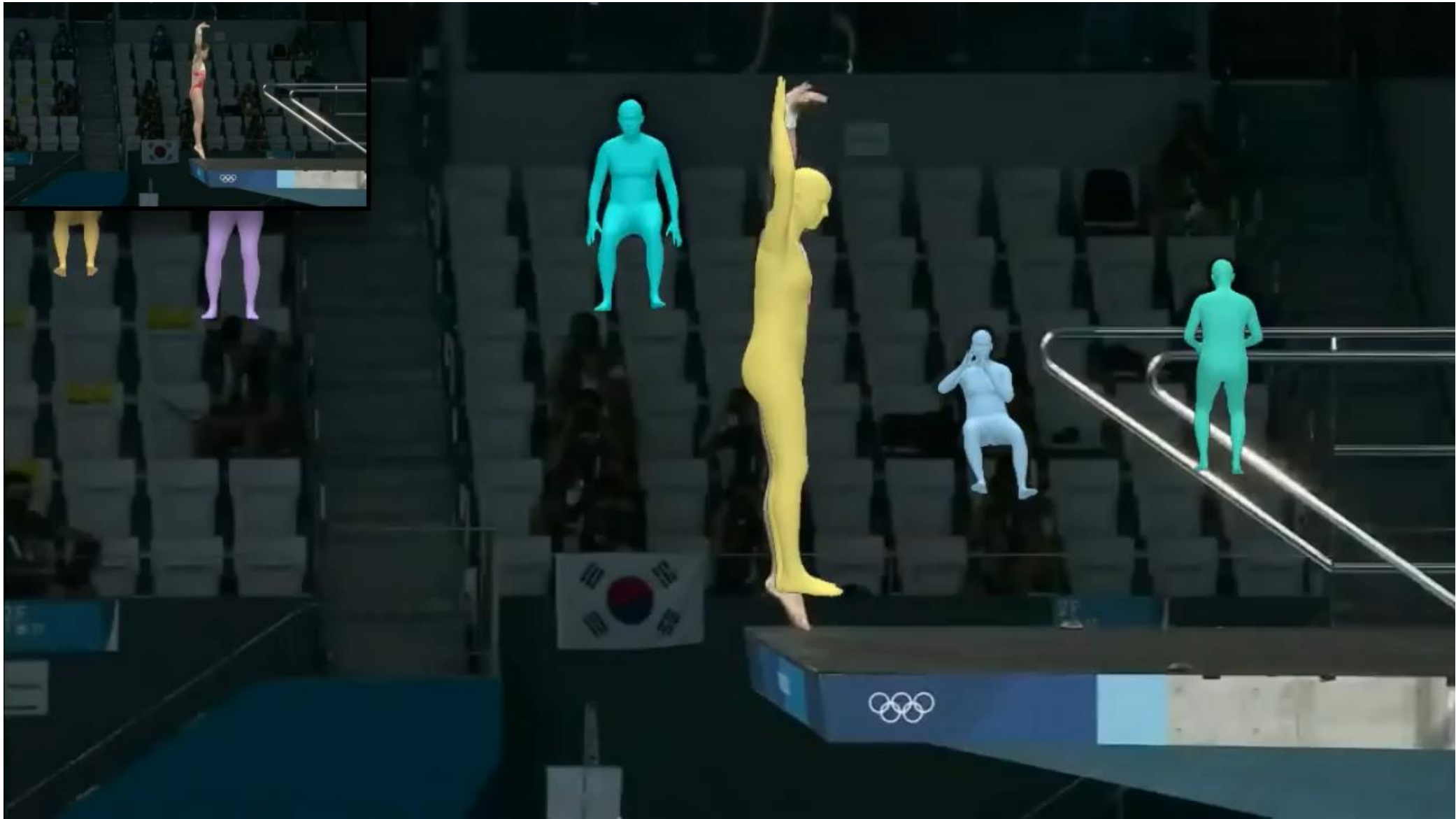
θ Pose
 β Shape
 π Camera



Finally, distill into a network!

4DHumans: HMR2.0 & PHALP++

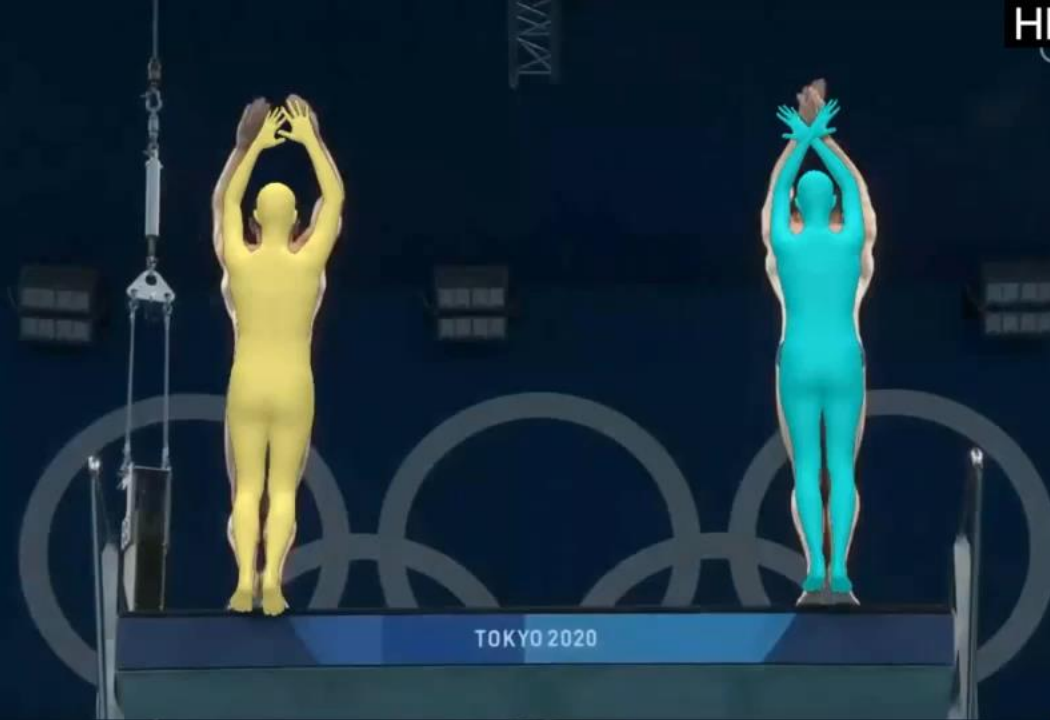






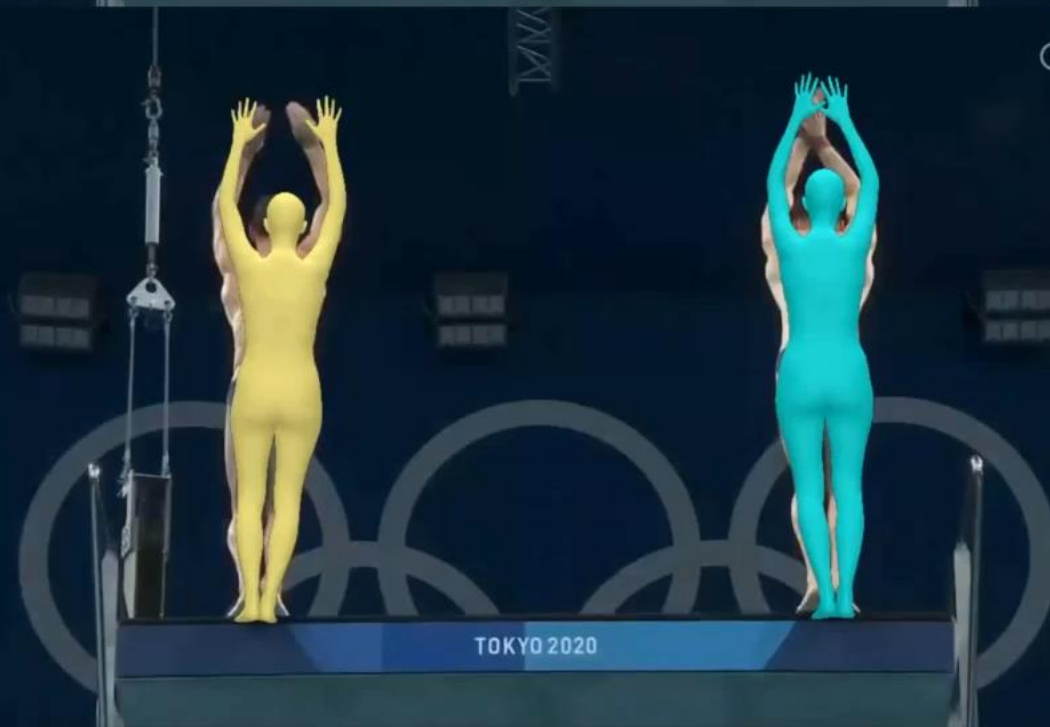
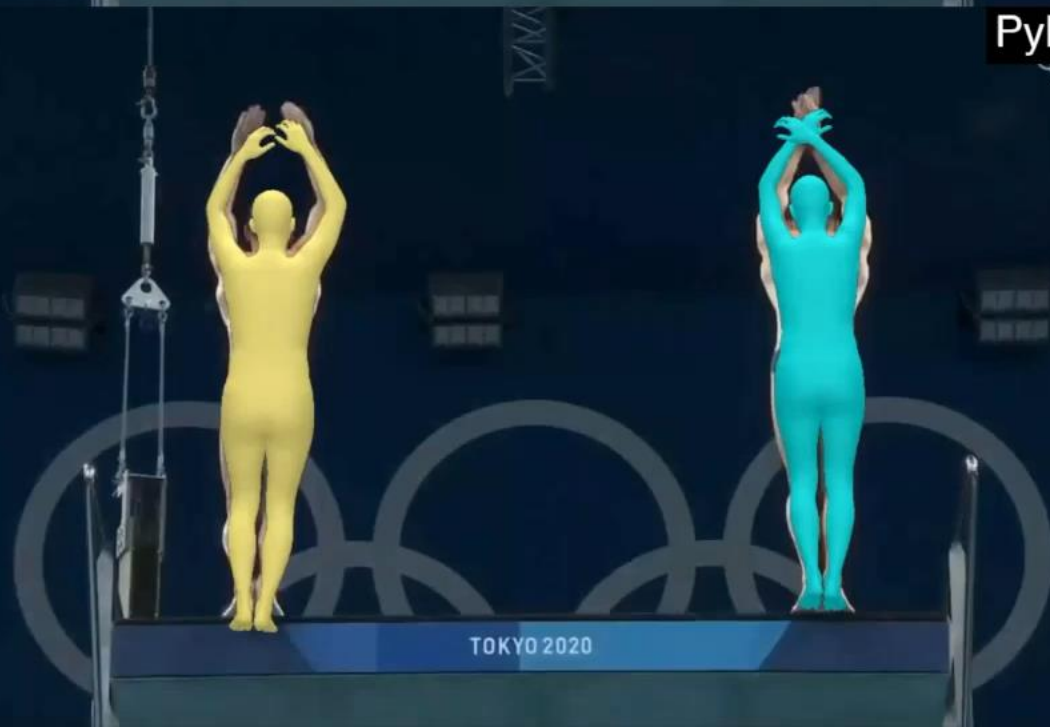
RGB Input 0

HMR 2.0



PyMAF-X 0

PARE





RGB Input

Per-frame
estimation



Camera View

Side View



HaMeR - Hand Mesh Recovery



George Pavlakos

CVPR 2024











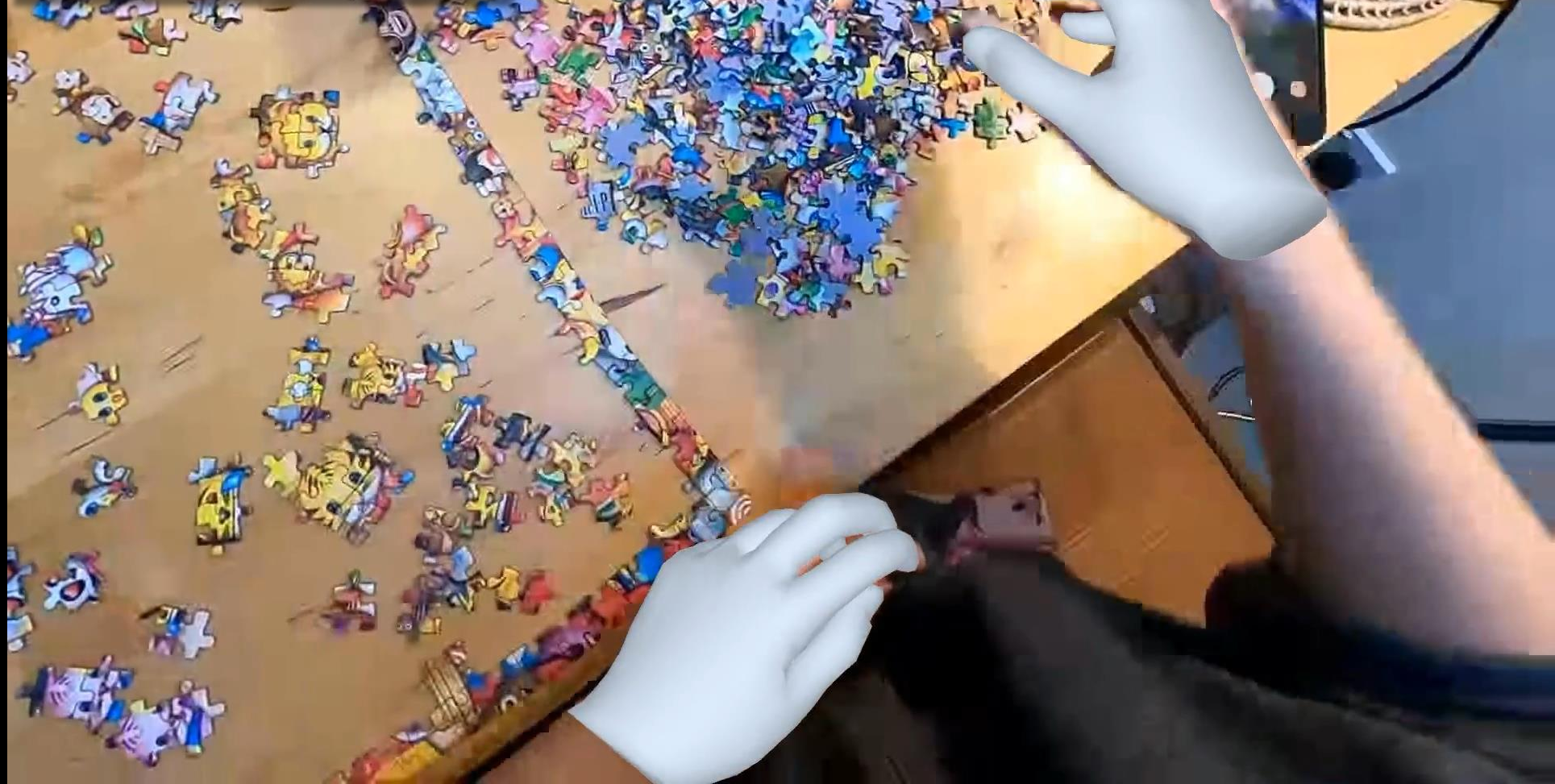
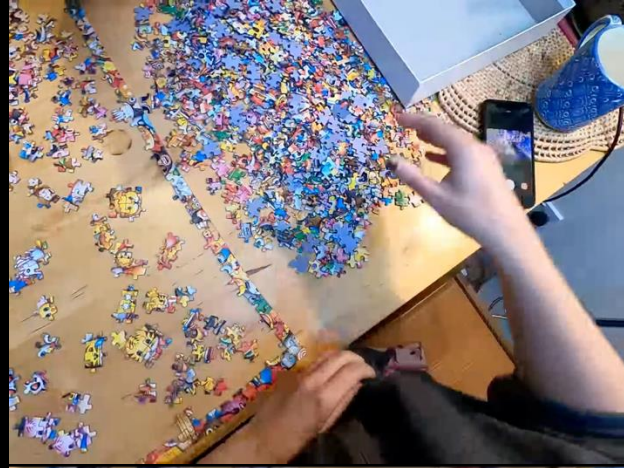


ITALIAN HAND GESTURE:

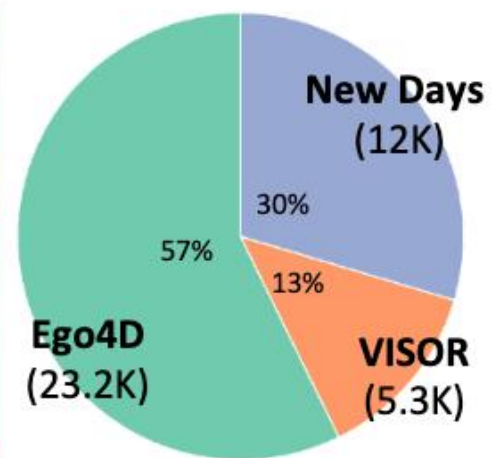
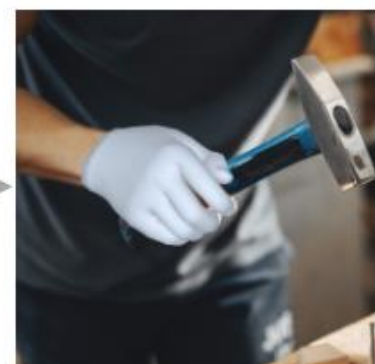
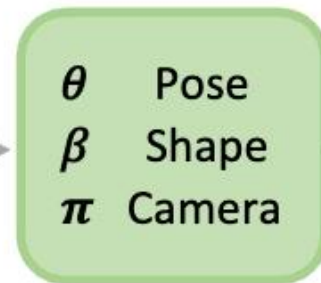


ITALIAN HAND GESTURE:









RGB Input



HaMeR



Mesh Graphormer



FrankMocap





ITALIAN HAND GESTURE:



ITALIAN HAND GESTURE:

