# Fabric Material Recovery from Video Using Multi-Scale Geometric Auto-Encoder

Junbang Liang[1] and Ming Lin[2]

[1] Amazon
[2] University of Maryland, College Park

**Abstract.** Fabric materials are central to recreating realistic appearance of avatars in a virtual world and many VR applications, ranging from virtual try-on, teleconferencing, to character animation. We propose an end-to-end network model that uses video input to estimate the fabric materials of the garment worn by a human or an avatar in a virtual world. To achieve the high accuracy, we jointly learn human body and the garment geometry as *conditions to material prediction*. Due to the highly dynamic and deformable nature of cloth, general data-driven garment modeling remains a challenge. To address this problem, we propose a two-level auto-encoder to account for both *global* and *local* features of any garment geometry that would directly affect material perception. Using this network, we can also achieve smooth geometry transitioning between different garment topologies. During the estimation, we use a closed-loop optimization structure to share information between tasks and feed the learned garment features for temporal estimation of garment materials. Experiments show that our proposed network structures greatly improve the material classification accuracy by 1.5x, with applicability to unseen input. It also runs at least *three orders of magnitude* faster than the state-of-the-art [66, 68]. We demonstrate the recovered fabric materials on virtual try-on, where we recreate the entire avatar appearance, including body shape and pose, garment geometry and materials from only a single video.

**Keywords:** Fabric material estimation, Synthetic dataset

## 1 Introduction

Human appearance reconstruction is one of the key techniques for building a vivid, interactive virtual world. However, especially for fabric material estimation, has been under-explored due to the complexity and the diversity of cloth dynamics and coupled interaction with an avatar body. Image features are often sparse, containing many noisy signals regarding the fabric materials worn on the body. An effective way to amplify useful signals is to estimate garment geometry from images as a by-product. However, this is an open challenge due to several reasons. First, garments have highly dynamical geometry that is not easy to capture and model. Previous works on garment modeling [23, 44, 62] and estimation [7, 18, 43] often propose solutions on one single type of garment, mostly t-shirts. Although the methods are also applicable to other garments, lack of generalization in capturing different garment geometries presents a considerable barrier for virtual try-on applications: users can only choose one of few pre-trained garment types and are not able to import new ones easily. Second, accurate estimation is

often hindered by view projection, body occlusion, and limited availability of 3D scanning. For example, the human-body estimation network may disagree with the garment reconstruction network in skeleton orientation due to the projection ambiguity (*e.g.* an arm is posed forward vs. backward), resulting in prediction misalignment. Thus, without a general garment representation and an accurate geometry estimation, it is very difficult to regress the fabric materials solely from images.

In this paper, we introduce a learning model that addresses these issues, achieving both garment geometry and fabric material estimation simultaneously from commonly available video inputs for virtual try-on. To handle the dynamic geometry and different topologies of the garments and to provide a unified parametric model for the garments, we propose a two-level auto-encoder network. The key observation is that classical point cloud encoders such as PointNet [46] are great for capturing global shapes, but not suitable for encoding the local details. Multi-scale feature extraction decomposes the problem into smaller partitions and also decouples global and local features to enable larger coverage on local shape learning and capturing local topology transitions. During the estimation, we couple the human body inference with the garment recovery to maximize the estimation accuracy of the two correlated tasks. Other than traditional multi-tasking, we further introduce a closed-loop structure so that the garment features of different scales can guide the body estimation to improve the accuracy for both. Based on the temporal change of garment features, we can perform accurate material classification accordingly. Our key contributions include:

- The first neural network for fabric material recovery of a garment from a RGB video (Section 3);
- A novel two-level auto-encoder for learning the latent space of garments through multi-scale feature coupling, resulting in higher accuracy for material parameter estimation (Section 4);
- Joint estimation of human body and apparels through a close-loop iterative optimization that can account for arbitrary topologies of garments and ensure geometric consistency (Section 5);
- A large dataset of garment motion sequences with wide variations of human body, fabric materials, textures, and lightings for virtual try-on (Appendix A).

Our experiments show that the proposed network structure effectively increases the performance and accuracy of the virtual fabric material estimation. By using only a few frames of a person wearing a garment, our model can faithfully reconstruct the garment fabric material(s), using the recovered shape and motion of both the garment and the avatar body as the conditioning in virtual fabric material estimation.

