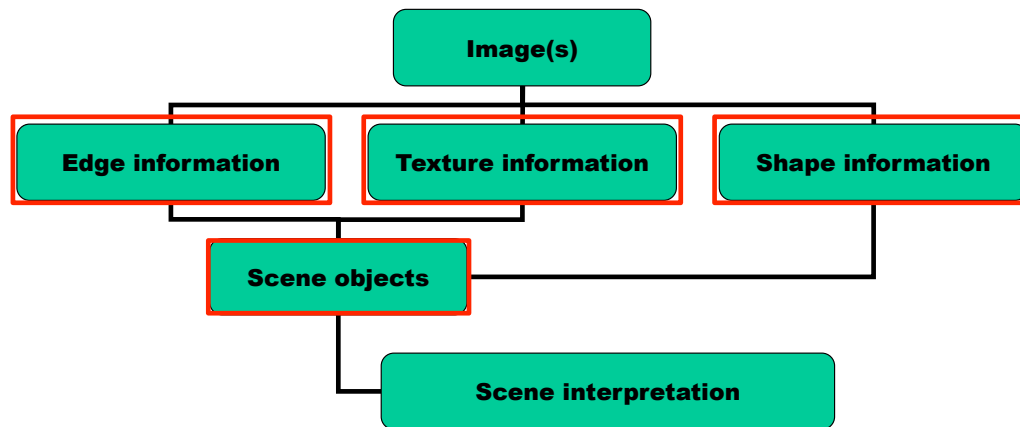# Modeling People and their Clothes
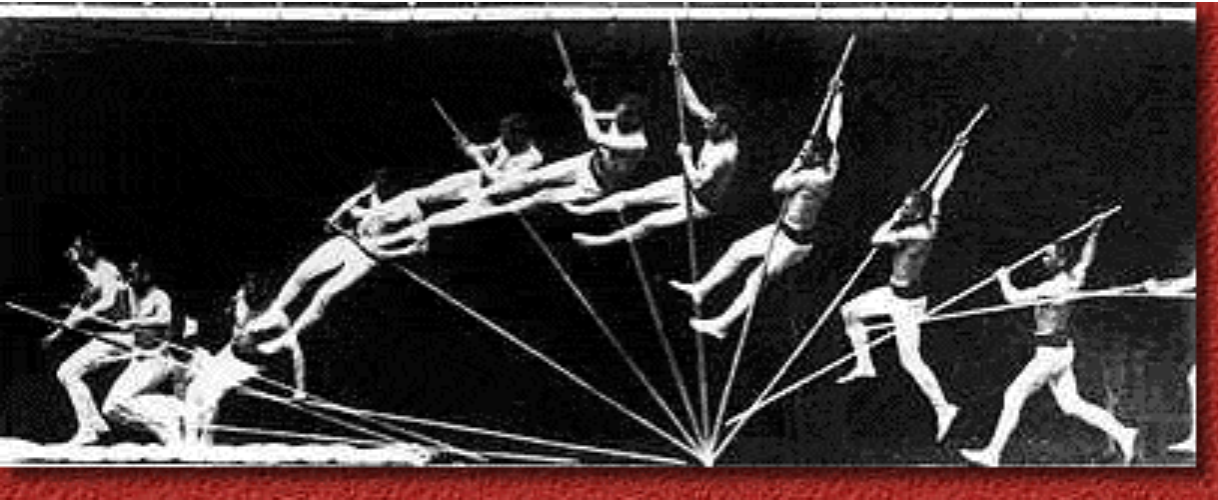
P. Fua
IC-CVLab
EPFL

# A Teachable Scheme



- Capturing the body by itself.
- Modeling the clothes in relation to it.
- Handling motion and deformations.
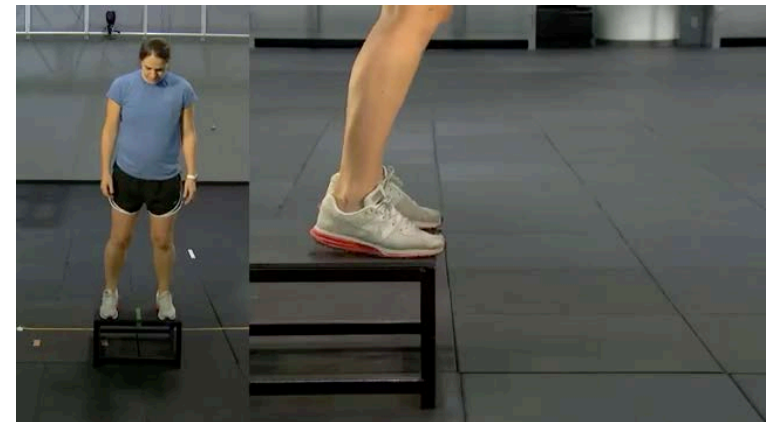
# Human Motion Analysis


Muybridge, circa 1890


EPFL, CVLab, 2025

# Applications



- Movies
- Fashion Design
- Sports Coaching
- Medicine
  - Enhancing performance
  - Injury prevention
  - Reeducation

# Articulated Body Model
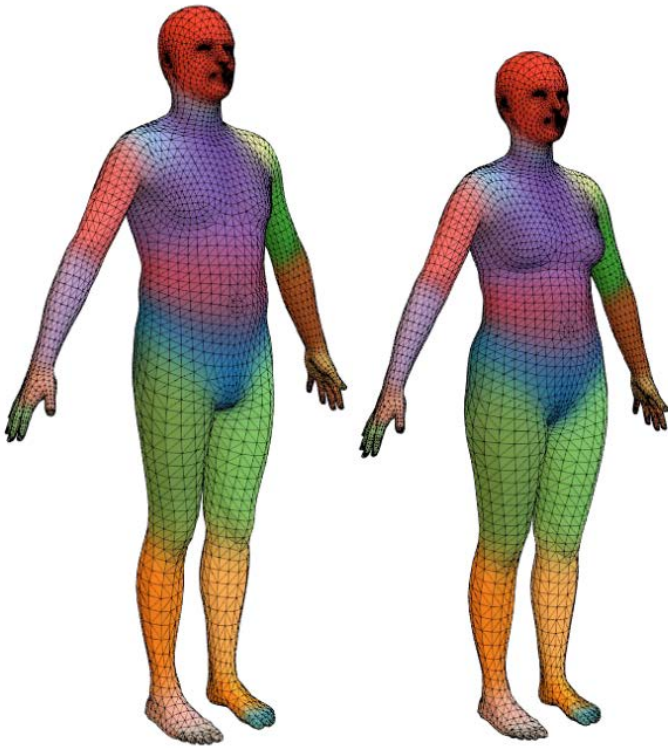
**Model**

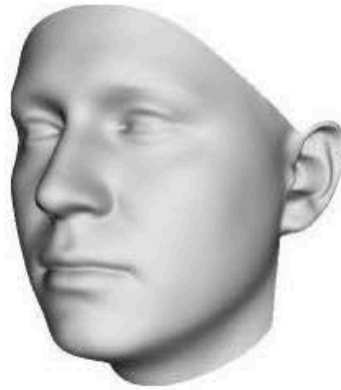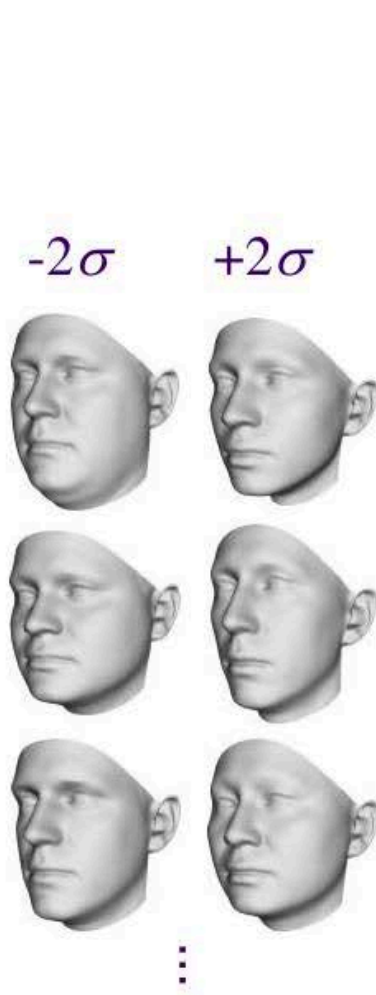**Model** with texture



| Pose | Shape | Dynamics | Texture |

- A model $M(\theta, \beta, \delta, A)$ takes as input a "small" number of pose, shape, and texture parameters and returns a 3D mesh.
- These parameters can be inferred from images and videos.