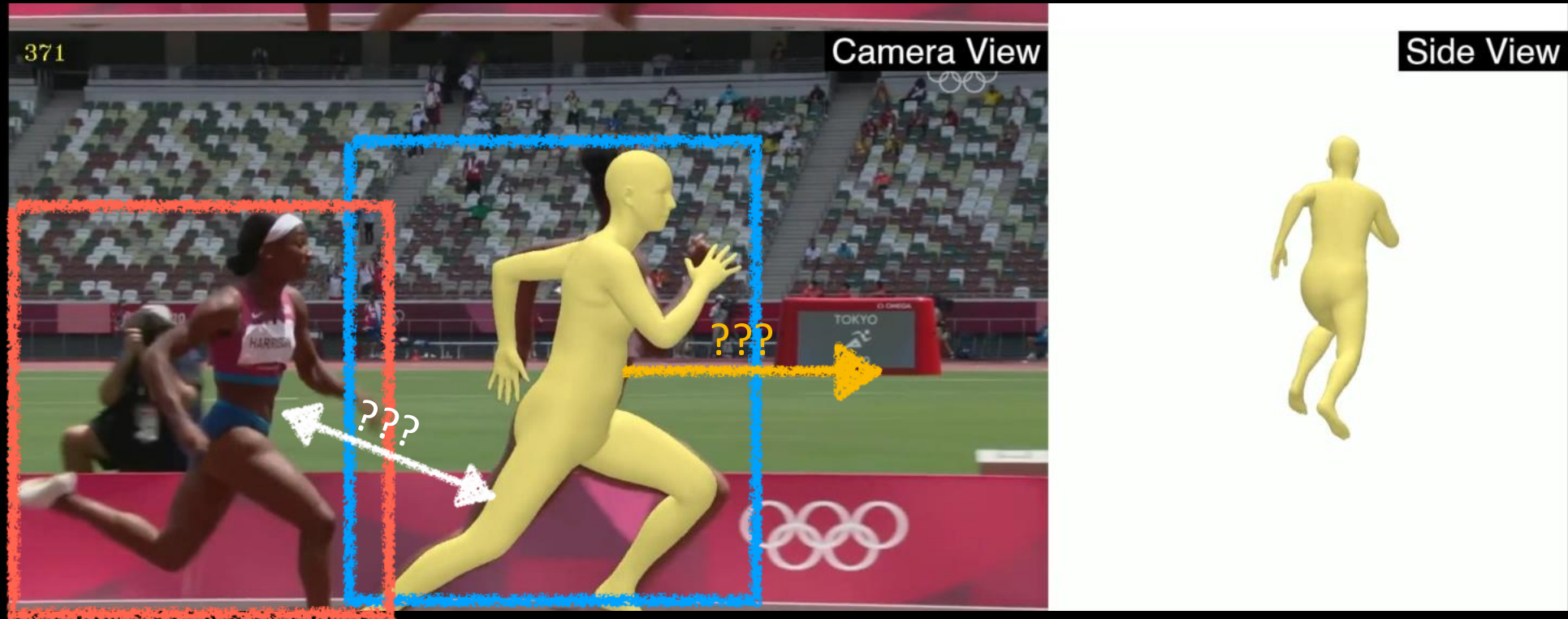
Right hand
Top view



Left hand
Top view



Caveat: Local Pose



Decoupling Human and Camera Motion from Videos in the Wild



Vickie Ye



Georgios Pavlakos



Jitendra Malik



Angjoo Kanazawa

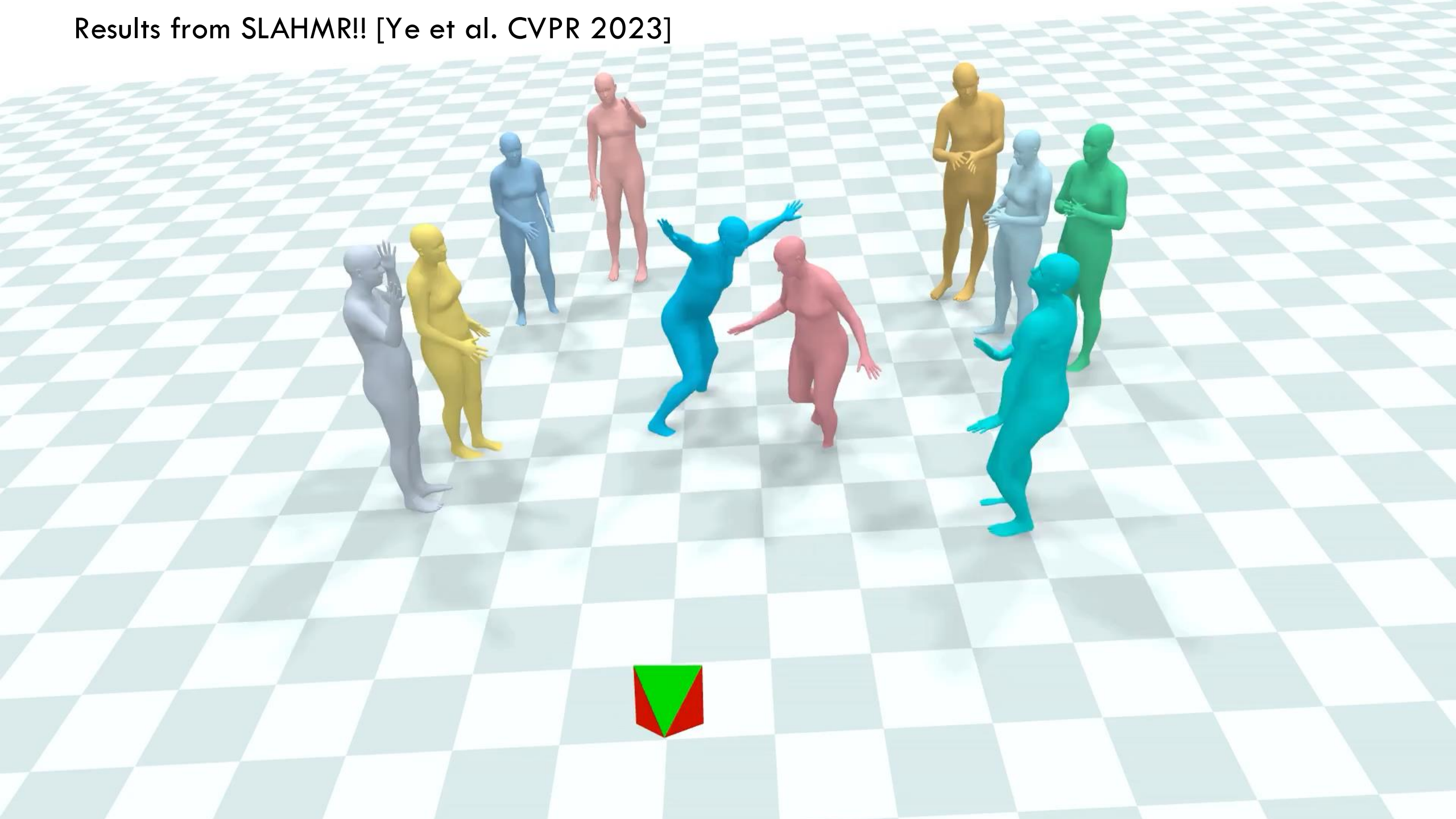
CVPR 2023

UC Berkeley

We live in a world that is 3D and dynamic.

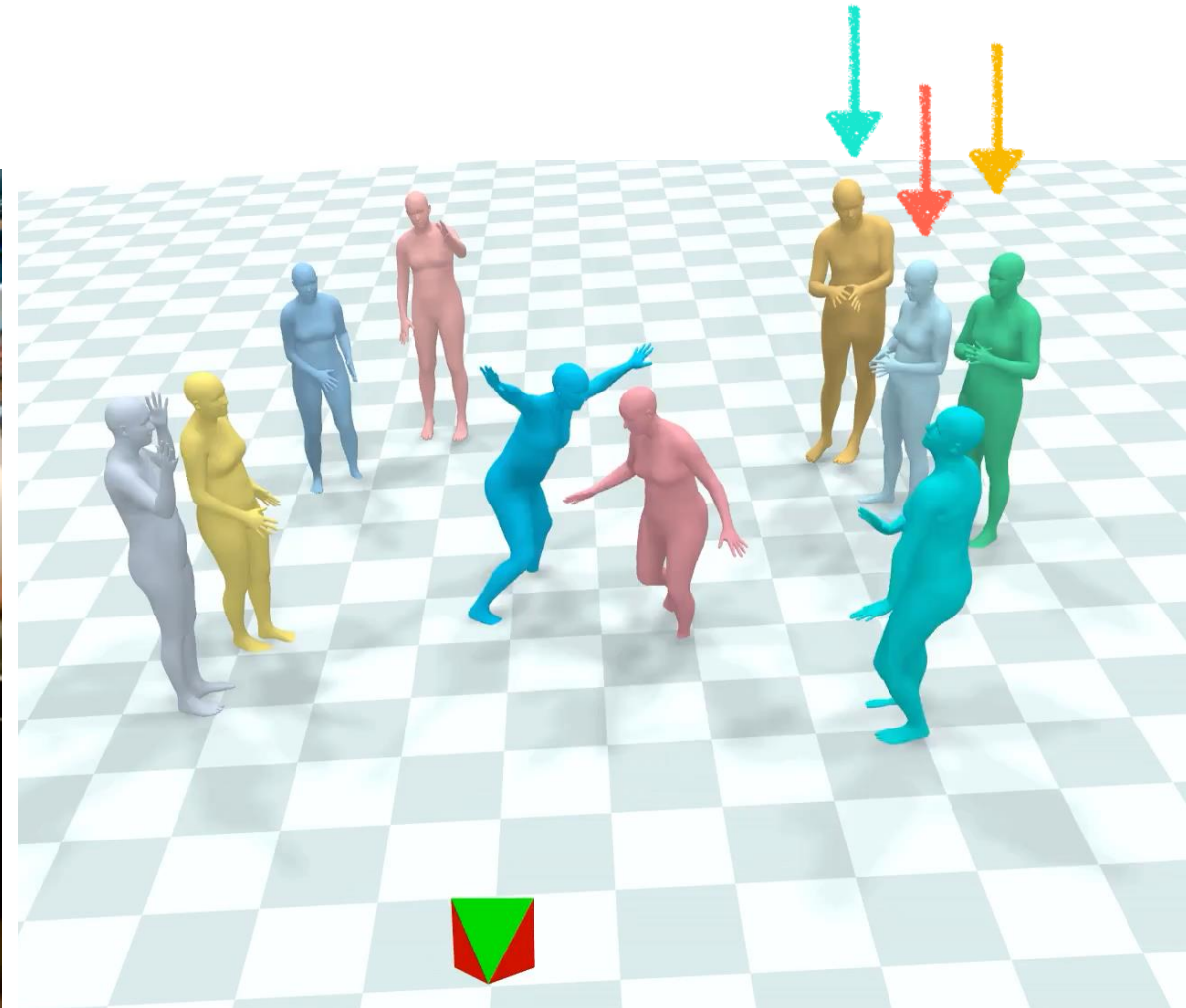
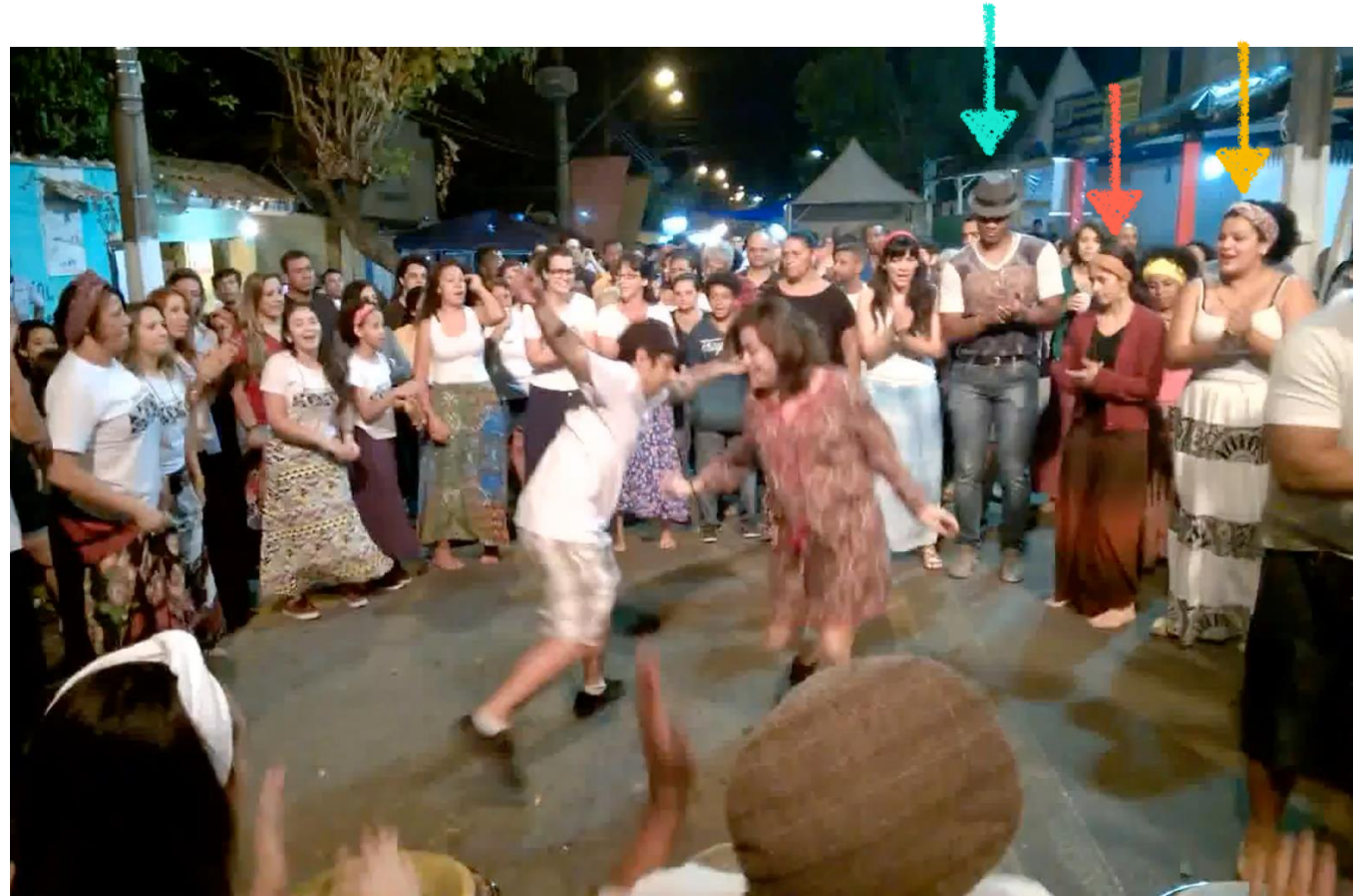


Results from SLAHMR!! [Ye et al. CVPR 2023]







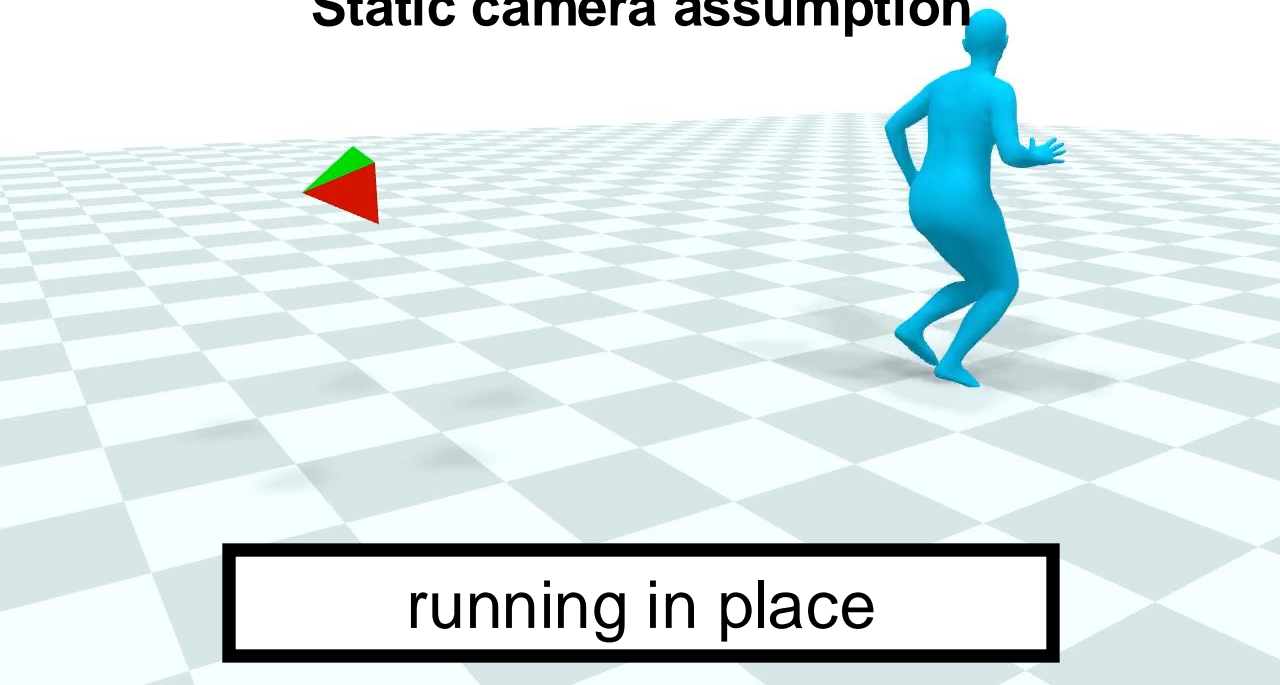


SCTV



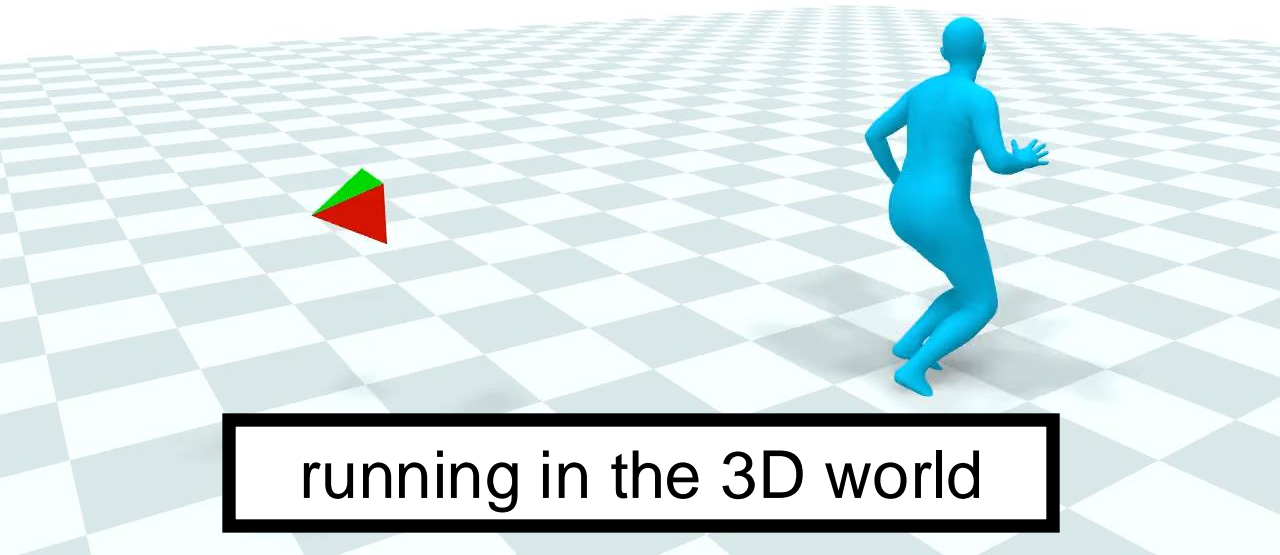


Static camera assumption



running in place

Modeling camera motion



running in the 3D world

Input:



Output:

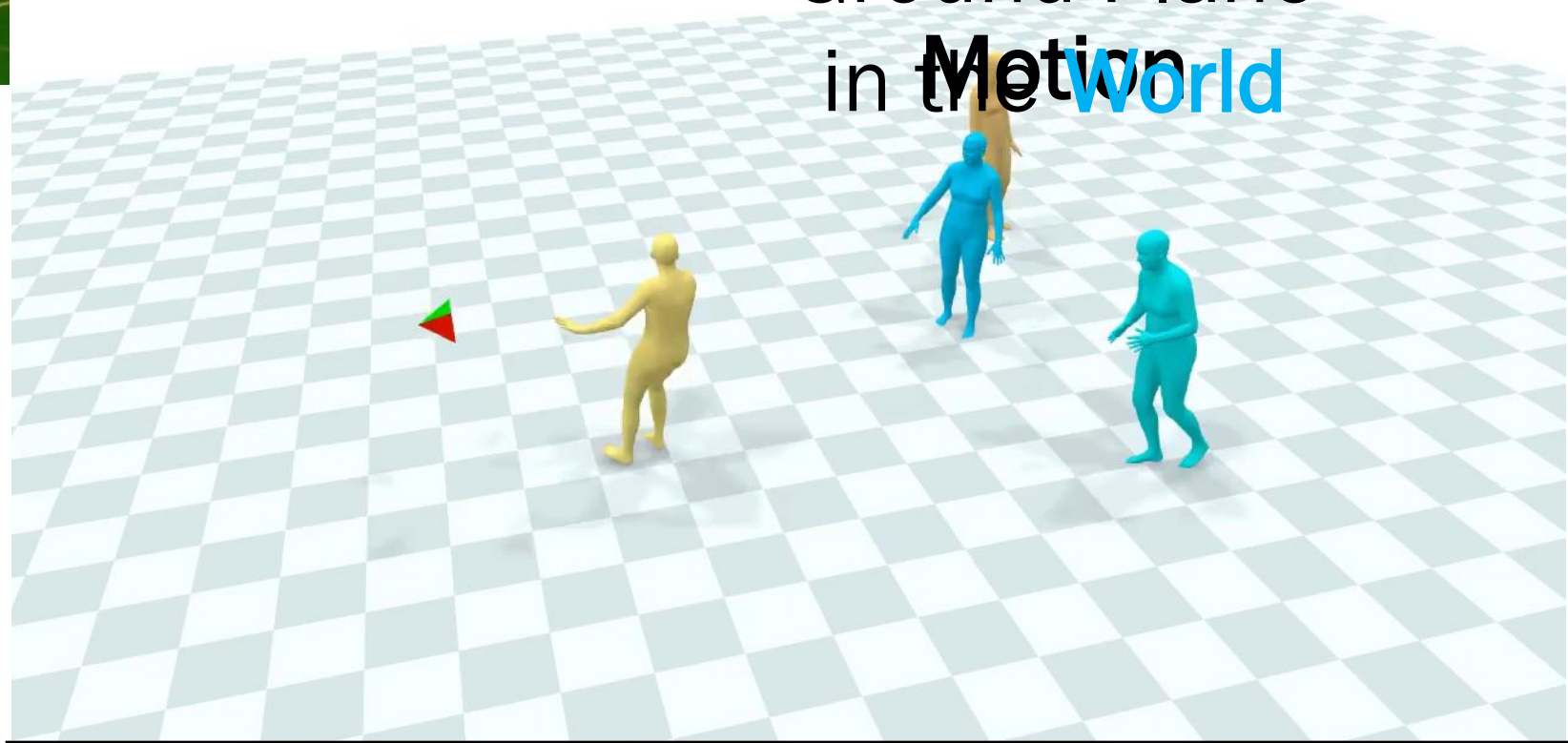
$$\{\text{world } \mathbf{P}\}$$

$$\{R, \alpha T\}$$

$$\{g\}$$

Tracked People
in **World** Frame

Camera
Ground Plane
in **World**



$$\{\text{cam } \mathbf{P}\}$$

$$\{R, T\}$$

Tracked
People in
Unscale
Camera
Frame
Camera
Motion



SLAHMR: Simultaneous Localization and Human Mesh Recovery

Key signal: Motion Prior

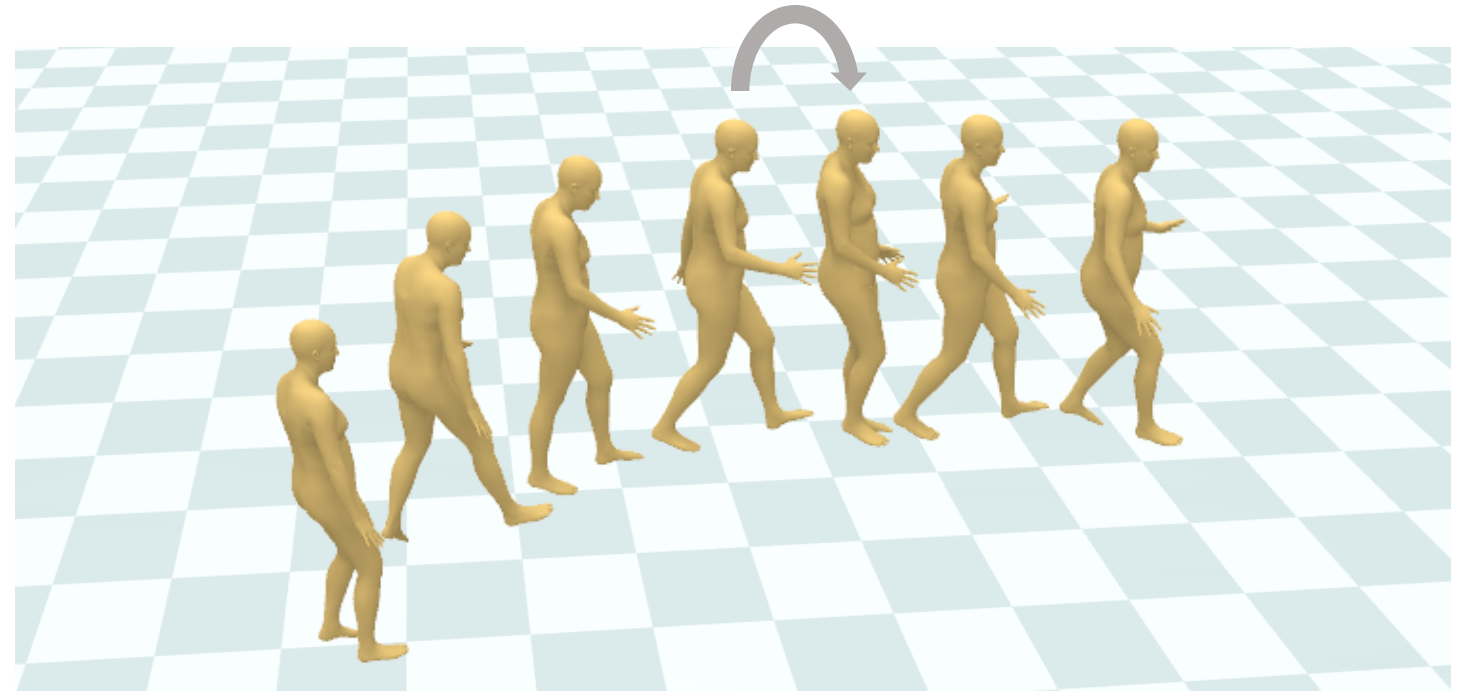


Recover World that gives most probable motion

A data-driven motion prior: HuMoR [Rempe et al. ICCV 2021]



$$p(\text{Human} | \text{Human})$$



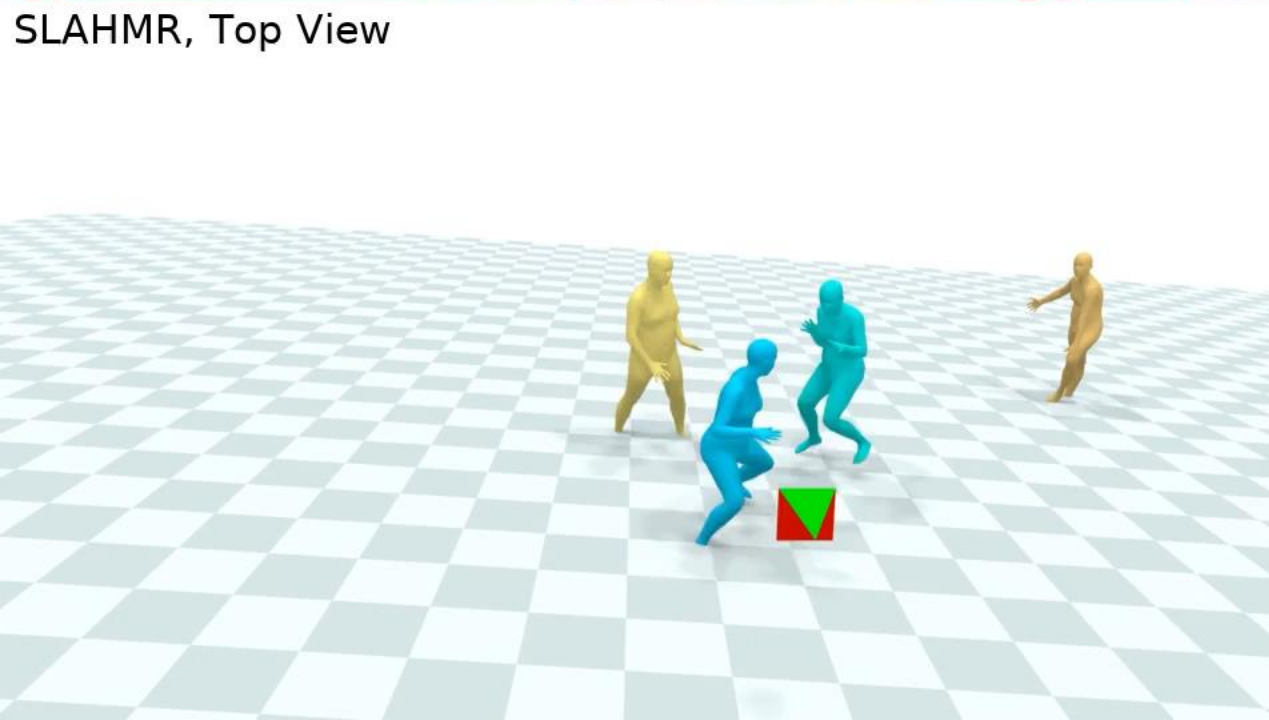
Input View



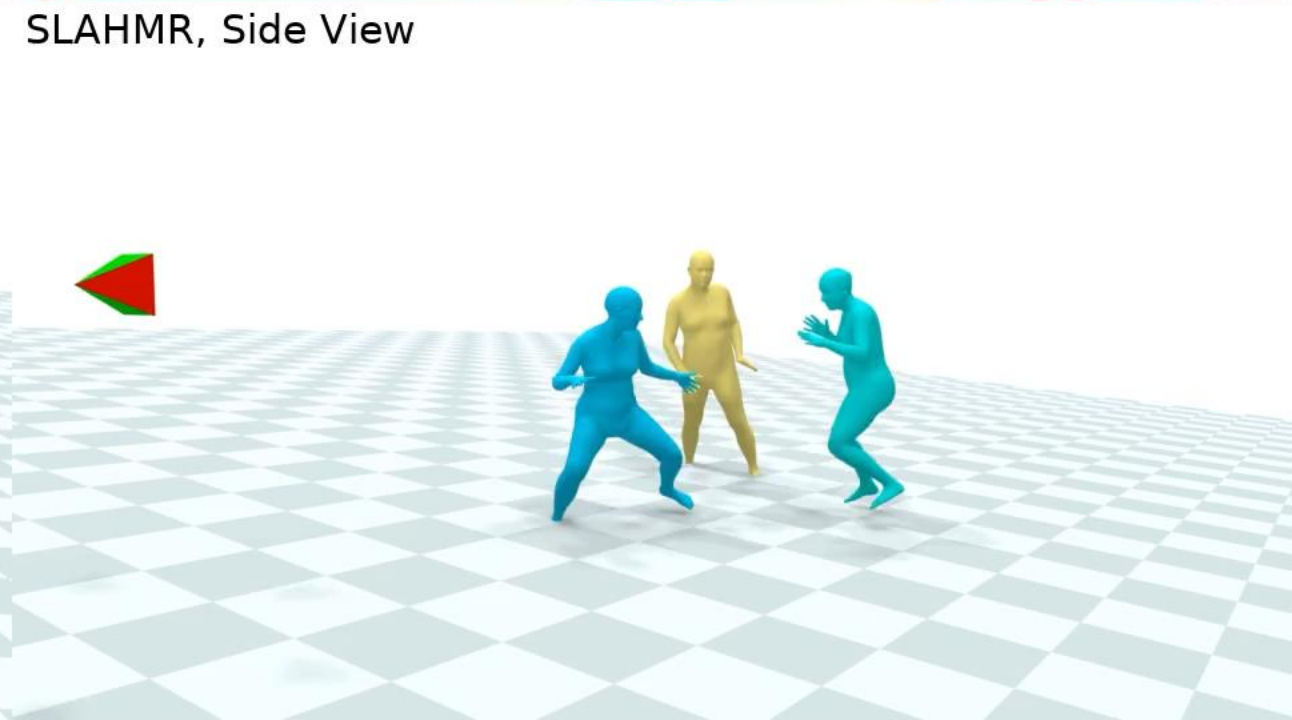
SLAHMR, Input View



SLAHMR, Top View



SLAHMR, Side View



Input View



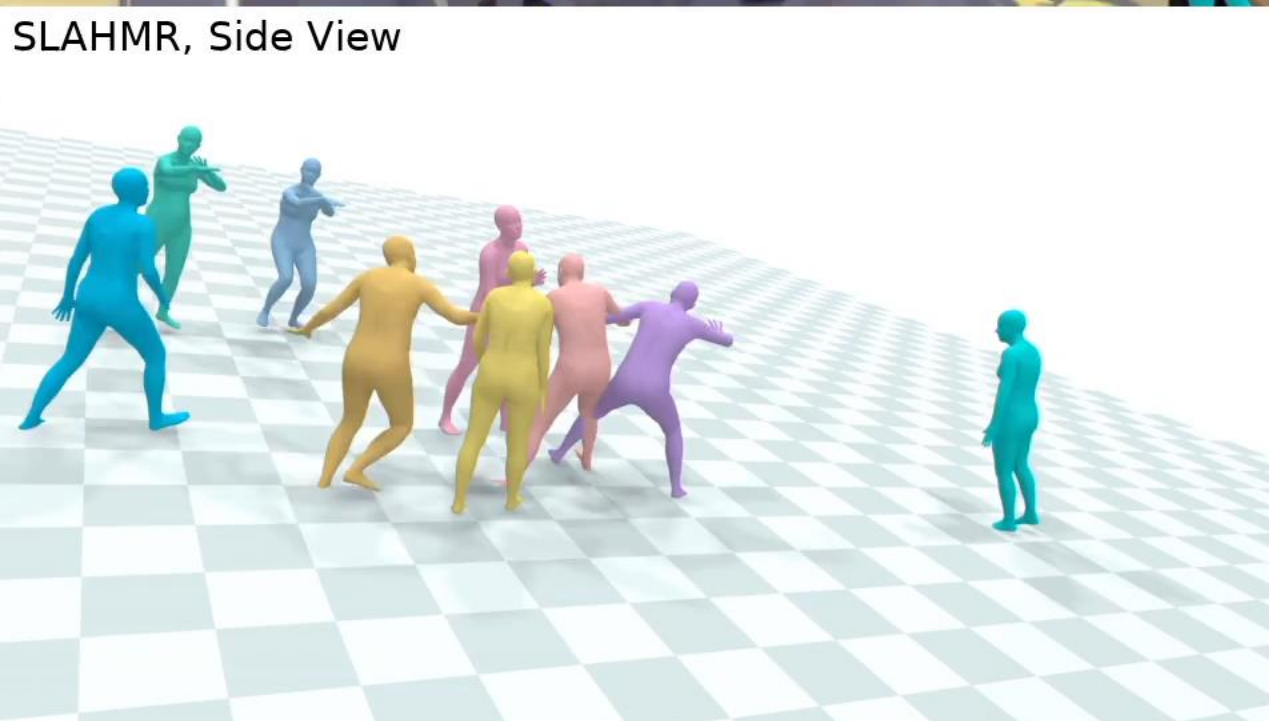
SLAHMR, Input View



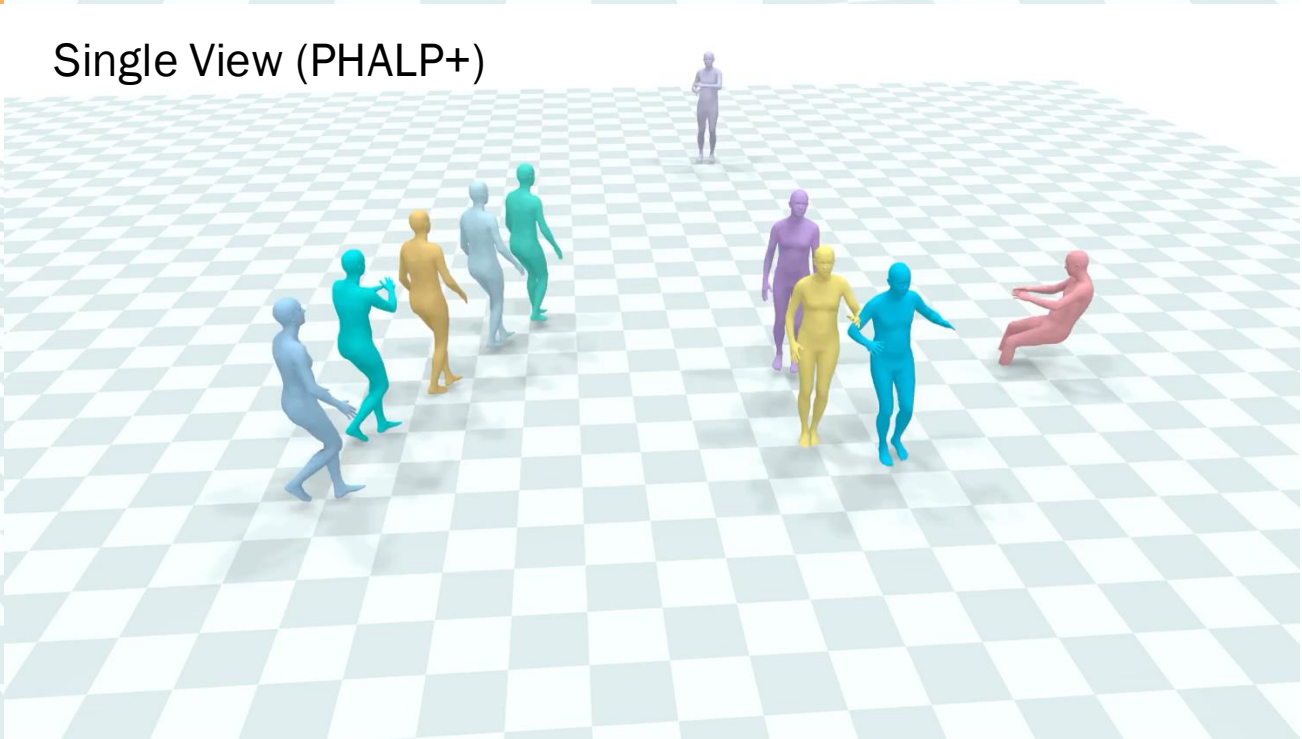
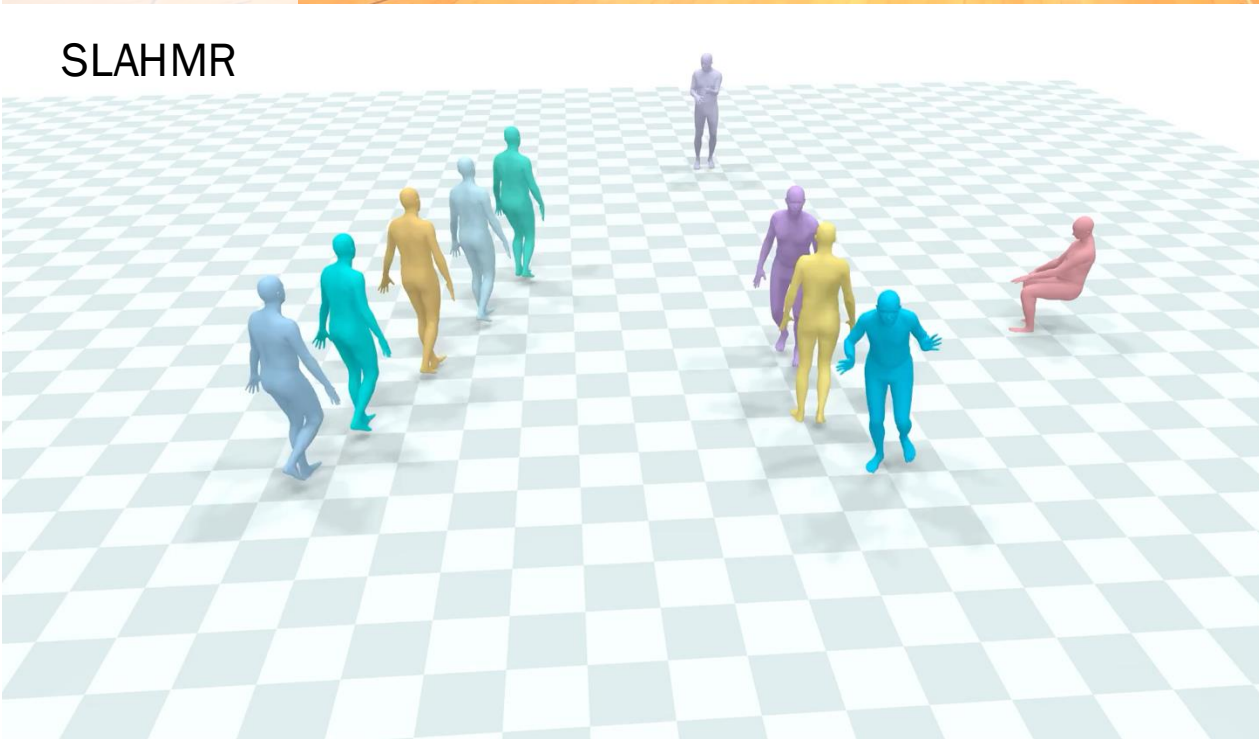
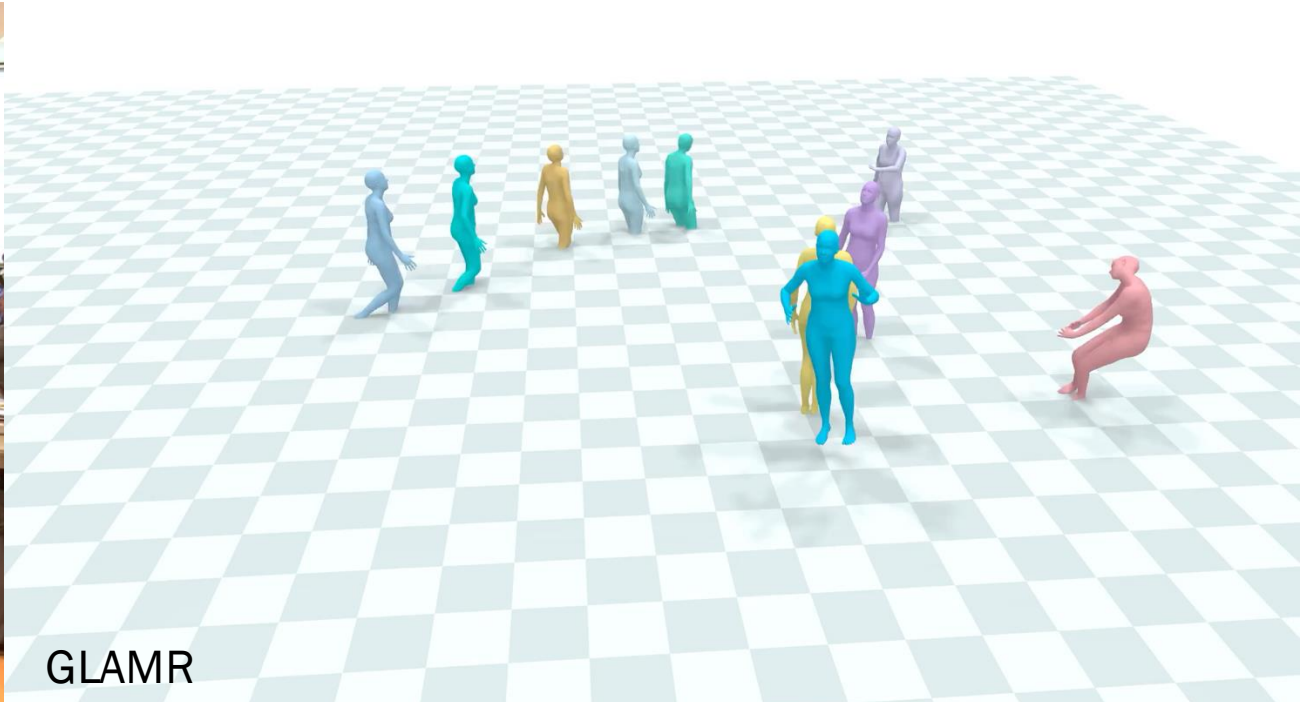
SLAHMR, Top View