## 2    Related Work

**Fabric material estimation.** Researchers have been tackling different inverse problems, including inverse cloth design [14], combinatorial material design [10], BRDF parameter capturing [60], weaving pattern reconstruction [22], human material perception [8, 9], and frictional coefficient estimation [42, 48]. Cloth material estimation is among the most challenging due to cloth's highly dynamic motions. Previous works study the task in a simplified and constrained scenario, and recover the materials using statistical observation [11, 16], optimization [41, 66], or learning [6, 67]. In contrast,

our method learns fabric materials from videos of a human wearing garments in more general and widely applicable framework, assuming commonly available inputs like image sequences and videos. More importantly, *our method makes use of the estimated multi-scale garment latent codes as input signal that is shown to be more effective in recovering overall garment geometry with local details than merely image features.*

**Garment modeling and estimation.** Garment geometry capturing or recovery has been widely studied: non-learning methods using symmetry and user input [72], optimization [27,66,68], or binocular data [12]. Recently, methods using deep learning have been proposed for faster speed and more convenient usage [4,7,15,18,23,24,26,28,43,69,74]. In addition, direct garment modeling methods have also been proposed using spherical parameterization [44] for estimation or displacement map [57] for retargetting.

Different from displacement-based cloth representation [7], PCA-based models [62], and mesh-CNN-based methods [23], our garment model is universal to all topologies, applicable to those not homotopic to human surfaces (*e.g.* long dresses), enabling semantic interpolation between different garments. Compared with [44], our model generates a stand-alone garment mesh that is easy to export and retarget. Our method is the *first network that jointly estimates the garment material and geometry* for virtual try-on. Our method is substantially different from most garment capturing or generation methods [34, 52, 69] regarding model input, output, and assumptions. *It does not required 3D scanning, which is often not easily available, but only videos that can be easily captured using mobile devices.*

**Point cloud encoder and decoder.** PointNet [46] was among the first network model for encoding an unordered point set. Follow-on improvements include spatial partition [32, 35, 47], edge convolution [63], local region filtering [45, 73], and analogous convolutional operators [36]. Although these recent works have utilized hierarchical structure to some extent, their methods are not sufficient for auto-encoding the garment geometry or topology. The key difference between garment auto-encoding and rigid gadgets auto-encoding is that there are a large number of local details (*e.g.* wrinkles) due to cloth's highly deformable and dynamical nature. As a result, latent codes for local details are necessary.

Recently, [39, 70] use similar ideas on point-based garment geometry modeling. The main differences between these methods and ours include (1) our garment latent code is independent and does not need a body point-cloud to morph on, (2) their latent code is randomly assigned before training, which is not ideal for learning a compact latent space for smooth interpolation, and could result in noisy estimation results when applied to downstream task, and (3) their local and global patches depend heavily on human body surface, which makes it difficult for the resulting point cloud to represent loose dresses.

**Human reconstruction from images.** Human estimation using RGB images has been a popular research topic in deep learning for its importance in virtual reality and computer animation. While early works propose network models for only 2D/3D body skeletons [13, 40, 64], more recent works introduce techniques to regress the entire human body – either using a parametric human model [3,29] or voxel-based representation [50, 58,71]. Given the fact that the annotations in most real-world datasets contain only joint positions, the learning process has been refined in various ways [2, 33, 37, 54, 65].

In order to estimate the fabric material, we need to recover the garment shape on the human body, which is an important problem rarely addressed in avatar reconstruction. In our pipeline, we use state-of-the-art human body predictions *as a strong prior* for the garment estimation module. Given the focus of this paper, we assume to use video input for only garment material recovery. We refer the interested readers to a recent survey [21] for comprehensive review on video-based pose estimation instead.

## 3   Method Overview

We first give the formal problem definition. Given a video clip showing a person moving (*e.g.* walking, jumping, bending, etc), we estimate the fabric materials of the garment worn by the person. We assume that the garment worn is made of the same material. By *fabric materials*, we refer to the physical material parameters used in cloth simulation. We adopt the same material parameter definition introduced in [61], which consists of 24 parameters for stretching stiffness and 15 parameters for bending stiffness.

Given the fact that the differences of the material parameter values do not intuitively reflect the human visual perception, we follow the previous work [67] to discretize the material parameter space based on the amount of deformations due to external forces. Using sensitivity analysis [51], the stretching stiffness is split into 6 classes and 9 for the bending. Combining both dimensions will yield 54 different material classes. As confirmed by [67], these 54 classes cover most of the common materials, including polyester, cotton, nylon, rayon, and their combinations. For example, one type of materials named *'white-swim-solid'* consisting of 87% nylon and 13% spandex, as measured by [61], fits in the discrete classification model with the stretching label of 2 and the bending label of 3.

In this paper, we introduce a deep neural network (Fig. 1) for simultaneously estimating the garment geometry and its material type(s), along with the human body. Our key idea is that image features are not sufficient for inferring garment materials; it is necessary to extract the garment geometry as well for a more accurate estimation. To support different topologies of garments, we choose point clouds for its geometry representation. To better account for the highly dynamic garment surfaces, we train a two-level point cloud auto-encoder (Sec. 4) so that it can learn the global shapes and local features of the garment to reduce the total number of degrees of freedom. We use the SMPL model (see [38] for its rigorous math definition) to represent the human pose and shape.