# Bodies as 3D Meshes (SMPL)



- The whole body can be represented as a low-resolution 3D mesh with 7000 vertices.

- That represents 71'000 parameters to infer from images.

- But these parameters are highly correlated.

➡ The model must encode these correlations.

# Reminder: PCA Face Model



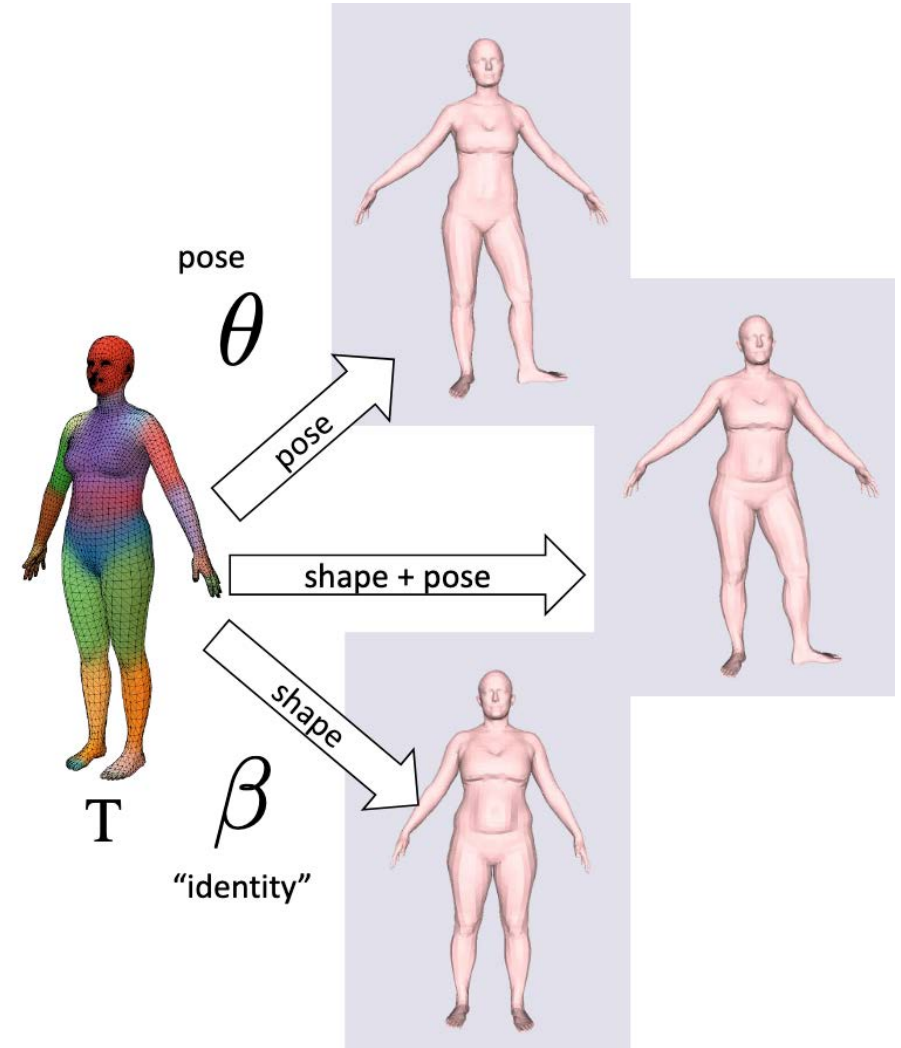$$S = \bar{S} + \sum_{i=1}^{99} \alpha_i S_j$$

$\bar{S}$ :   Average shape

$S_i$ :   Shape vector

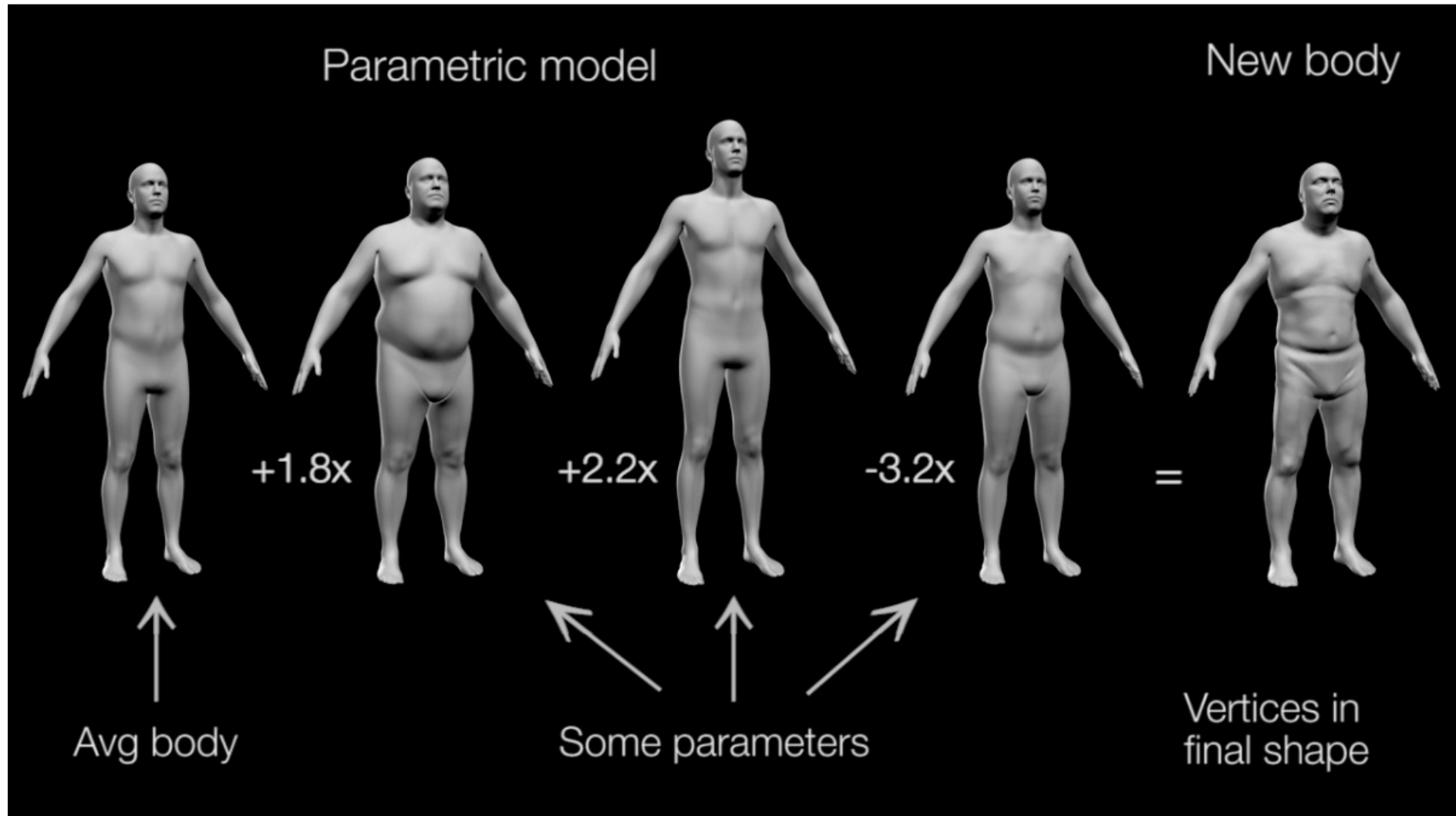$\alpha_i$ :   Shape coefficients

# Factored Model (SCAPE)

- The model parameterizes deviations from a template mesh.