SLAHMR, Side View



Comparisons

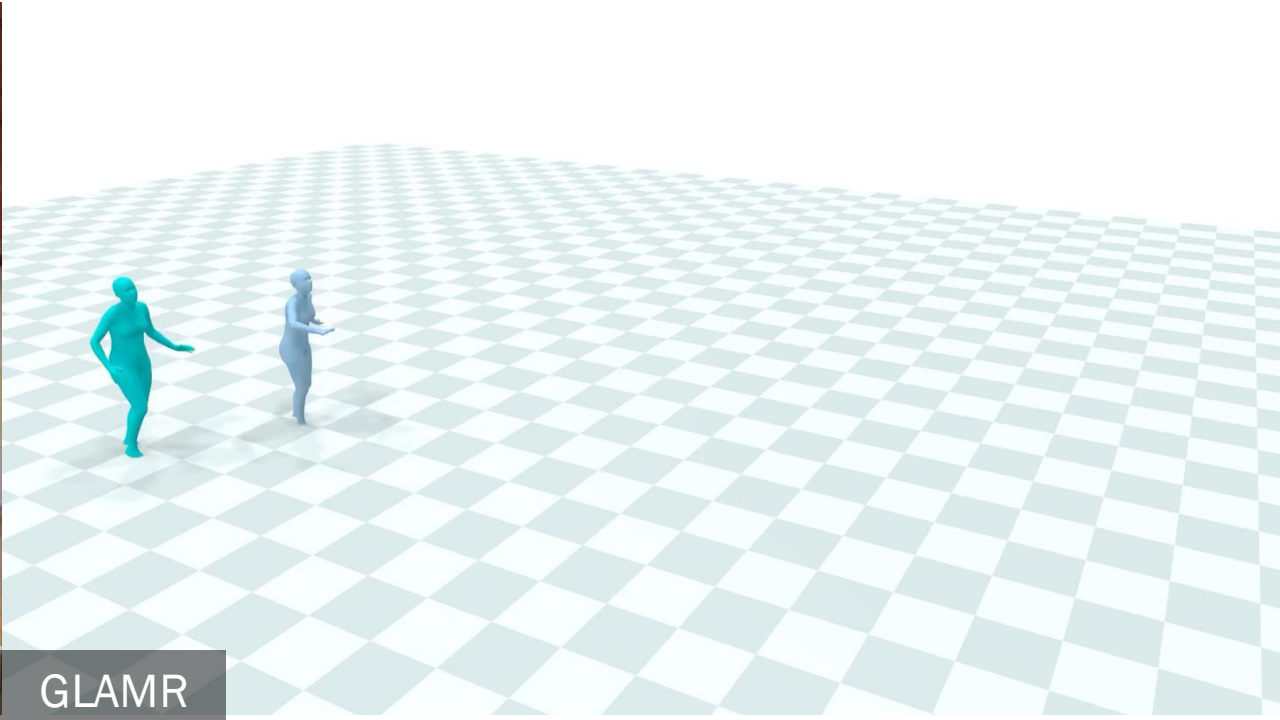




Comparisons

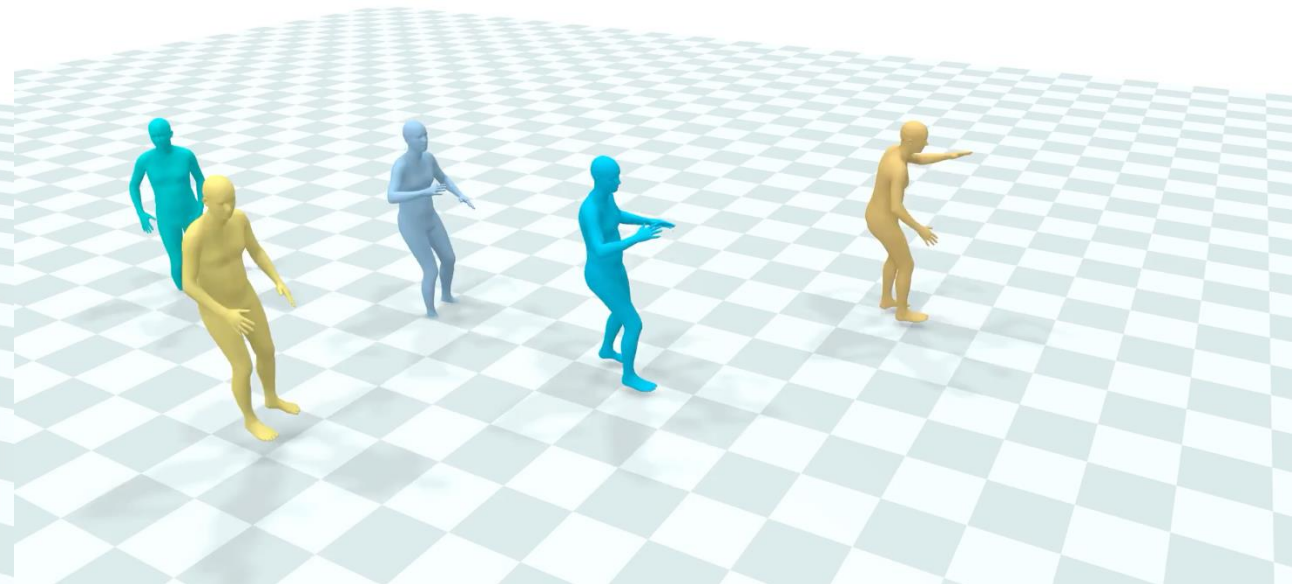
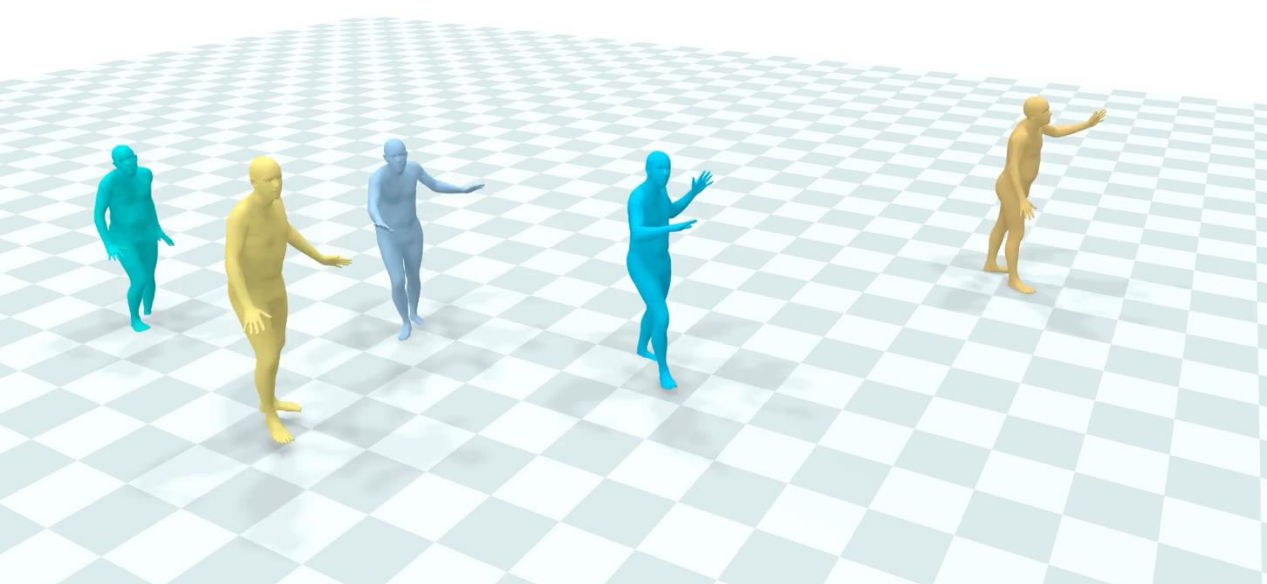
Input View

SLAHMR



GLAMR

Single View (PHALP+)

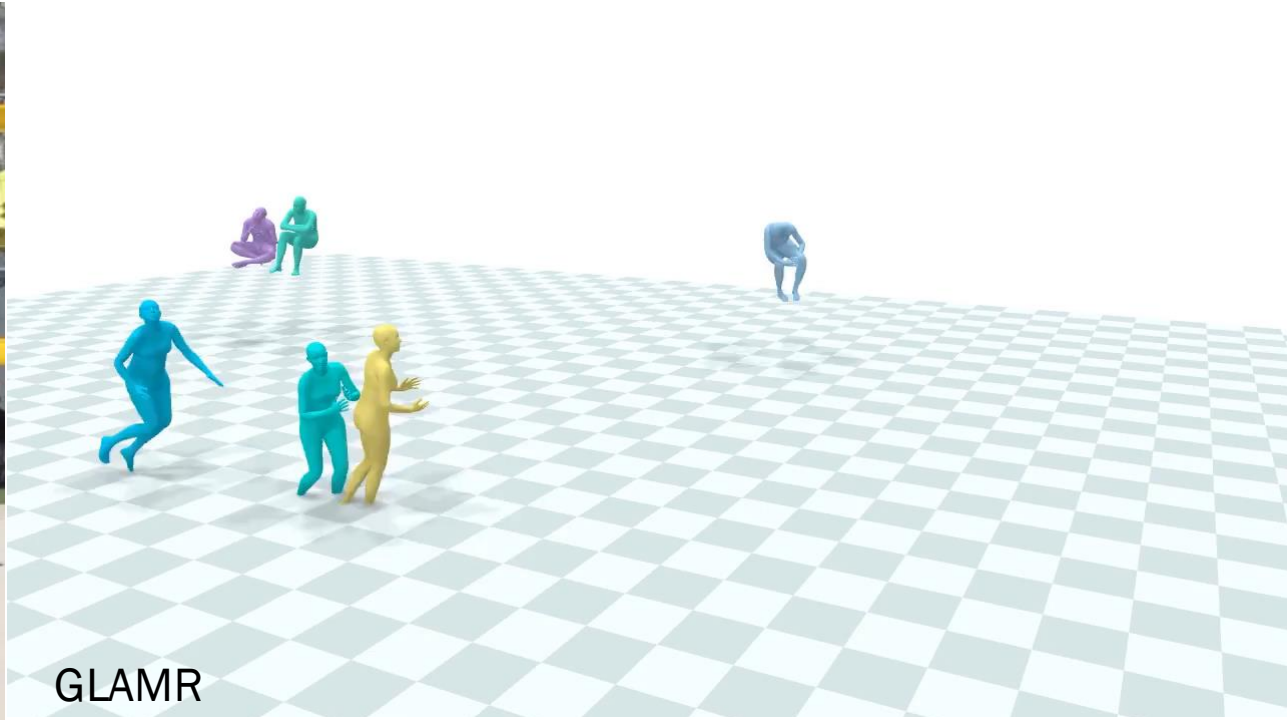


Comparisons

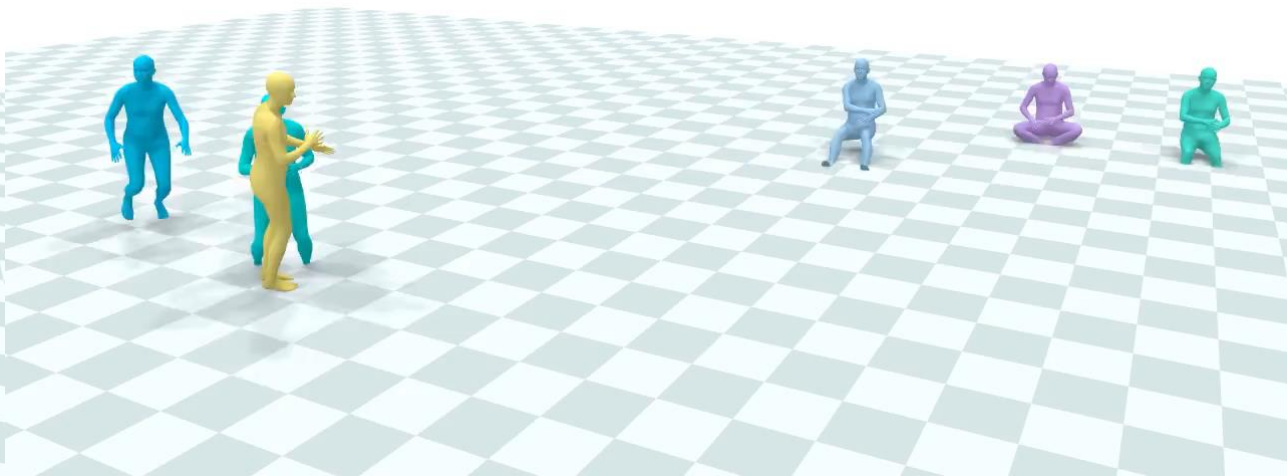
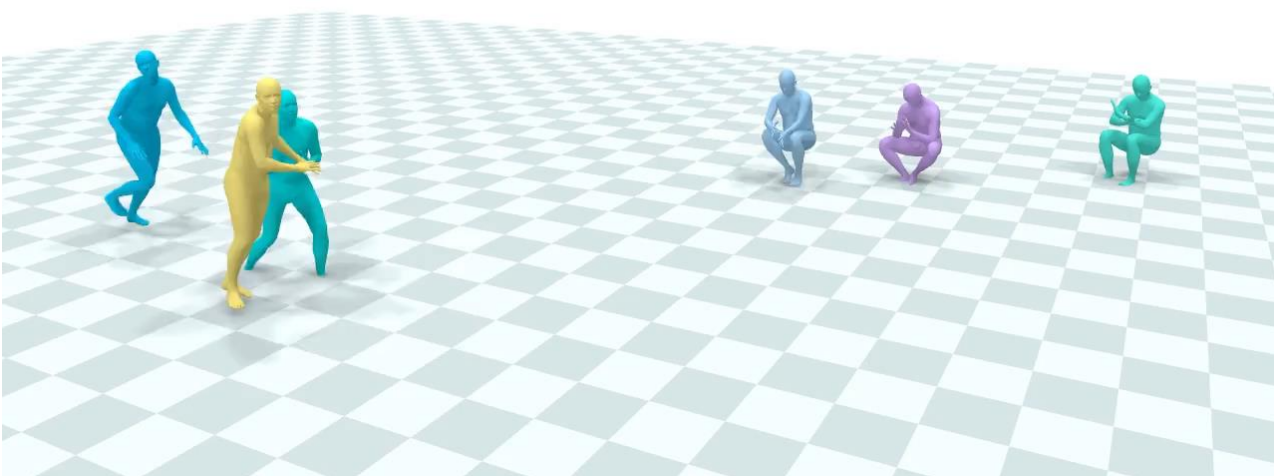


Input View

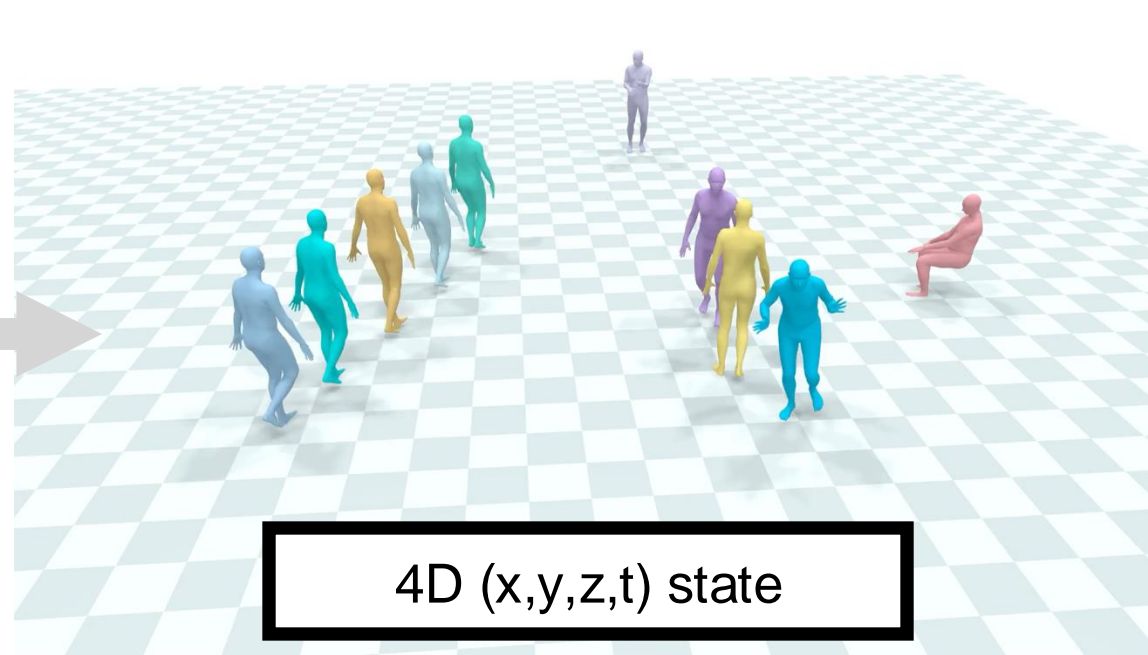
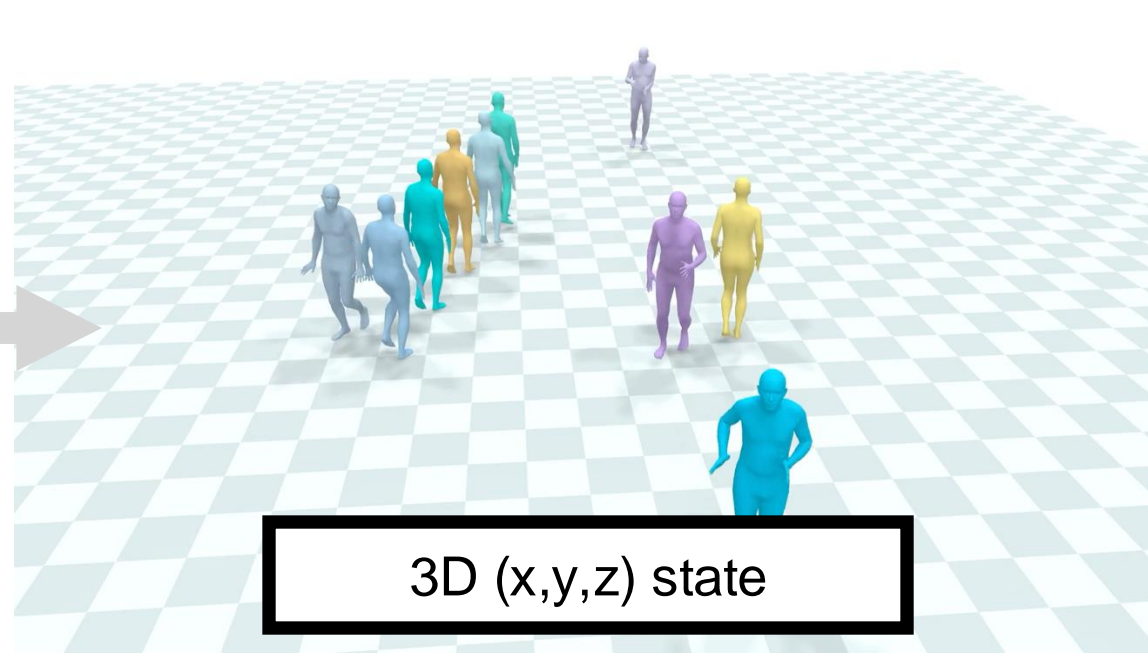
Ours



GLAMR
Single View (PHALP+)



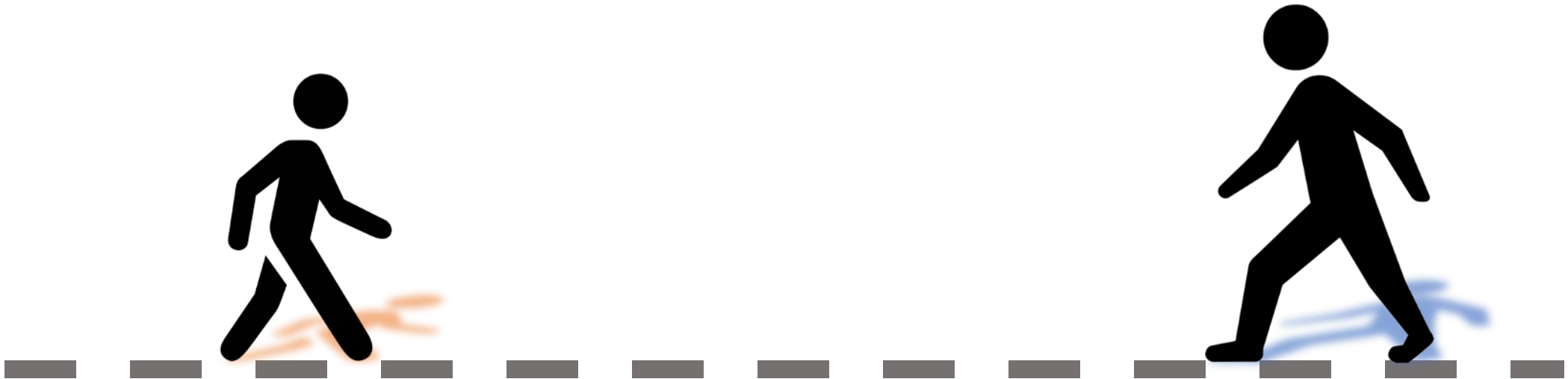
4D Reconstruction



Prereq: tracking people across time



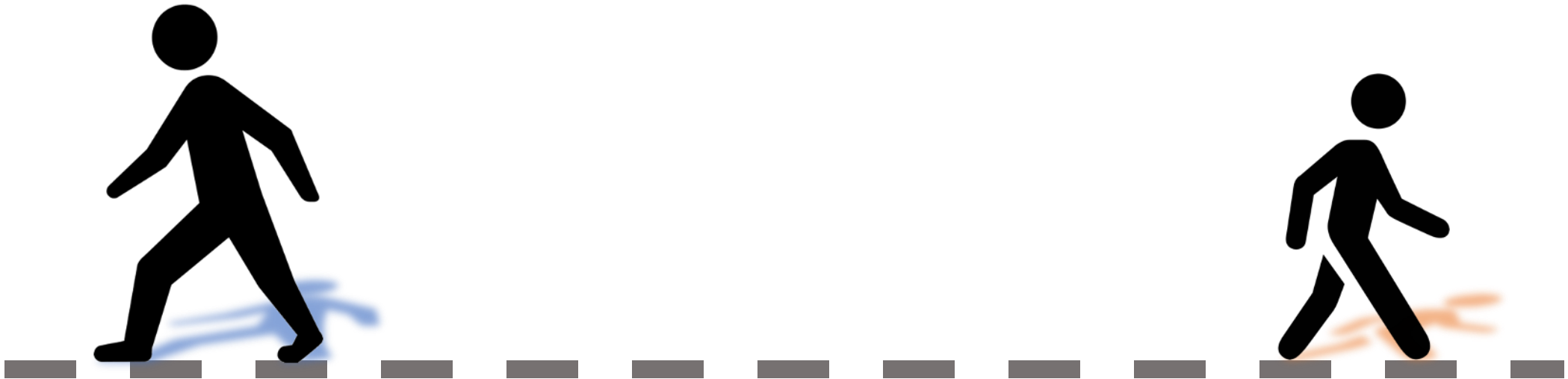
- In 2D, occlusion is hard to disambiguate



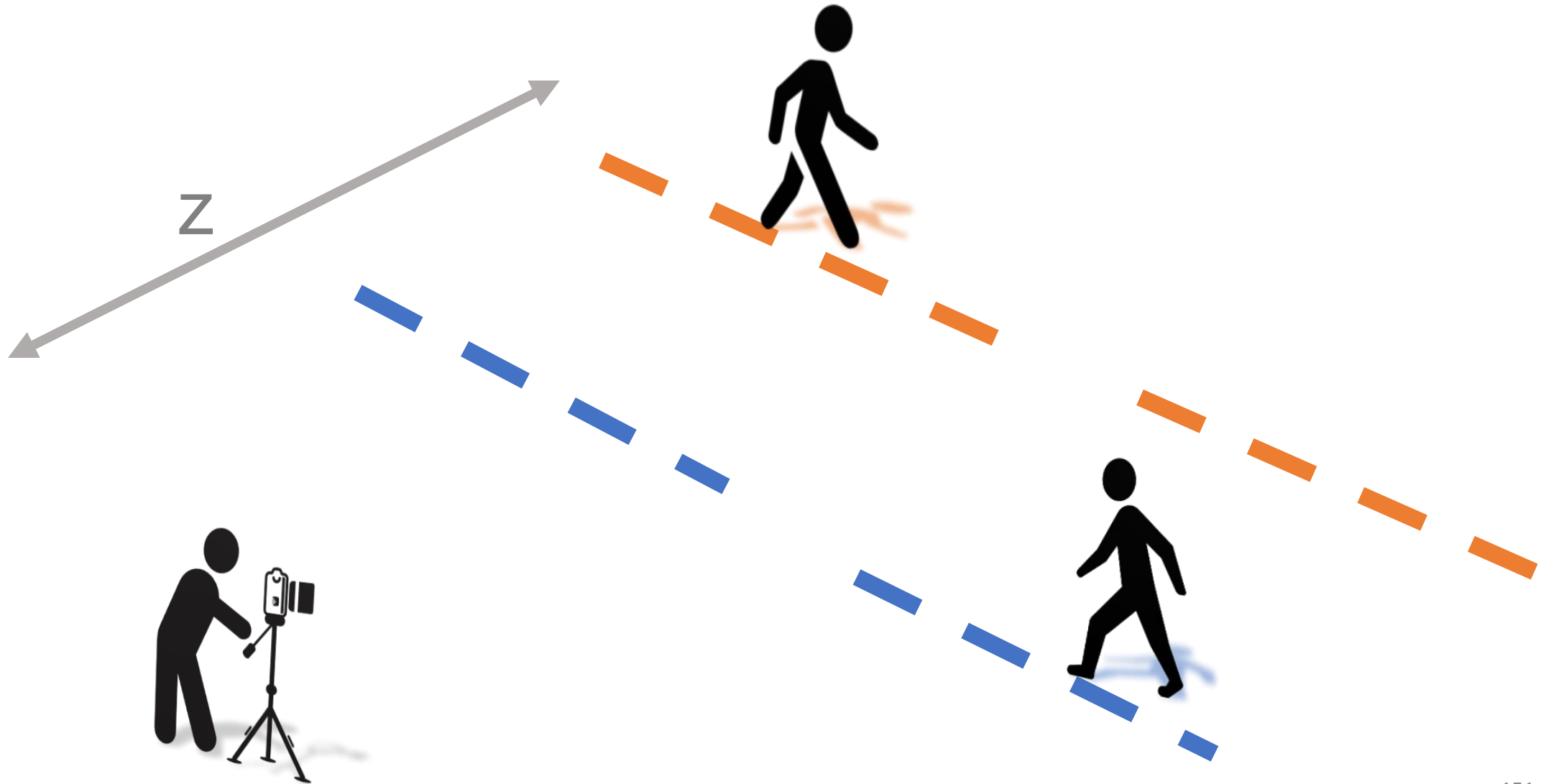
- In 2D, occlusion is hard to disambiguate



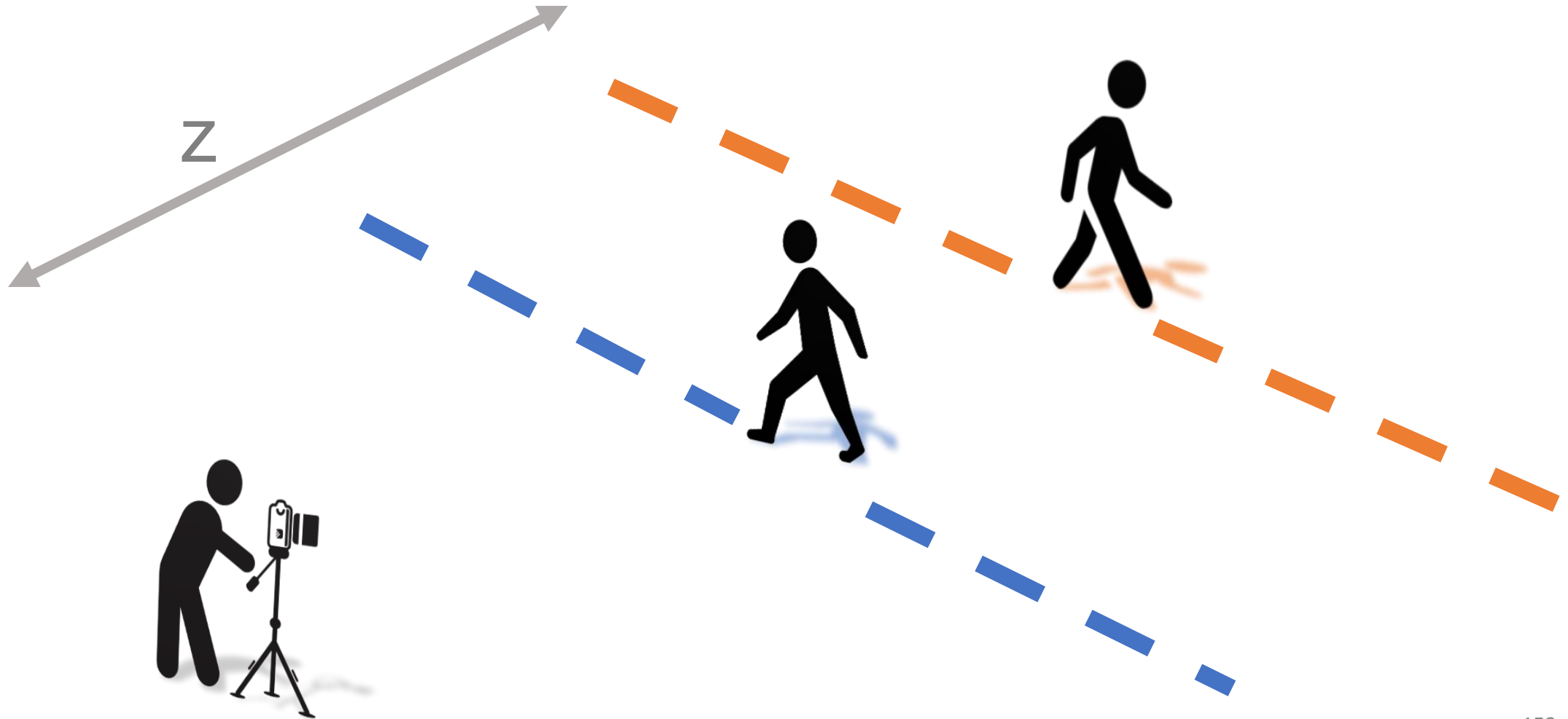
- In 2D, occlusion is hard to disambiguate



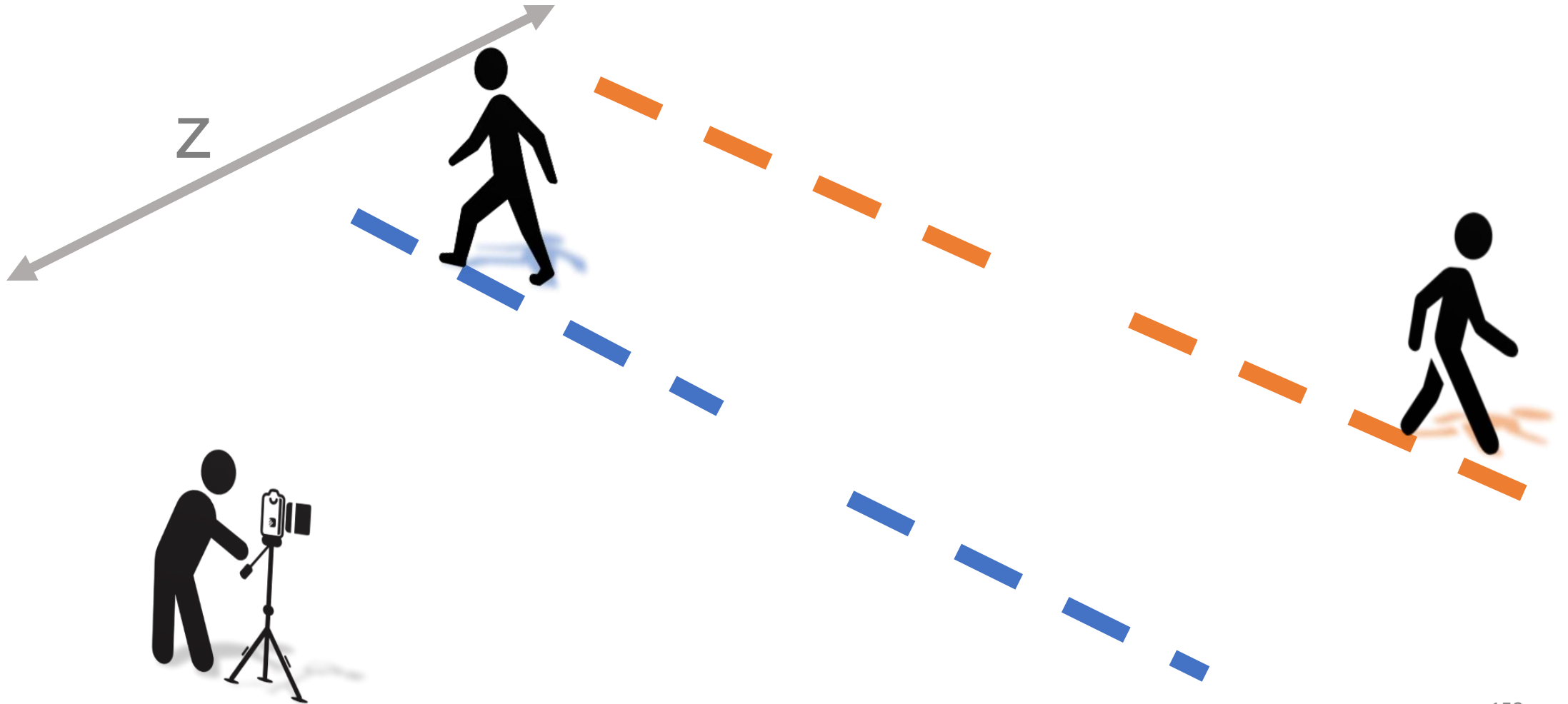
- In 2D, they overlap, but in 3D they don't!



- In 2D, they overlap, but in 3D they don't!



- In 2D, they overlap, but in 3D they don't!