- Uses the same kind of dimensionality reduction techniques as those used to create face morphable models.

- Requires a large training database for learning purposes.

➡ Simplifies learning and inference.



pose

$\theta$

pose

shape + pose

shape

$T$ $\beta$

"identity"

# Changing the PCA Coefficients (SMPL)



Parametric model

New body

+1.8x    +2.2x    -3.2x    =

Avg body

Some parameters

Vertices in final shape

# 4D Body Shapes

# 3D at 60 FPs



University of Tübingen / MPI-Informatics:
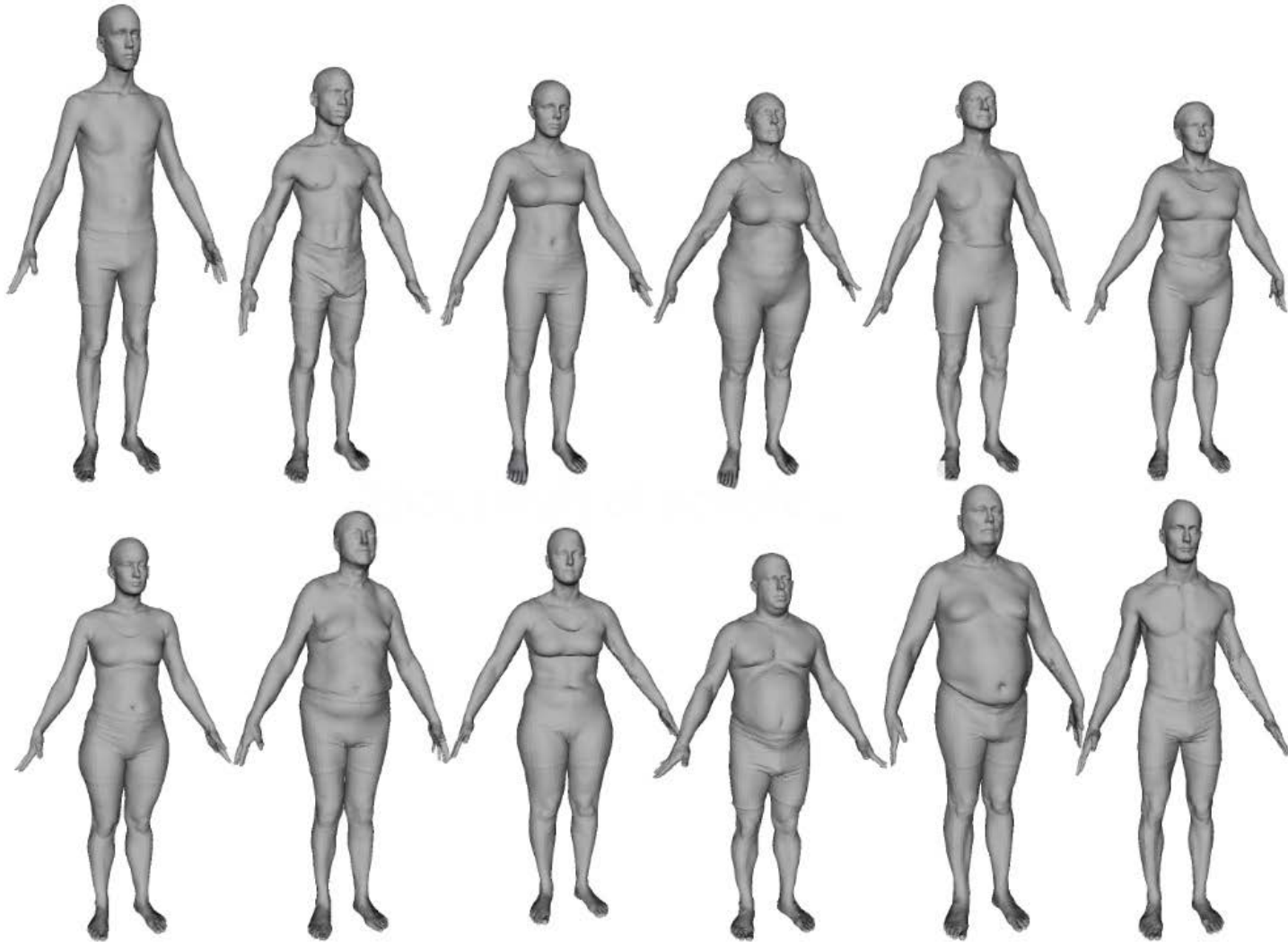- Thousands of people
- Thousands of poses

# Many Different Body Shapes

# Many Different Poses

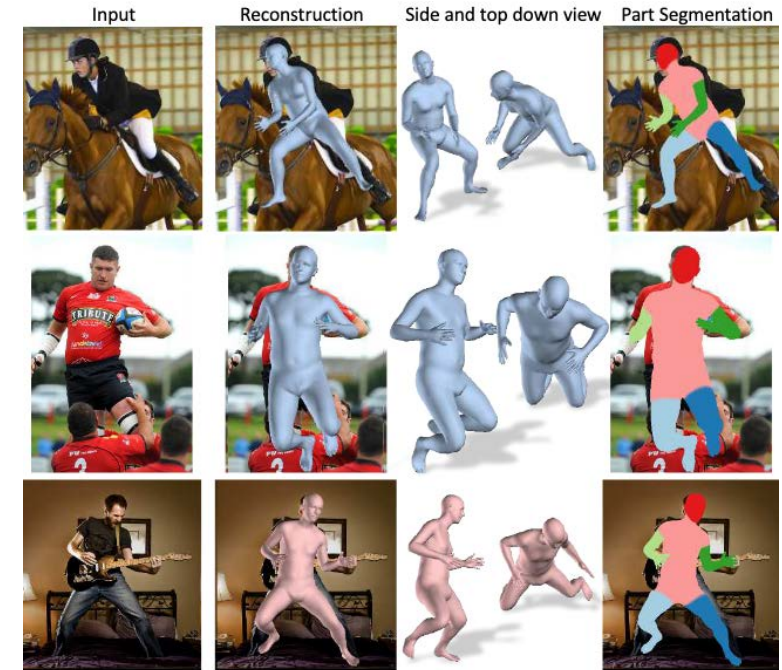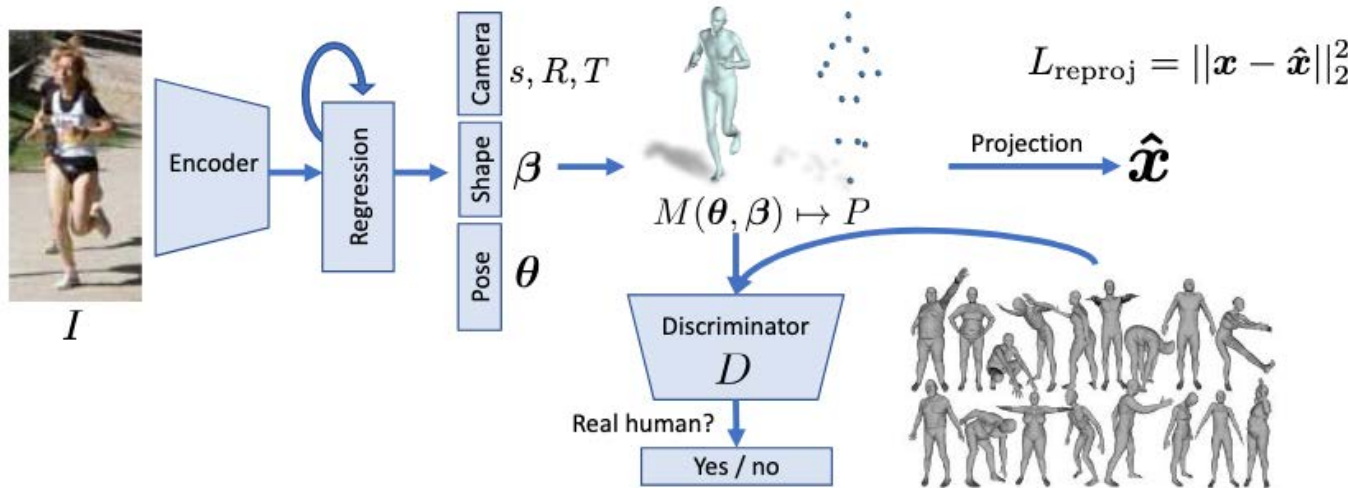# From Image to Body Pose