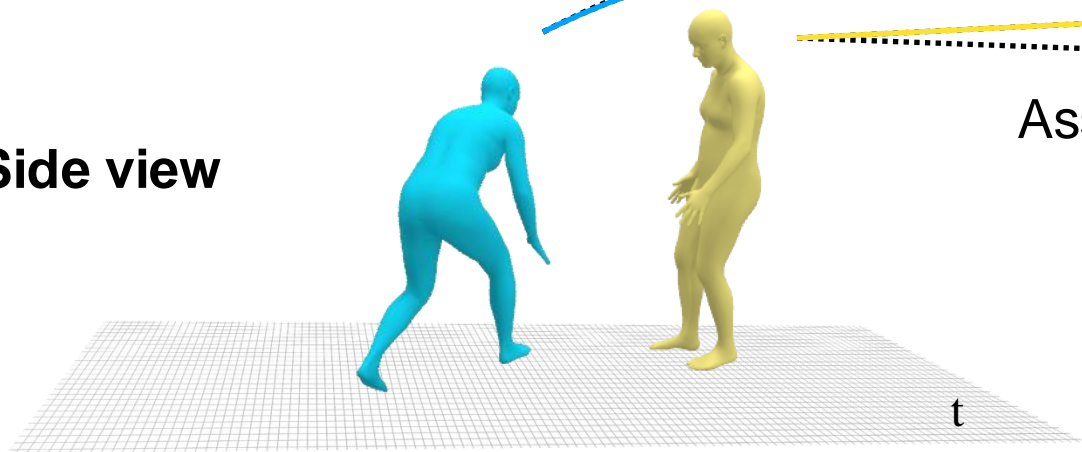


t



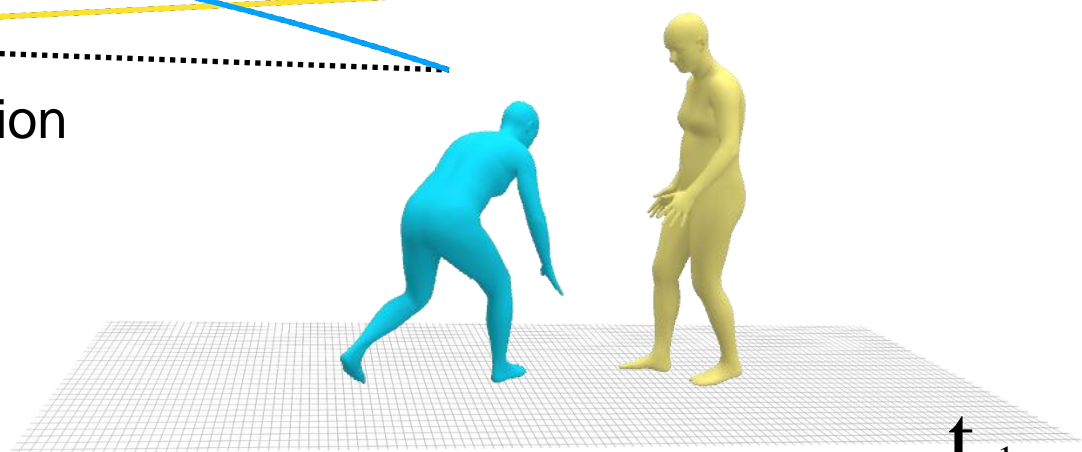
t_{+1}

Side view



t

Association



t_{+1}

Distance

=

3D location distance

+

3D appearance distance

+

3D pose distance

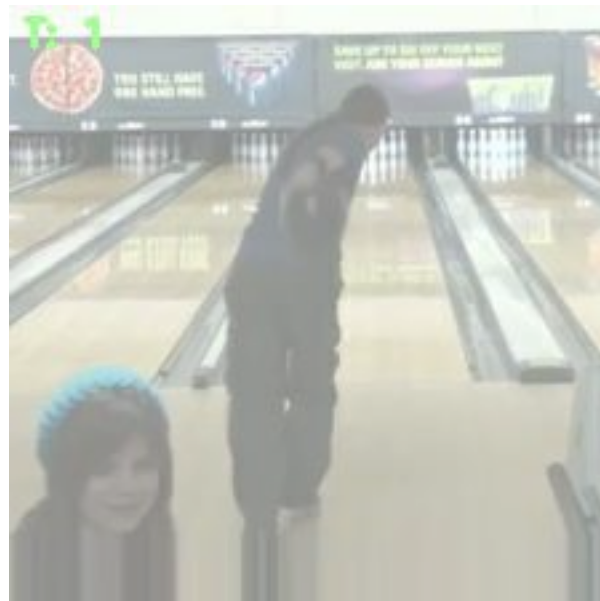
A benefit of video: Dynamics



Auto-regressive prediction of 3D motion from video



Input
Video



Ground
Truth
Video



Predicted
Future



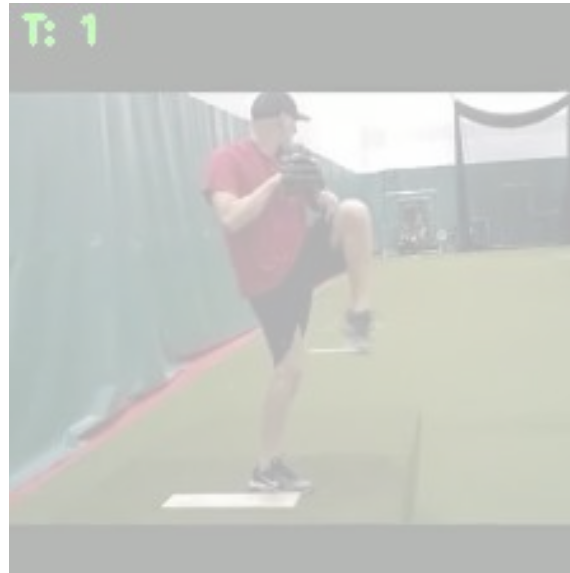
Different
Viewpoint



Test Time



Input
Video



Ground
Truth
Video

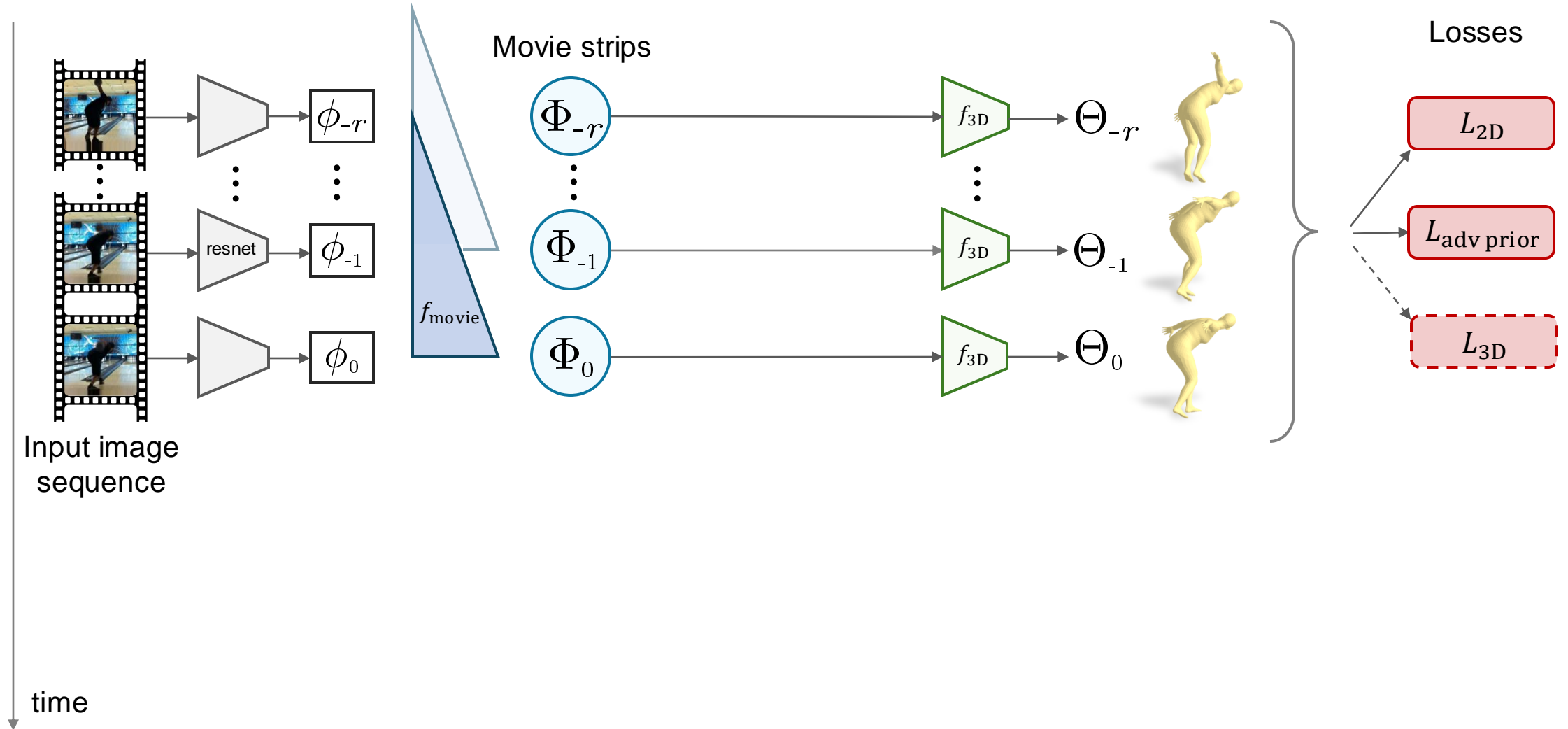


Predicted
Future

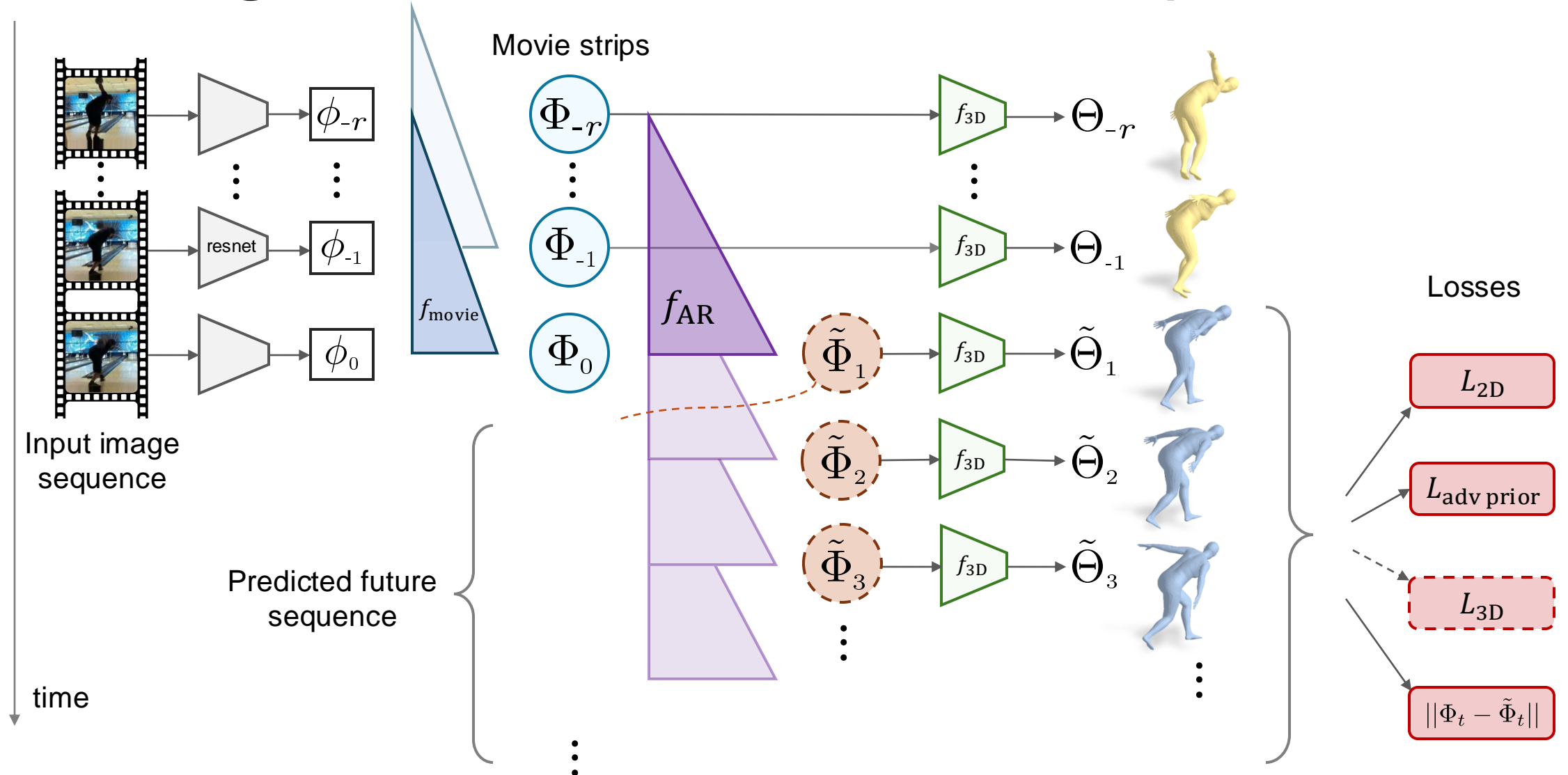


Different
Viewpoint

Overview



Autoregressive Prediction (Latent Space)



Yellow = Conditioning

Blue = Future Prediction from movie strip



Camera View



Alternate Viewpoint

Yellow = Conditioning

Blue = Future Prediction from movie strip



Camera View



Alternate Viewpoint

Yellow = Conditioning

Blue = Future Prediction from movie strip



Camera View



Alternate Viewpoint

Finally, a step towards this baby



SfV: Reinforcement Learning of Skills from Videos

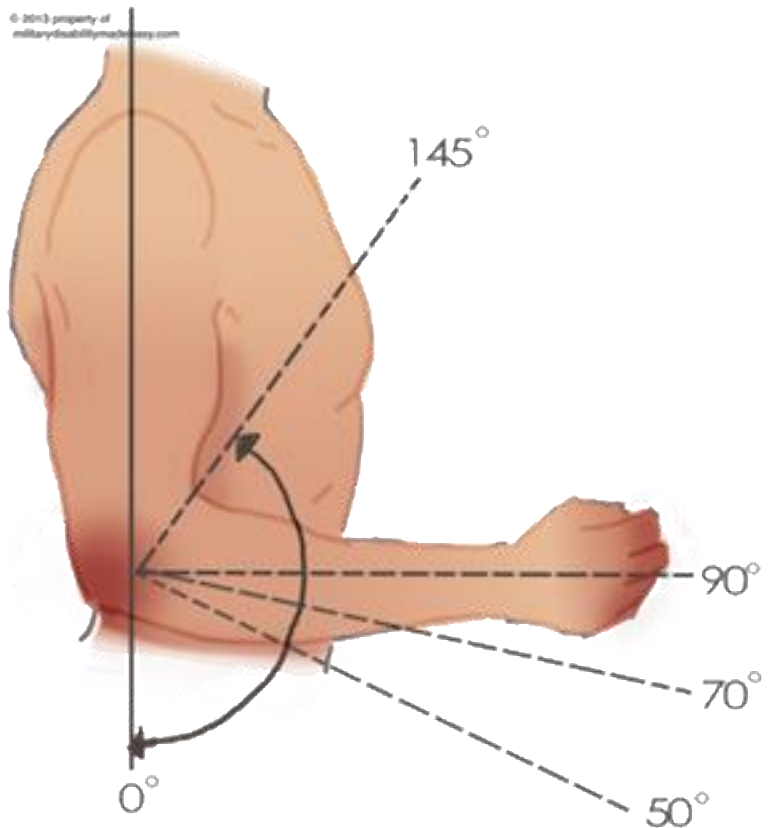
Xue Bin Peng, Angjoo Kanazawa, Jitendra Malik, Pieter Abbeel, Sergey
Levine

SIGGRAPH Asia 2018

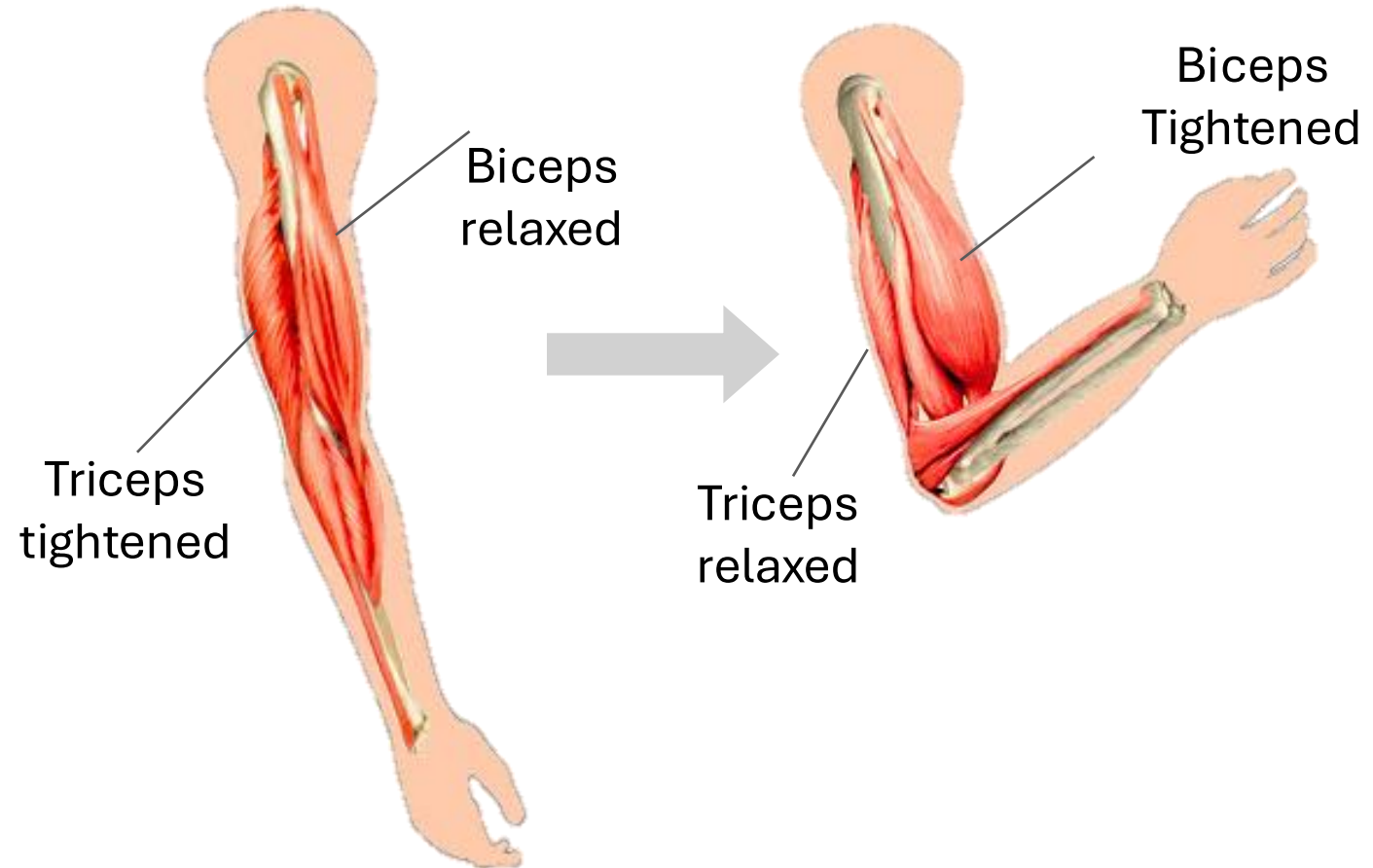


Perception is not the end of the story

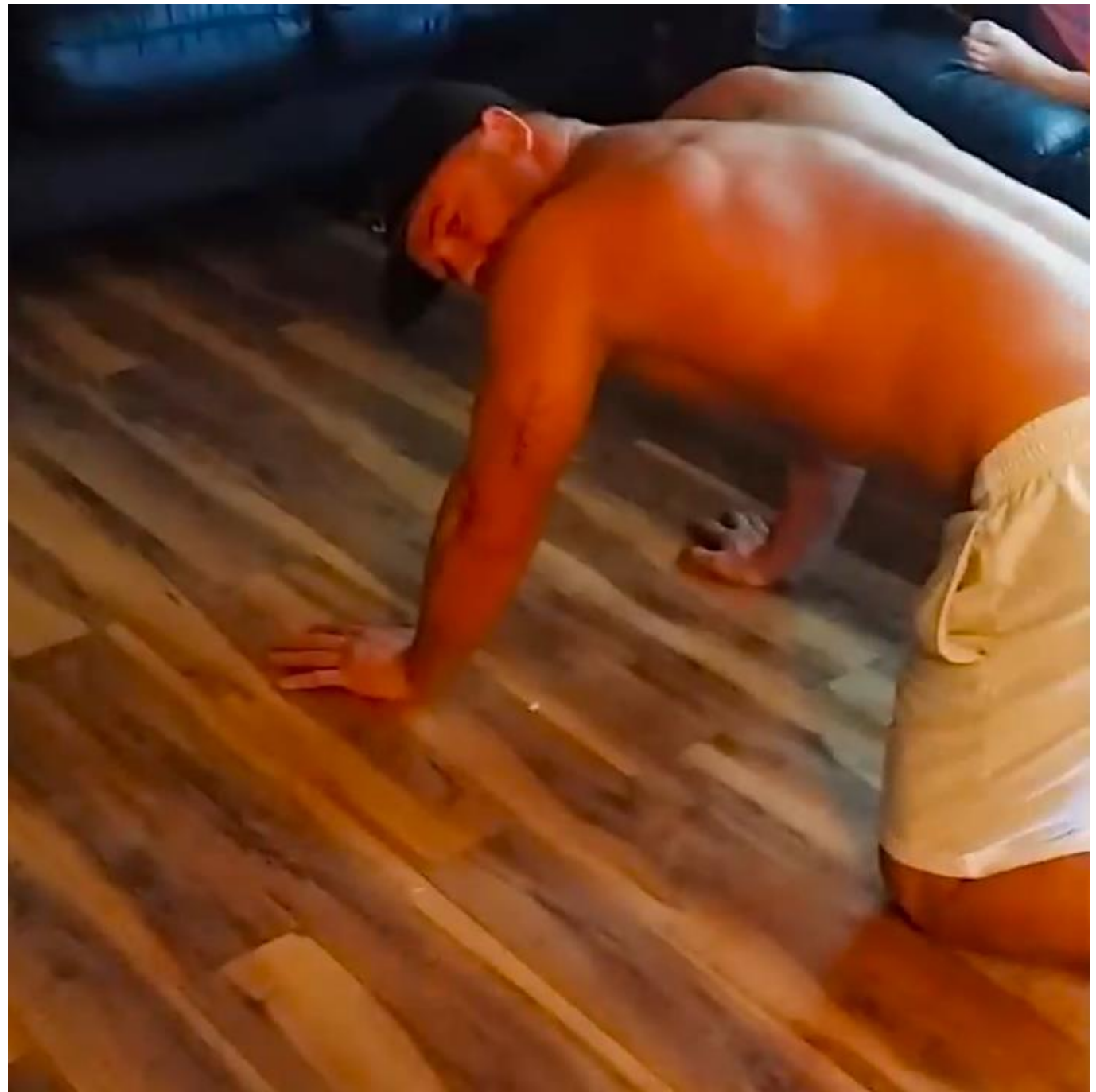
Perceiving the 3D pose



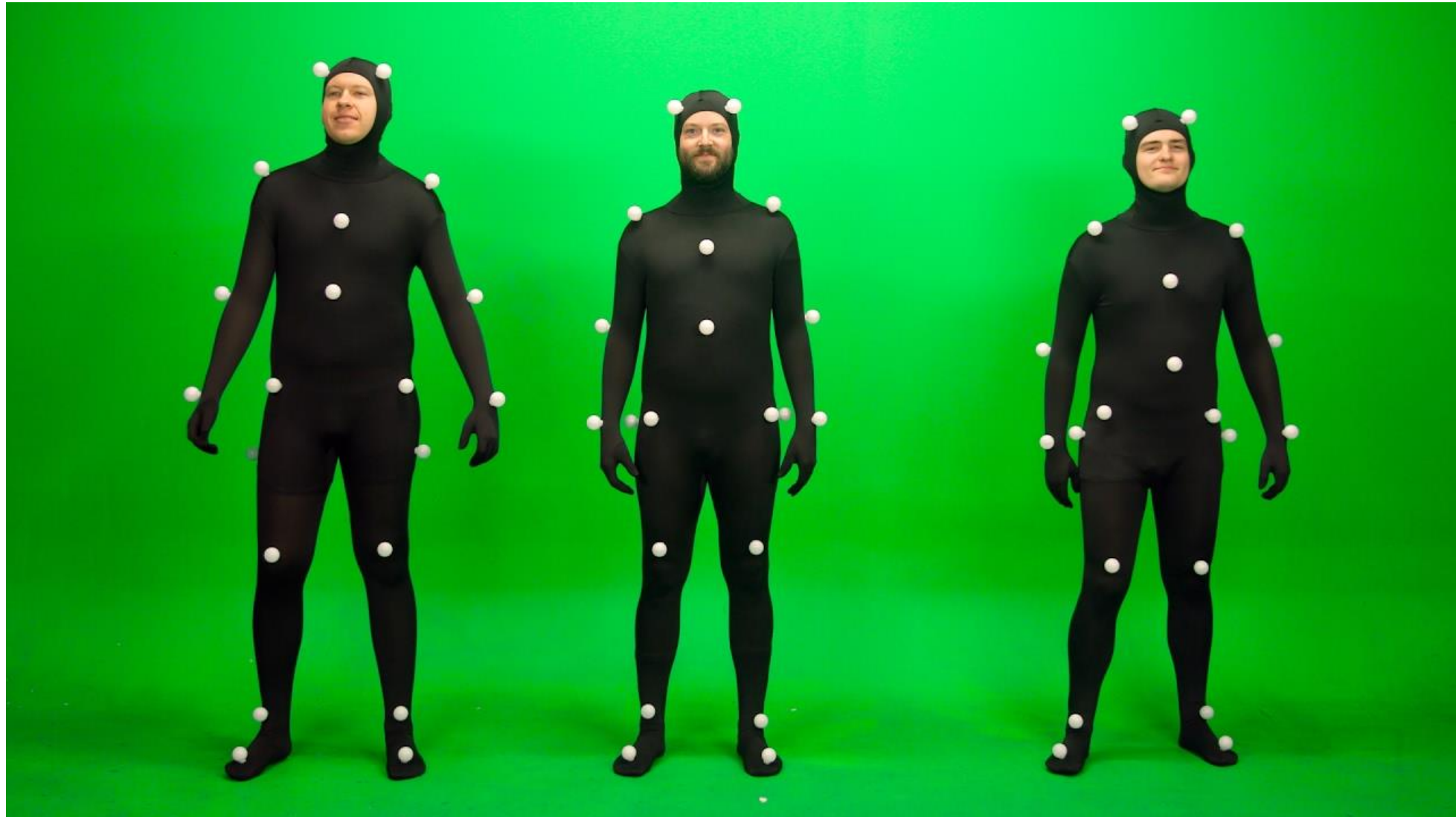
Actuating the muscle to get to the pose



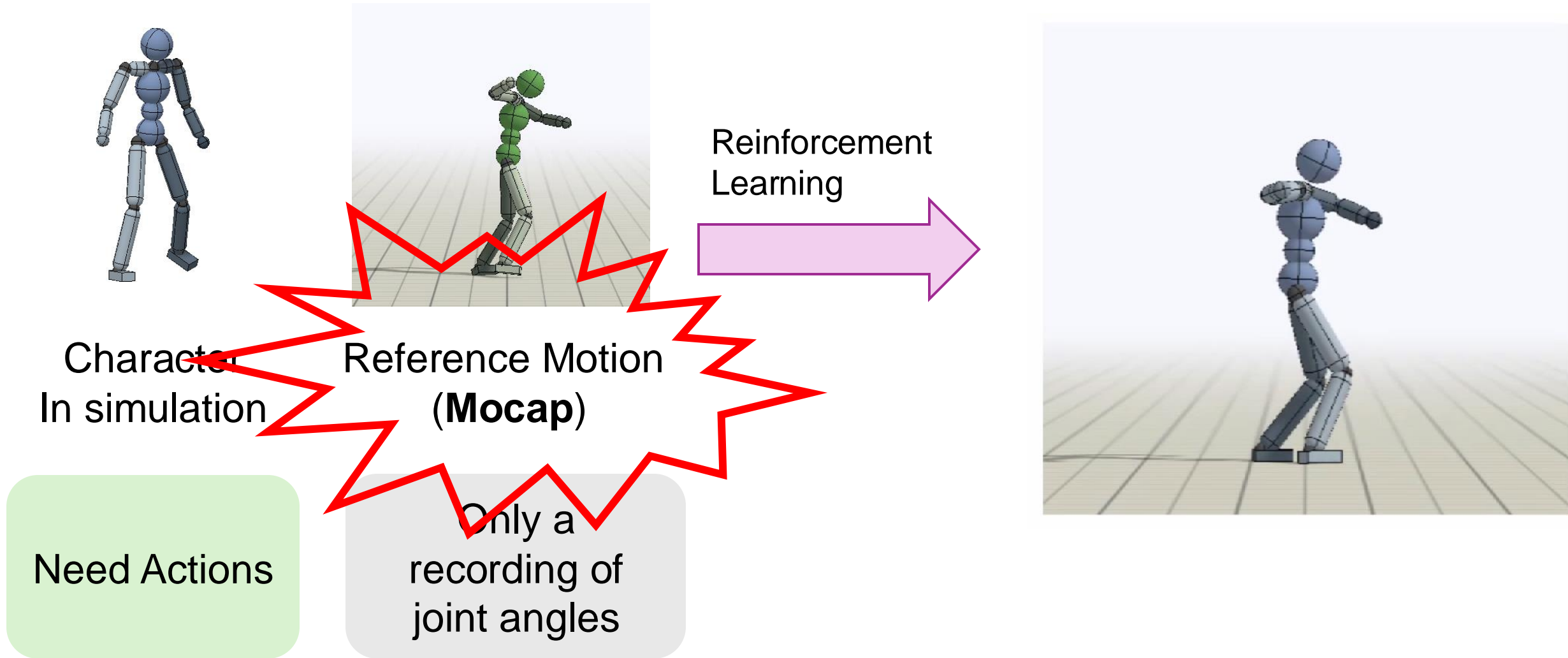
Control is not easy!



Past work on start from mocap



Deep Reinforcement Learning Based Motion Imitation



Expand the world to video



Learning Dynamic Skills from Videos



Video

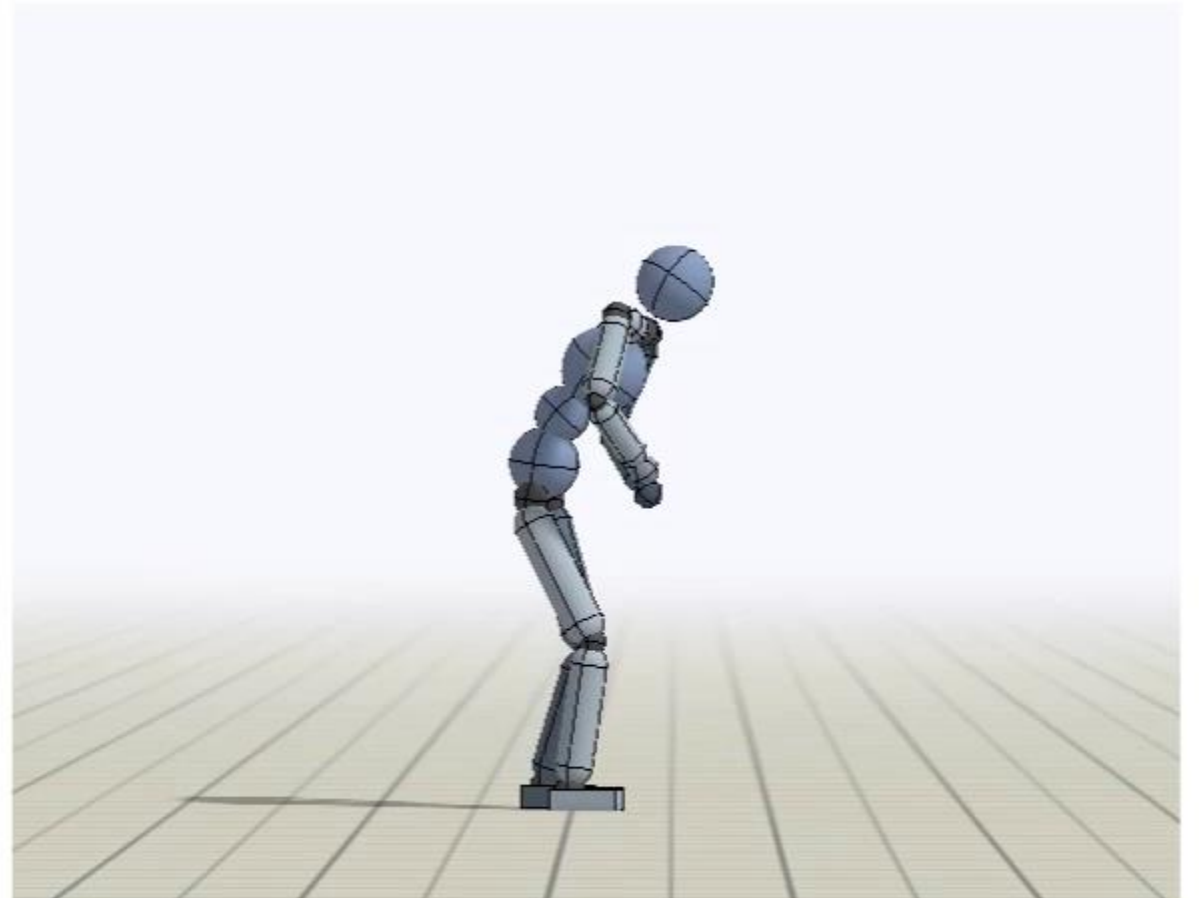


Temporally smooth mesh recovery

Use recovered 3D pose to train a physically simulated agent



Video

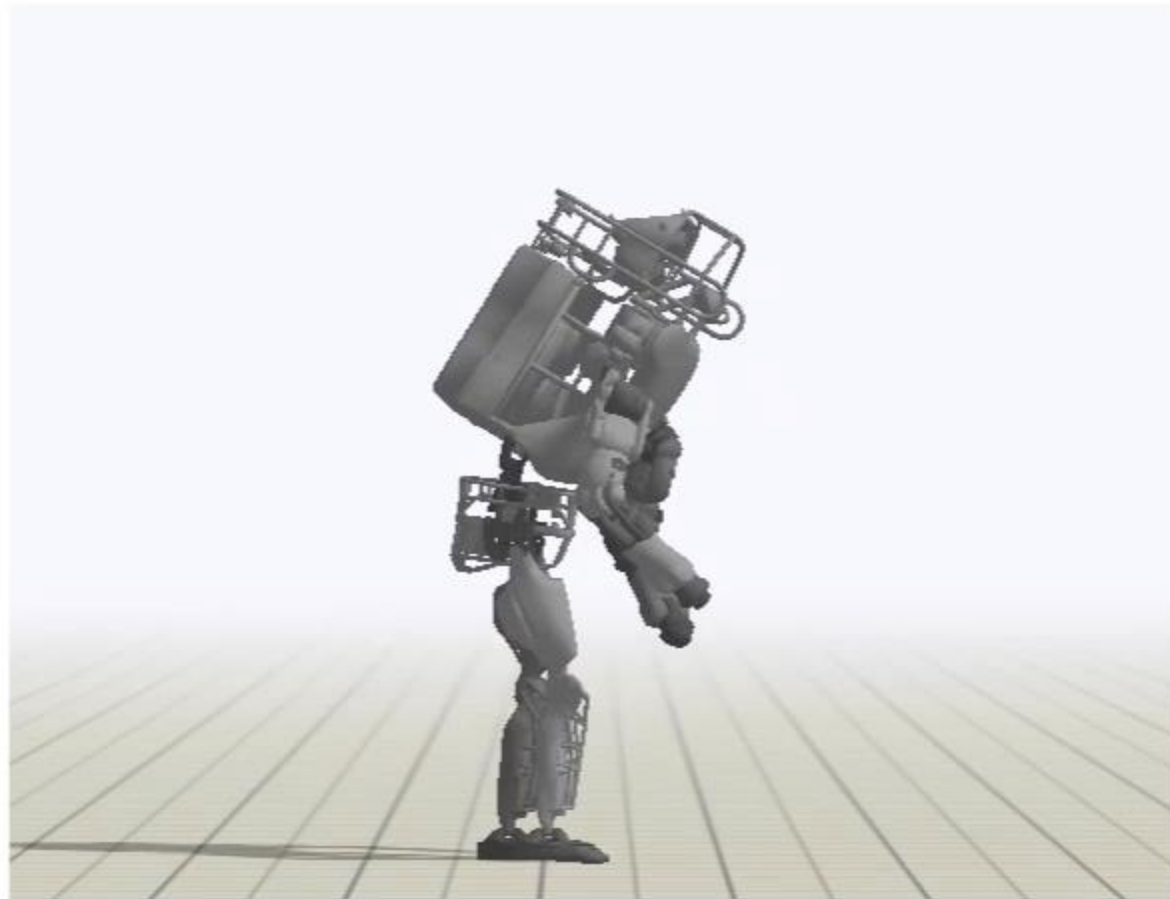


Policy

Train Atlas (169kg)