- Use a CNN to detect 2D joints.
- Infer SMPL parameters from those.

Not all joints can be expected to be visible!

# Increasing Robustness



$$L_{\text{reproj}} = ||\boldsymbol{x} - \hat{\boldsymbol{x}}||_2^2$$

$$M(\boldsymbol{\theta}, \boldsymbol{\beta}) \mapsto P$$

- Detect 2D joints.
- Infer SMPL parameters from those.
- Use adversarial training to ensure consistency.

# From Video to Body Motion



- Estimate SMPL parameters from each individual video frame while enforcing temporal consistency.
- Use an adversarial network to enforce realism, given a large motion training set.

# Hidden Joints



HMR 2.0

Input Image → Vision Transformer → Transformer w/ Cross Attn → MLP

SMPL Query Token → Transformer w/ Cross Attn

$\theta$ Pose
$\beta$ Shape
$\pi$ Camera

- Loose clothing can hide individual joints.
- Bring in the transformers!

➡ Direct regression from image to SMPL parameters.

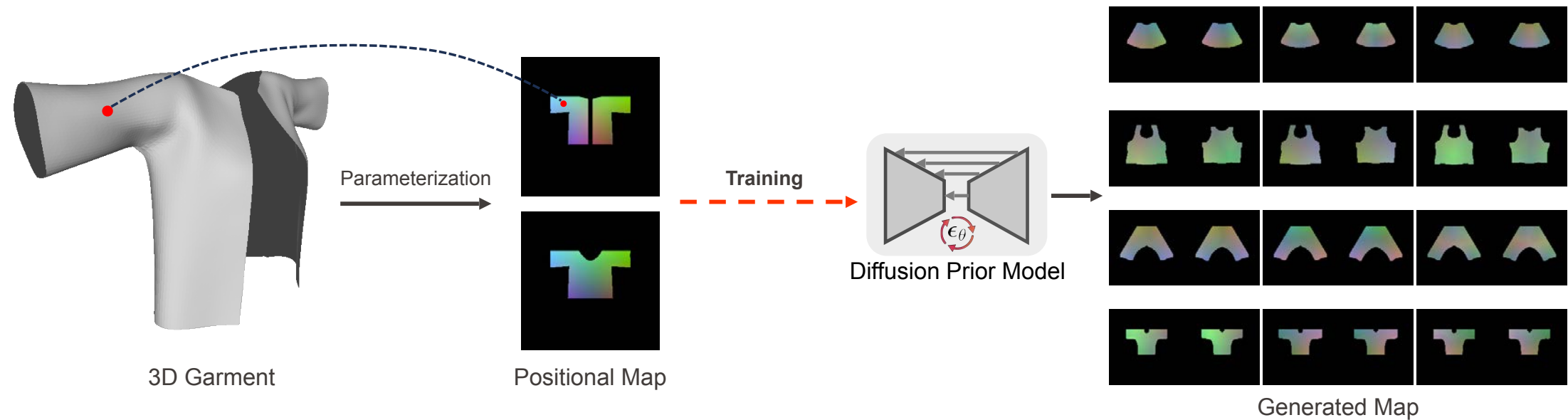# What About People Wearing Loose Clothing?
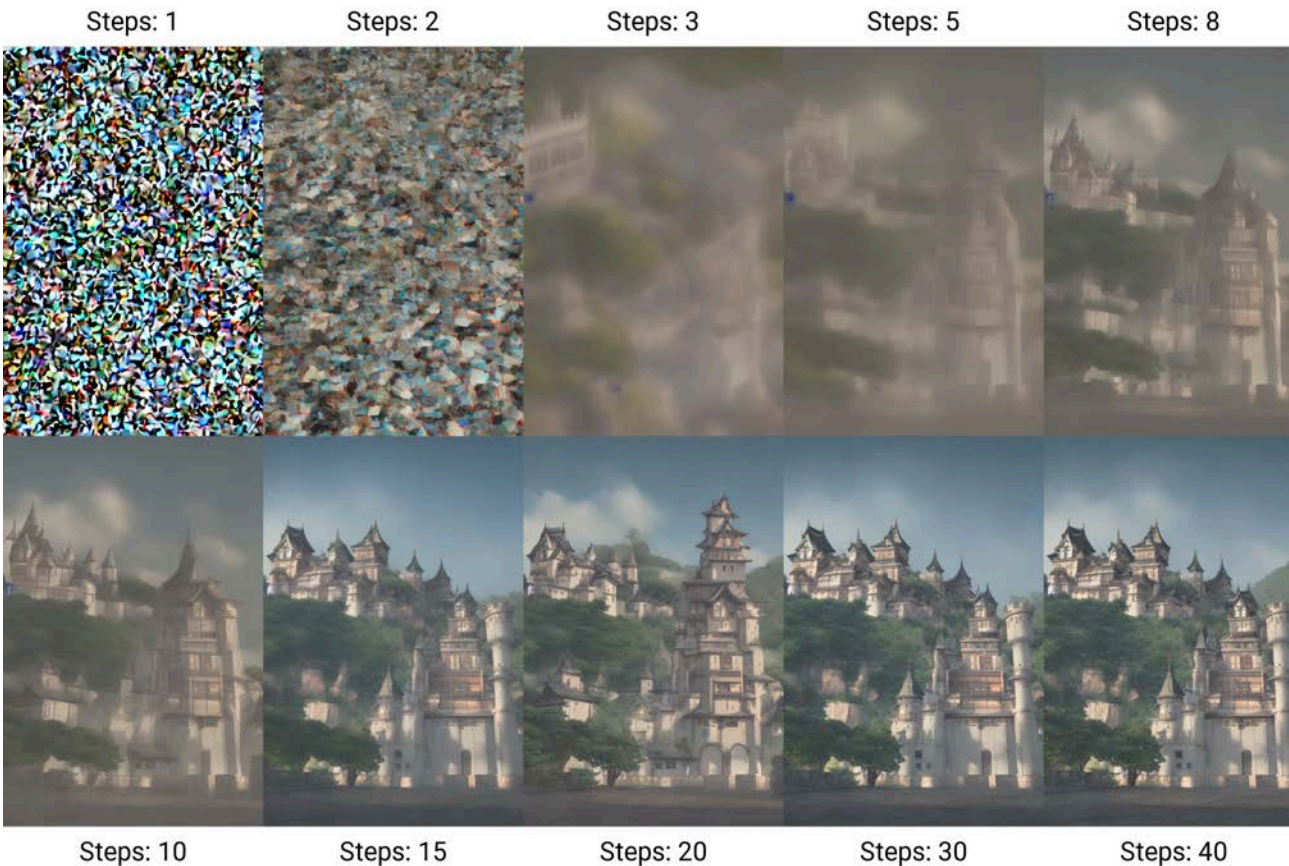


Can we also recover the shape and motion of the clothes in addition to that of the body?