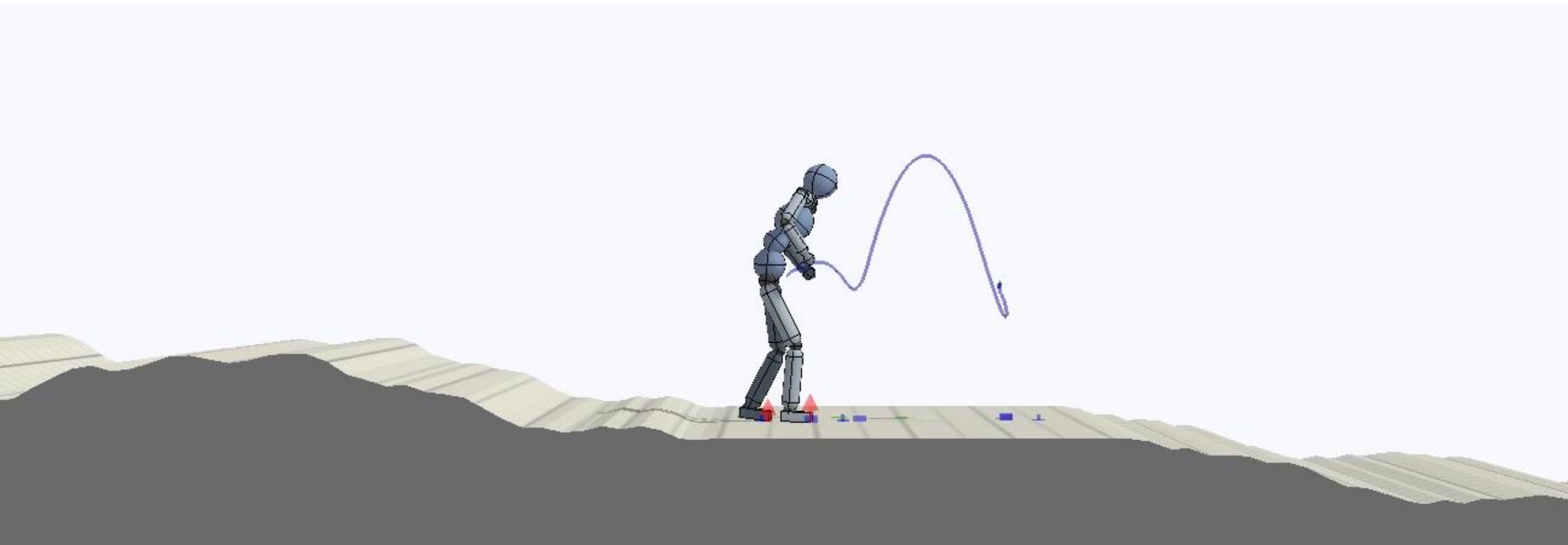


Video

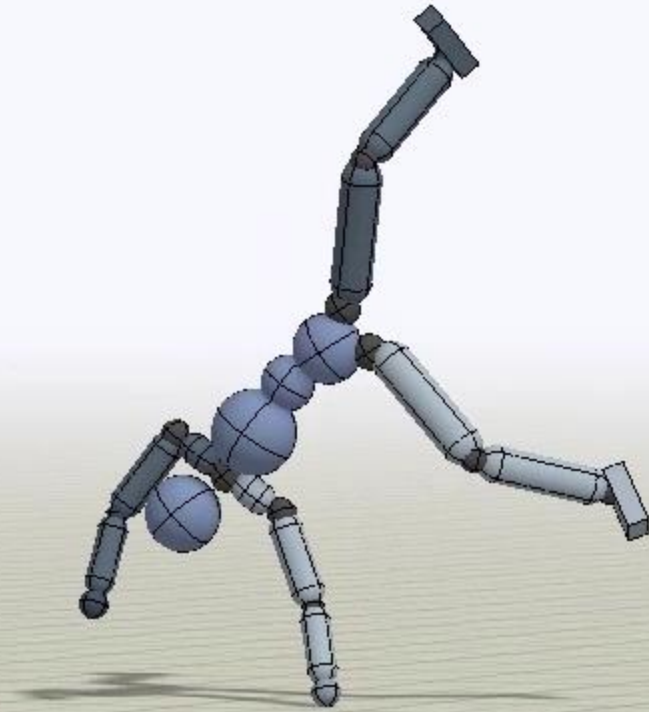


Policy

Environment Retargeting



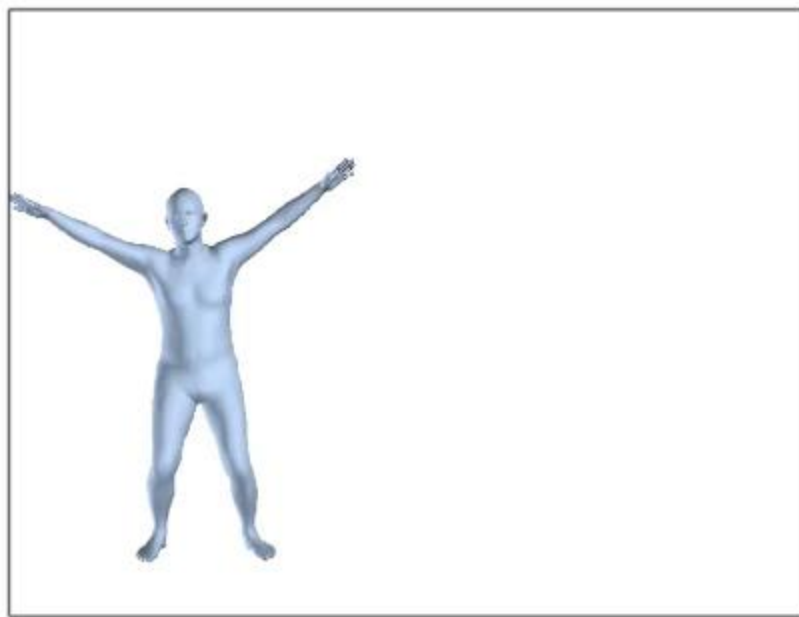
Robustness



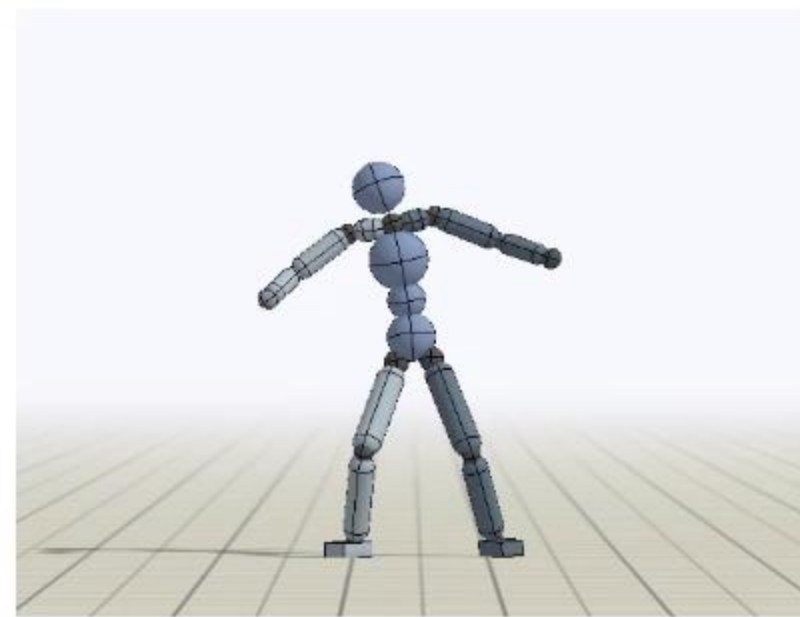
Humanoid: Cartwheel



Video



Recovered 3D Body

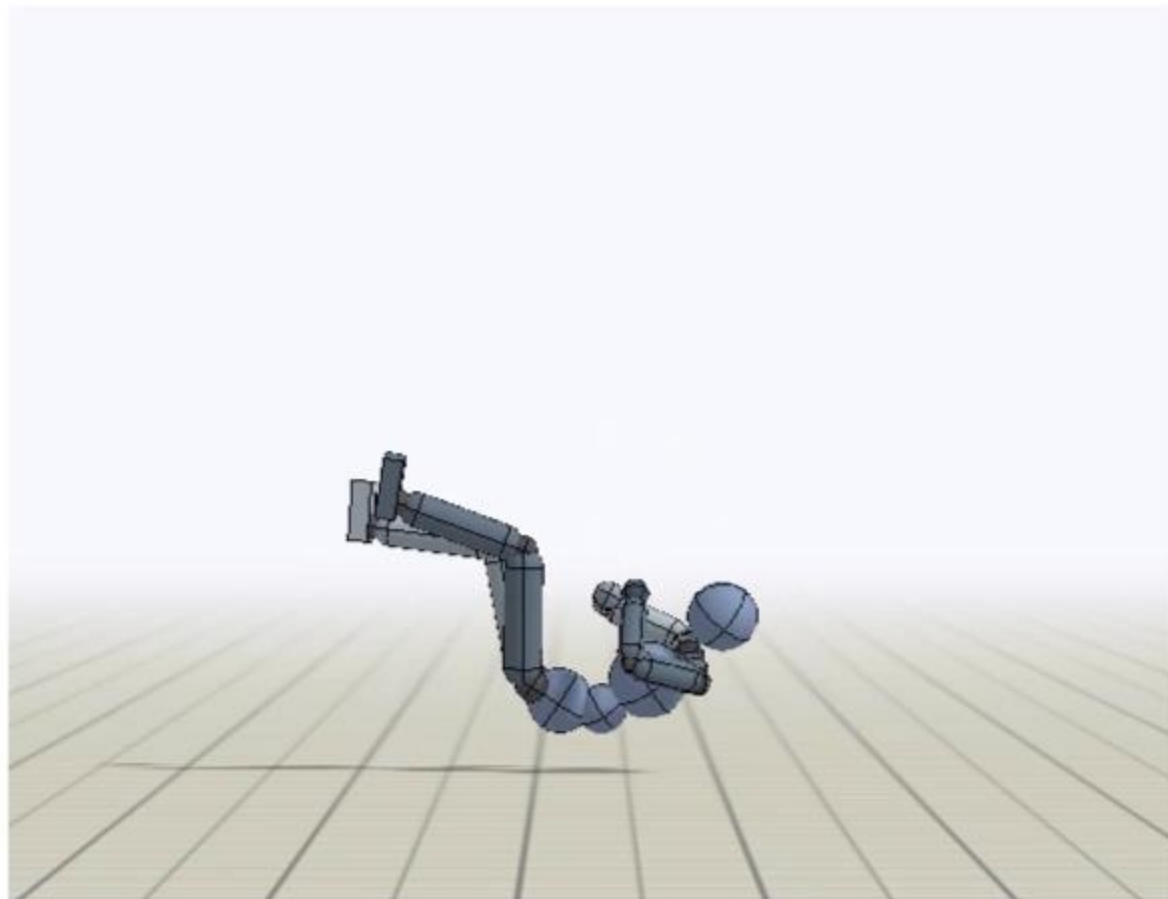


Policy

Humanoid: Kip-Up



Video: Kip-Up

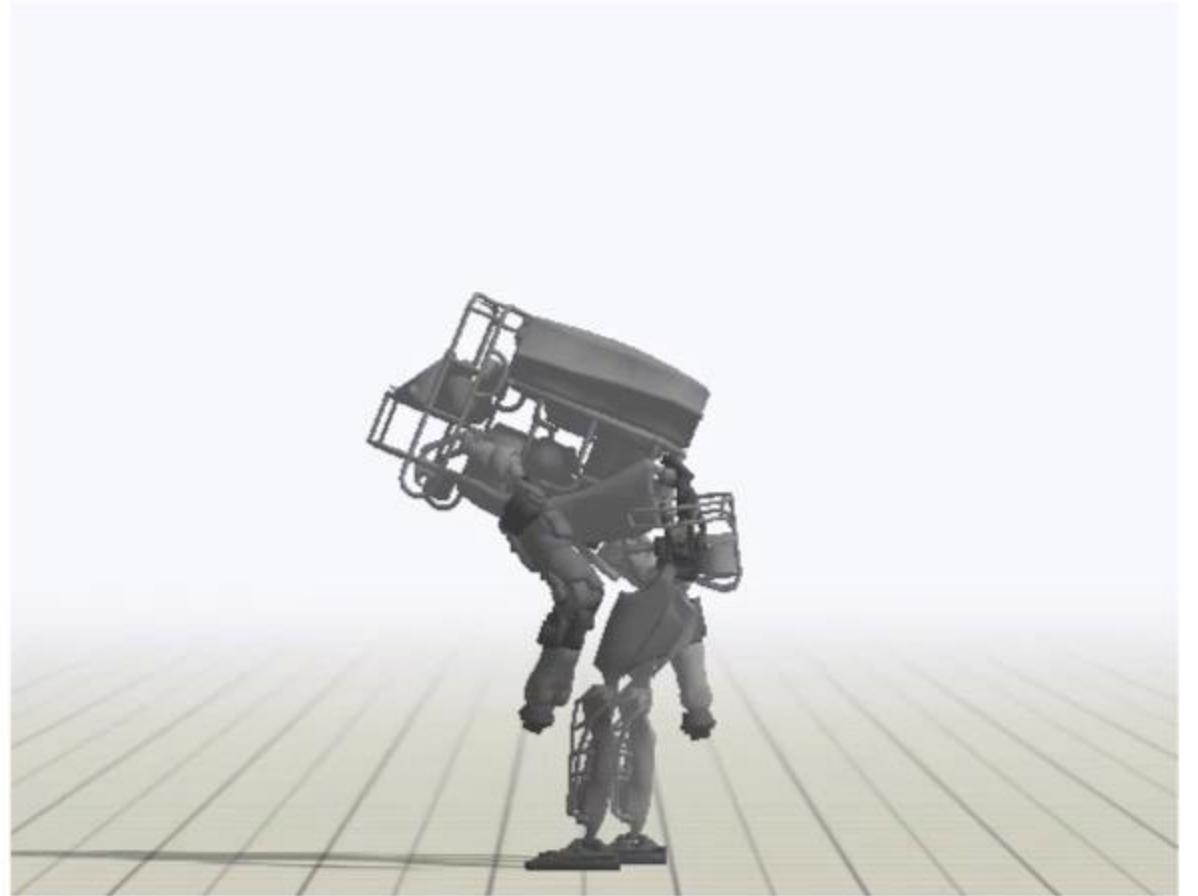


Policy

Atlas (169kg) : Handspring



Video: Handspring A

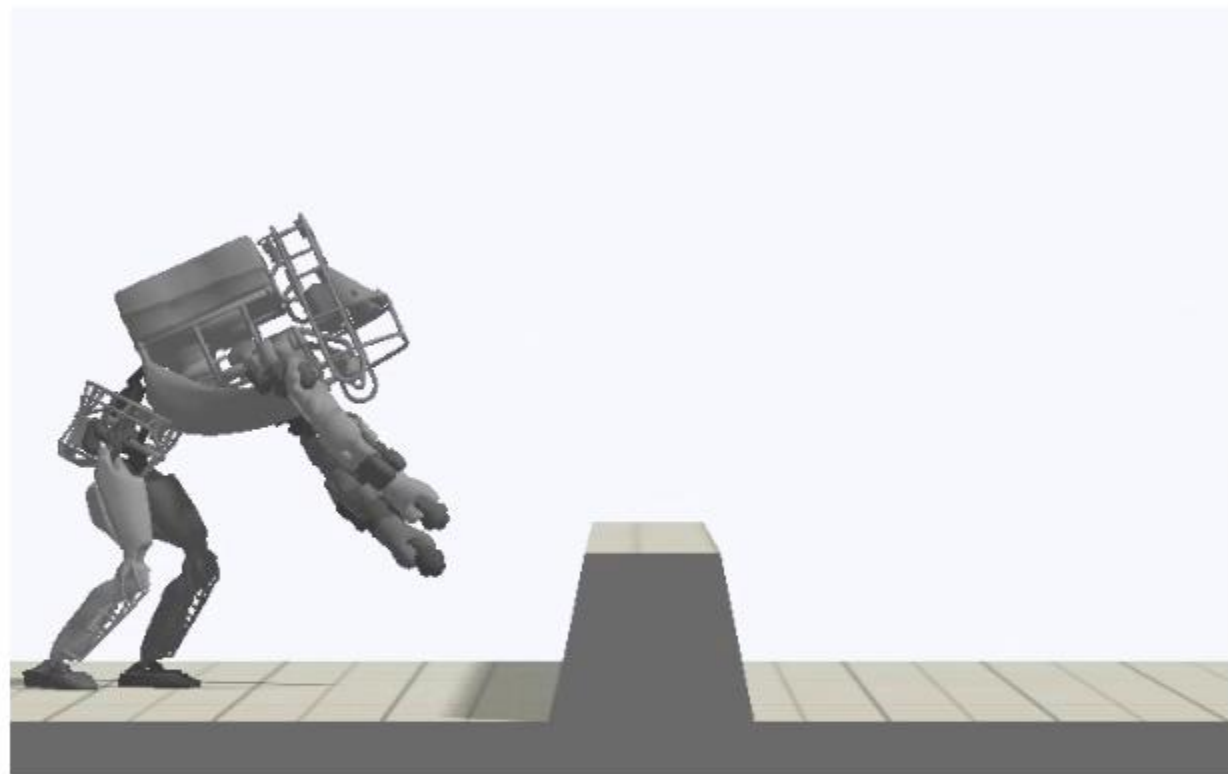


Policy

Atlas: Vault



Video: Vault

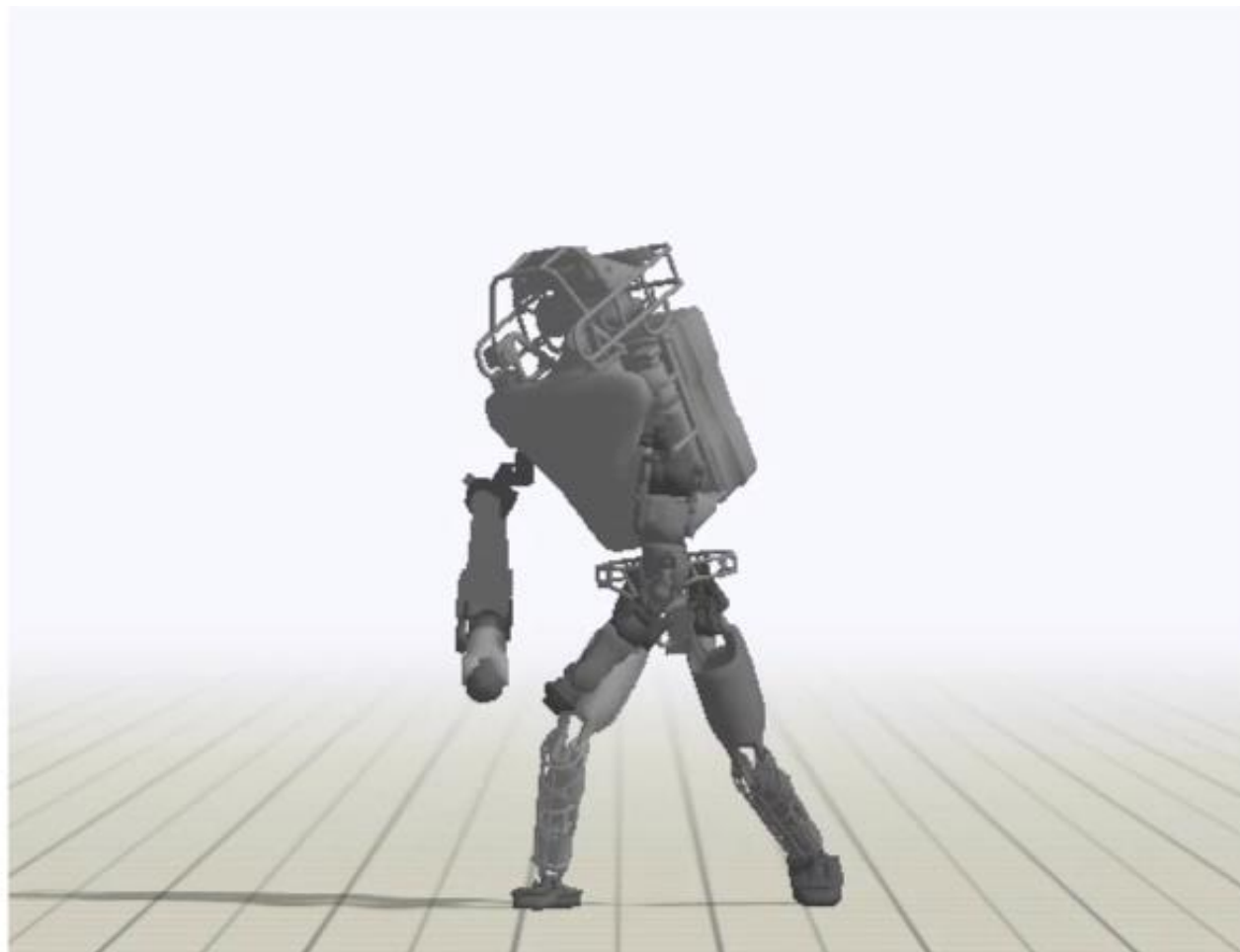


Policy

Atlas: Dance



Video



Policy

Failure Cases



Video



Recovered 3D Body



Policy