# Garment Parameterization



3D Garment → Parameterization → Positional Map → Training → Diffusion Prior Model ($\epsilon_\theta$) → Generated Map

- Parameterize a 3D garment as a set of 2D positional maps.
- Train a diffusion model on it to learn the shape prior.

# Diffusion / Flow Matching

Steps: 1     Steps: 2     Steps: 3     Steps: 5     Steps: 8

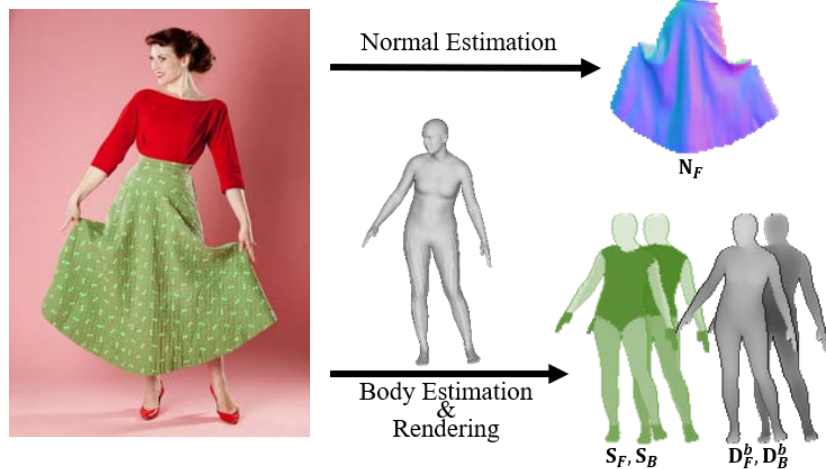Steps: 10     Steps: 15     Steps: 20     Steps: 30     Steps: 40

- Train a network to turn noise into a distribution that conforms to a specific prior.
- Can be "guided" to obey some constraints.

—> We use it to generate realistic clothing that matches the images.

# Reconstruction Pipeline

Given an image of clothed person, its garment normal estimation $\mathbf{N}_F$ and body part/depth estimation $(\mathbf{S}_F, \mathbf{S}_B, \mathbf{D}_F^b, \mathbf{D}_B^b)$, we
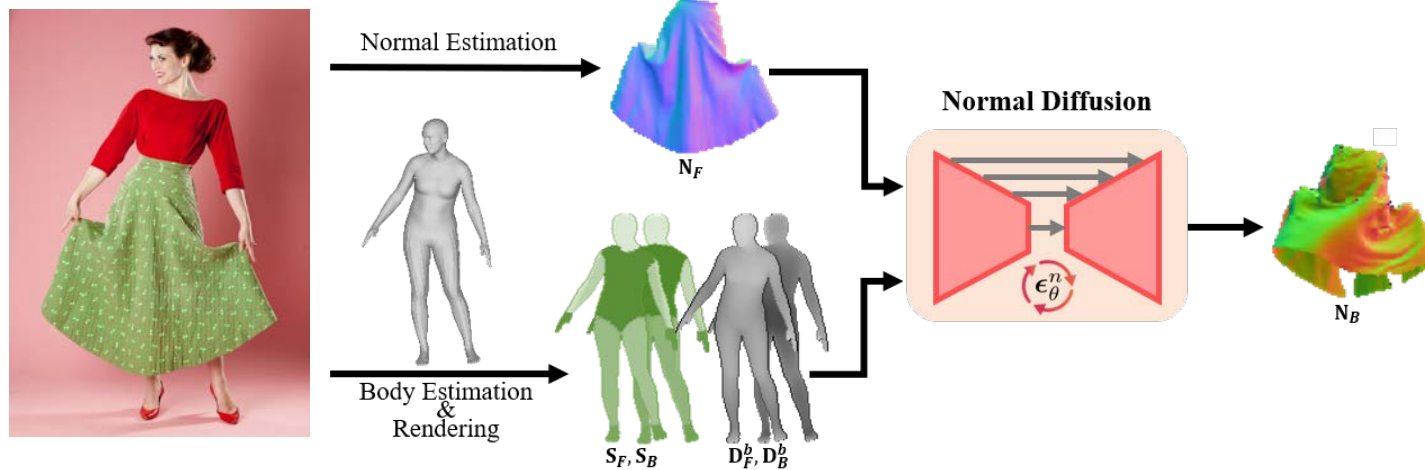
- synthesize the normal for the invisible back view $\mathbf{N}_B$
- predict the UV coordinates $(\mathbf{C}_F, \mathbf{C}_B)$ and depth $(\mathbf{D}_F^g, \mathbf{D}_B^g)$
- turn predictions to UV positional maps
- fit the prior to the positional maps for reconstruction

# Reconstruction Pipeline

Given an image of clothed person, its garment normal estimation $N_F$ and body part/depth estimation $(S_F, S_B, D_F^b, D_B^b)$, we

- synthesize the normal for the invisible back view $N_B$
- predict the UV coordinates $(C_F, C_B)$ and depth $(D_F^g, D_B^g)$
- turn predictions to UV positional maps
- fit the prior to the positional maps for reconstruction

# Reconstruction Pipeline

Given an image of clothed person, its garment normal estimation $\mathbf{N}_F$ and body part/depth estimation $(\mathbf{S}_F, \mathbf{S}_B, \mathbf{D}_F^b, \mathbf{D}_B^b)$, we
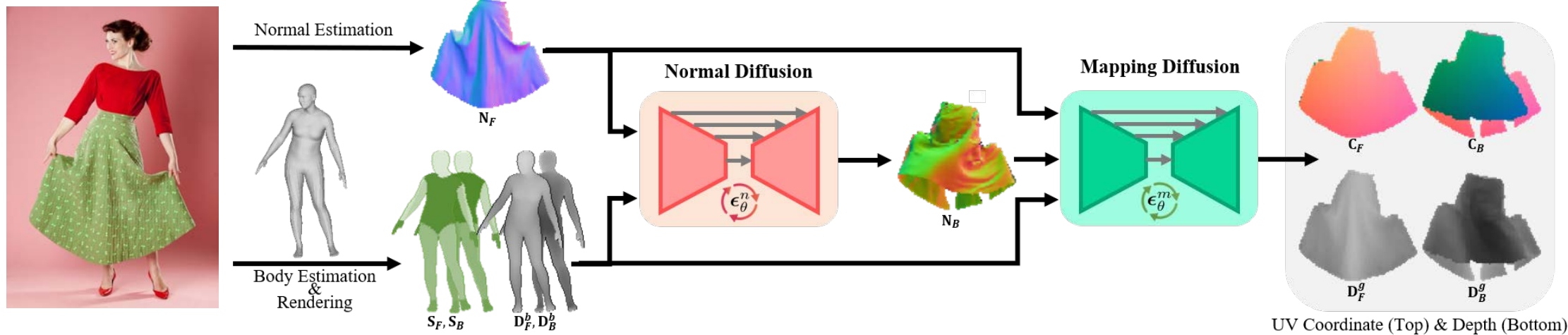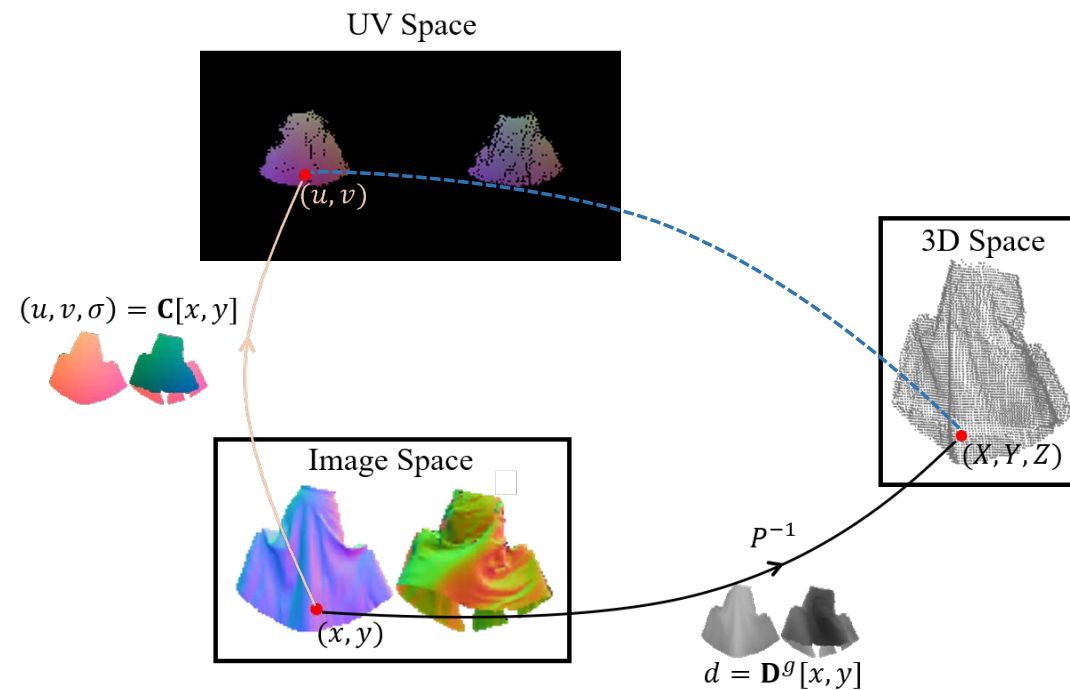
- synthesize the normal for the invisible back view $\mathbf{N}_B$
- predict the UV coordinates $(\mathbf{C}_F, \mathbf{C}_B)$ and depth $(\mathbf{D}_F^g, \mathbf{D}_B^g)$
- turn predictions to UV positional maps
- fit the prior to the positional maps for reconstruction



UV Coordinate (Top) & Depth (Bottom)

# Reconstruction Pipeline

Given an image of clothed person, its garment normal estimation $N_F$ and body part/depth estimation $(S_F, S_B, D_F^b, D_B^b)$, we

- synthesize the normal for the invisible back view $N_B$
- predict the UV coordinates $(C_F, C_B)$ and depth $(D_F^g, D_B^g)$
- turn predictions to UV positional maps
- fit the prior to the positional maps for reconstruction

UV Space

$(u, v)$

$(u, v, \sigma) = \mathbf{C}[x, y]$

3D Space

$(X, Y, Z)$

Image Space

$(x, y)$

$P^{-1}$

$d = \mathbf{D}^g[x, y]$

EPFL

iVLab

# Reconstruction Pipeline

Given an image of clothed person, its garment normal estimation $N_F$ and body part/depth estimation $(S_F, S_B, D_F^b, D_B^b)$, we
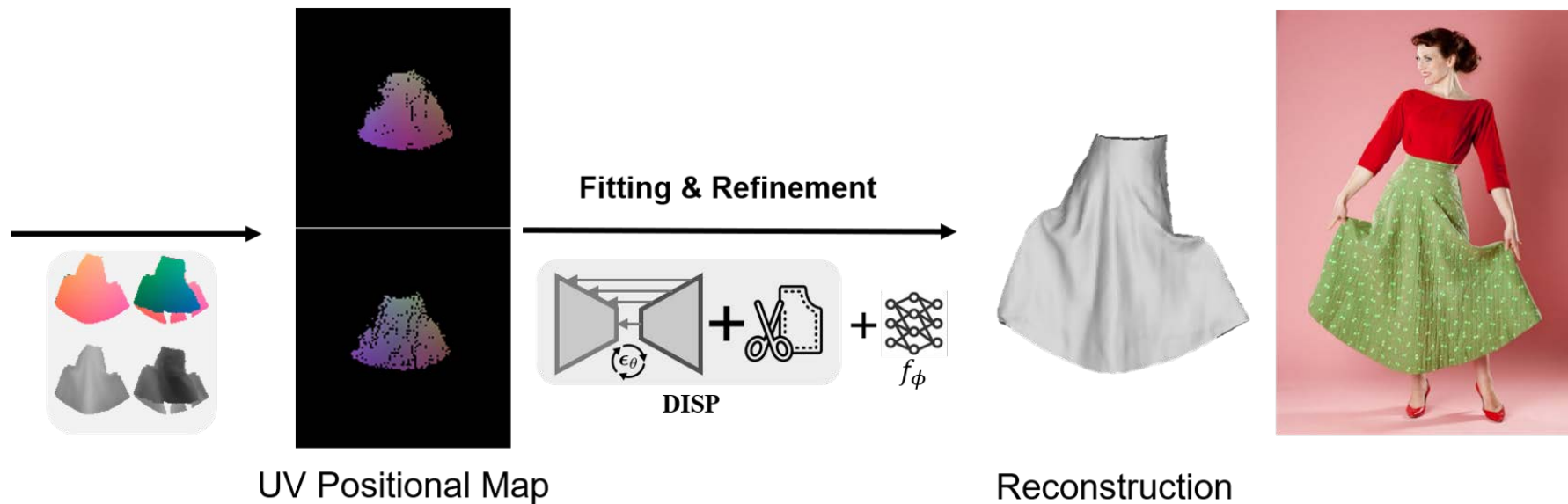
- synthesize the normal for the invisible back view $N_B$
- predict the UV coordinates $(C_F, C_B)$ and depth $(D_F^g, D_B^g)$
- turn predictions to UV positional maps
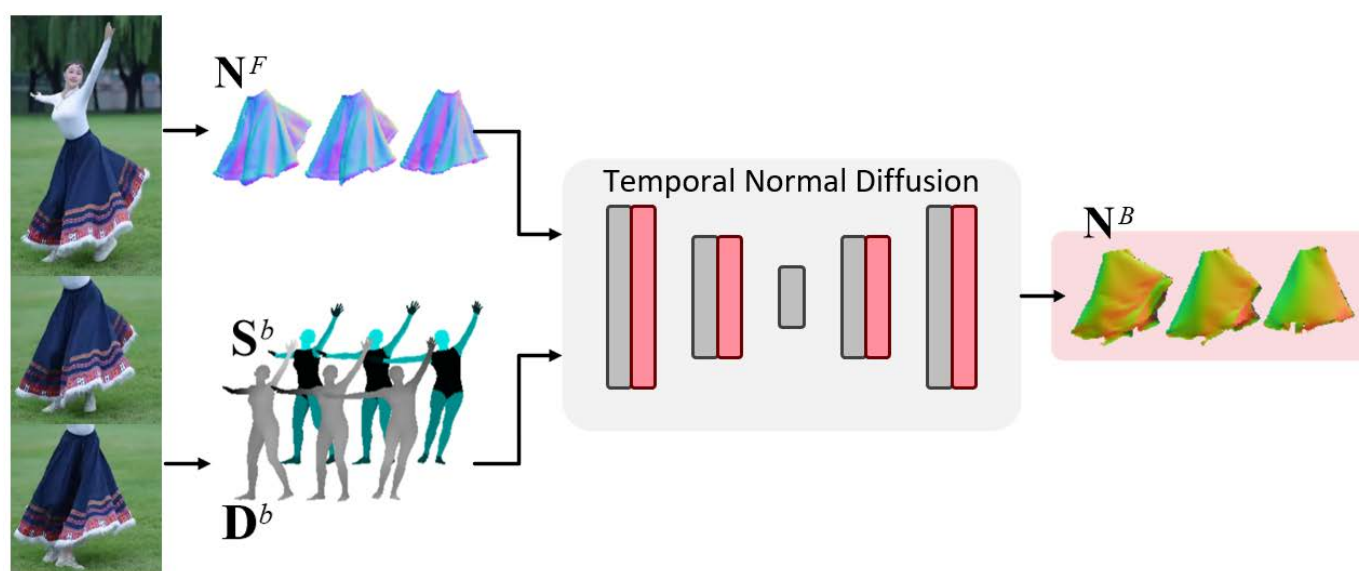- fit the prior to positional maps for reconstruction



**Fitting & Refinement**

**DISP**

$\epsilon_\theta$

$f_\phi$

UV Positional Map

Reconstruction

# Garment Recovery from Real Images

Li et al. SIGGRAPH'25
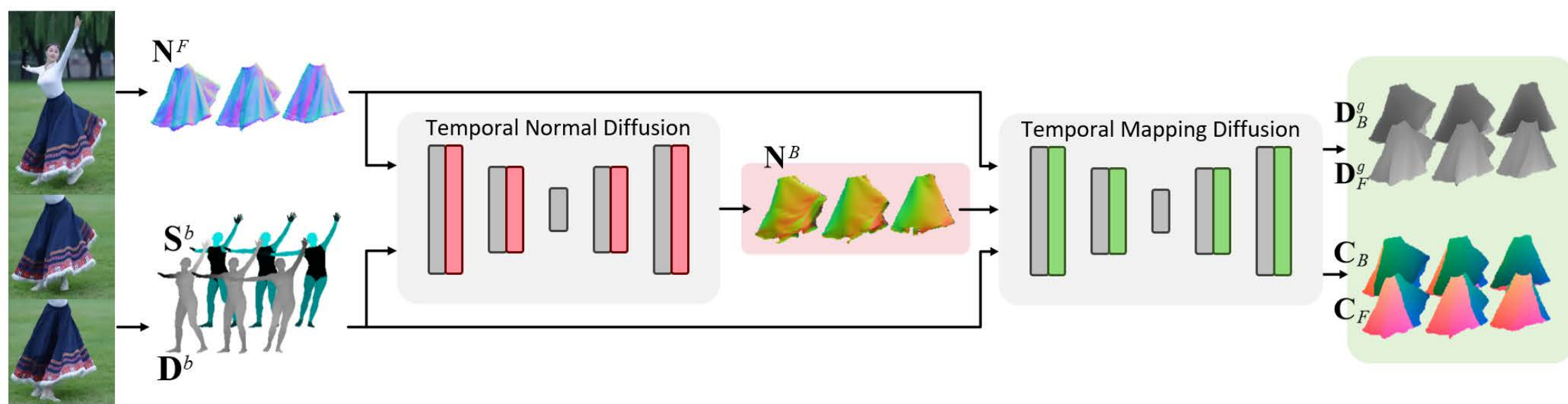
EPFL

# From Images to Video

To handle video cases, we
- introduce temporal diffusion models
- enforce geometric and temporal guidance
- temporal consistency guidance, depth-to-normal guidance, interpenetration-aware guidance
- fit the prior to the positional maps with projection-based constraint for reconstruction

# From Images to Videos
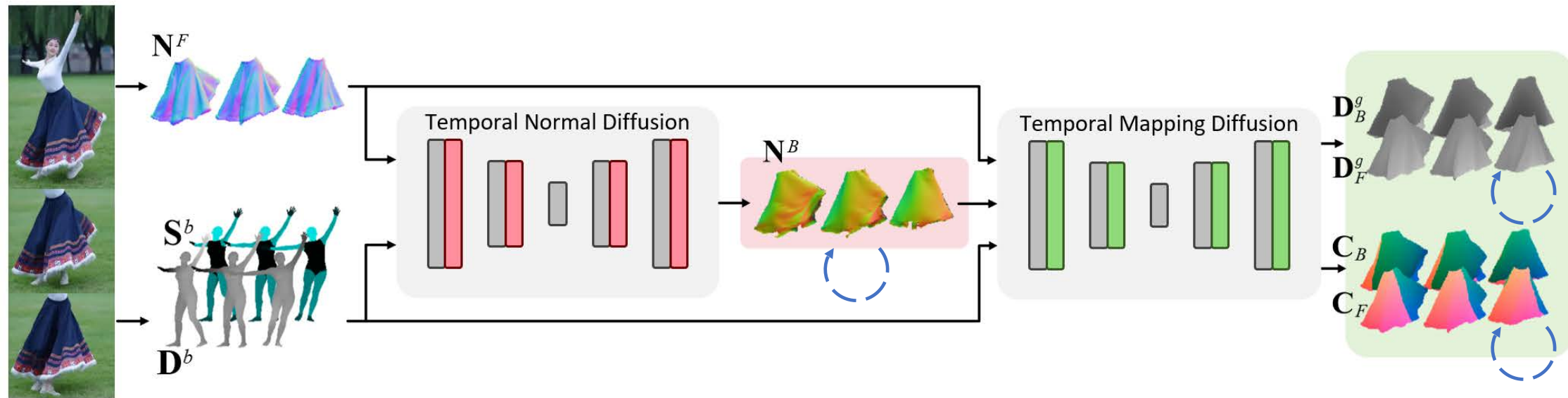
To handle video cases, we
- introduce temporal diffusion models
- enforce geometric and temporal guidance
- temporal consistency guidance, depth-to-normal guidance, interpenetration-aware guidance
- fit the prior to the positional maps with projection-based constraint for reconstruction

# From Images to Videos
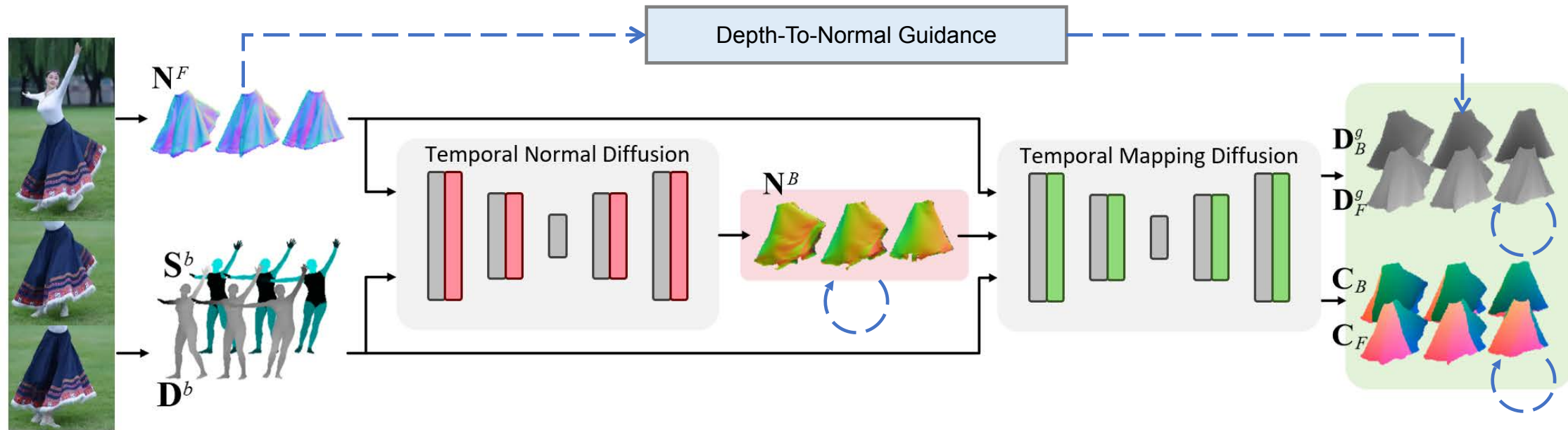
To handle video cases, we
- introduce temporal diffusion models
- enforce geometric and temporal guidance
- temporal consistency guidance, depth-to-normal guidance, interpenetration-aware guidance
- fit the prior to the positional maps with projection-based constraint for reconstruction

# From Images to Videos

To handle video cases, we
- introduce temporal diffusion models
- **enforce geometric and temporal guidance**
- temporal consistency guidance, **depth-to-normal guidance**, interpenetration-aware guidance
- fit the prior to the positional maps with projection-based constraint for reconstruction
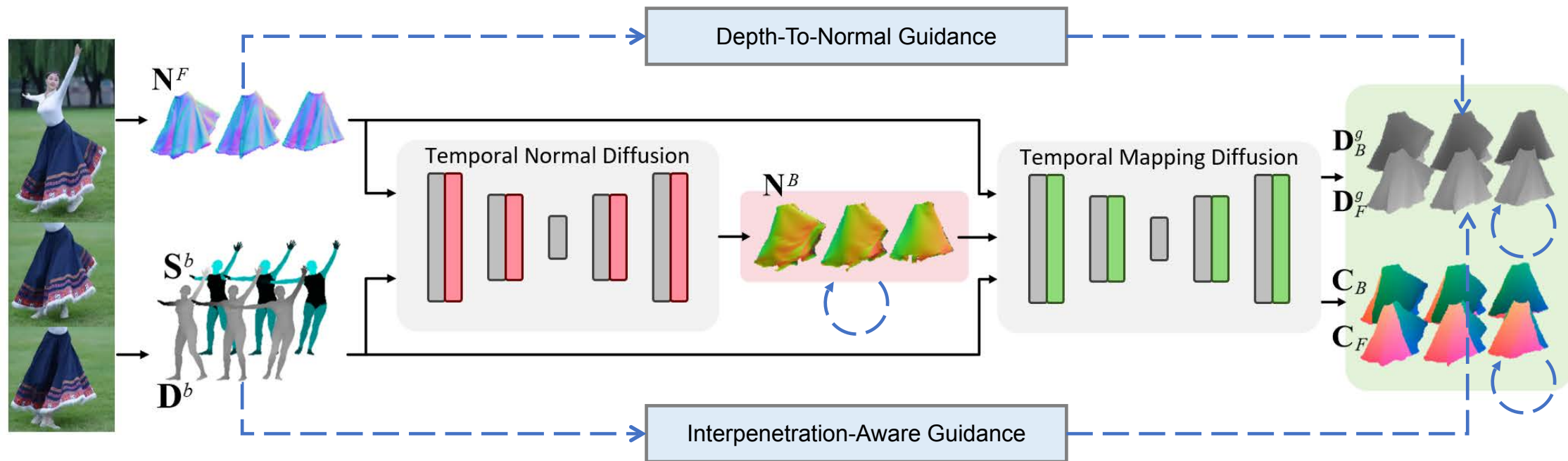
# From Images to Videos

To handle video cases, we
- introduce temporal diffusion models
- **enforce geometric and temporal guidance**
- temporal consistency guidance, depth-to-normal guidance, **interpenetration-aware guidance**
- fit the prior to the positional maps with projection-based constraint for reconstruction
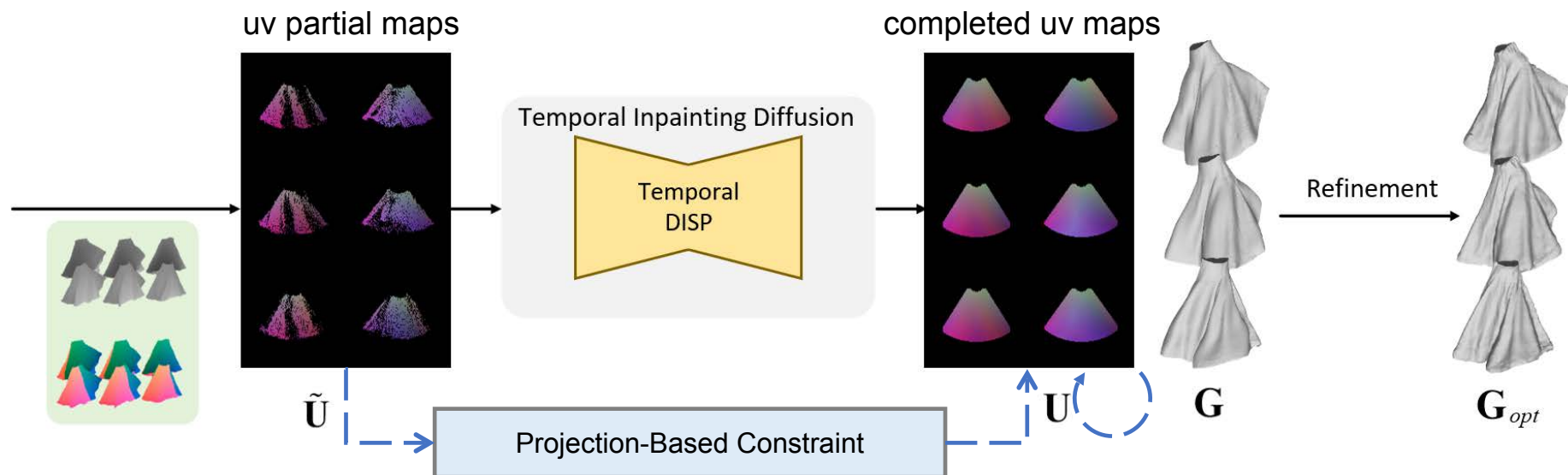
# From Images to Videos

To handle video cases, we
- introduce temporal diffusion models
- enforce geometric and temporal guidance
- temporal consistency guidance, depth-to-normal guidance, interpenetration-aware guidance
- fit the prior to the positional maps with projection-based constraint for reconstruction

# Recovered Garments from Videos



Input          Front-view          Back-view          Projection

# What about Long Robes?

# What about Long Robes?



- Quite good but not quite right when seen from the side.
- From a physical point of view, the 3D pose is not realistic.
➡ There still is work to do!

PhD anyone?