We divide the estimation pipeline into two phases. First, we estimate the human body and the cloth geometry in a frame-by-frame manner (Sec. 5.1). A closed-loop optimization structure is used to improve the estimation accuracy of these two correlated tasks. The garment geometry prediction module is conditioned on the human body parameters, and at the same time provides corrective feedback to the human body prediction module. We then feed the features of the image and the garment geometry from each frame together to a temporal neural network for the garment material estimation (Sec. 5.2). By sharing common features, providing corrective feedback, and conditioning on outputs of closely-related tasks, our network model can achieve higher estimation accuracy on all three tasks than independent estimation baselines.
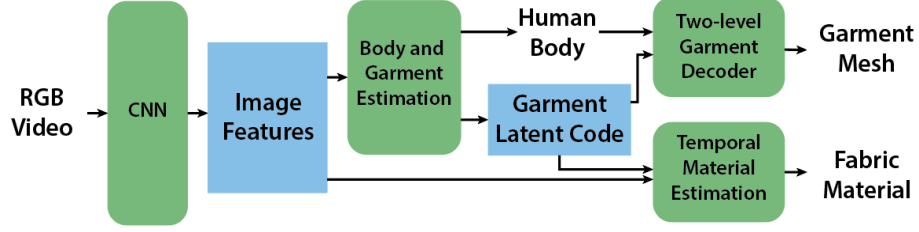
**Fig. 1. Overall network structure**. Given an RGB video, we extract its image features and estimate the body and garment shape frame by frame (Sec. 5.1). The latter is decoded to obtain a garment mesh (Sec. 4). The temporal sequences of image and garment are fed to an LSTM for material classification (Sec. 5.2).

## 4    Garment Auto-encoder

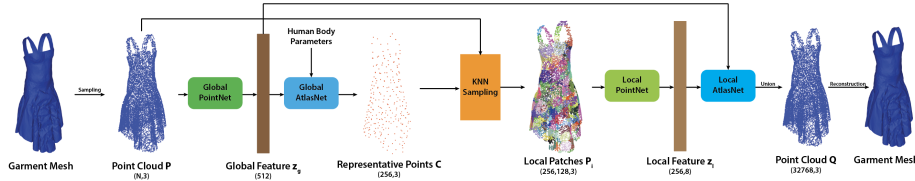We first set up an auto-encoder for the cloth model. Since the model is designed not to



**Fig. 2. The network structure of the garment auto-encoder**. The point cloud sampled from the original mesh is first fed to a global PointNet for coarse shape features. Its representative point set is then obtained by decoding the global features from an AtlasNet. From those points, we sample the local patches using K-nearest neighbor and pass them to a local PointNet for detailed shape features. The local decoder is then conditioned on the global latent code and the corresponding patch center to recover the patches that are stitched together to form the reconstructed point cloud.

assume fixed garment topology, we choose to use point clouds as the underlying representation. Other representations, such as graph-based [55] or displacement-based [7, 43], rely on either fixed graph structure, or fixed human surface, thus not applicable for generalization to different garments. The use of auto-encoder here is necessary because the degrees of freedom (DoF) for point clouds are too high for estimation. An encoder-decoder structure can effectively reduce the DoF and retain only the essential information, such as the global shape and the local details. More importantly, it clusters similar shapes to similar latent codes, which is beneficial to the estimation module. As later shown in the Appendix, our model provides smooth transitioning between different topologies by using simple interpolation between latent codes. We discuss details on recovery methods from point cloud to mesh in the Appendix.

### 4.1    Two-Level Encoder-Decoder Structure

Previous point cloud auto-encoders such as AtlasNet [19] use Multi Layer Perceptron (MLP) to transform a 2D patch to a set of 3D points in the space. Their method performs

well in point cloud datasets that include rigid objects, such as airplanes or chairs, since the deformations presented in those objects are simple and regular. However, it cannot be directly applied to learn garment point clouds, since garments have a much larger variance in point cloud distribution due to its dynamic nature. For example, a simple dress can create different wrinkle structures under different external forces. As a result, one global auto-encoder cannot account for all detailed structures, resulting in overly smoothed point clouds. Recently, [5] proposes a method to resolve patch overlapping and collapsing occurred in AtlasNet, but it still cannot account for arbitrary topologies and detailed wrinkles.

We propose a two-level auto-encoder for learning the latent space of the cloth. As shown in Fig. 2, we use a set of representative points $\mathbf{C}$ to express the global shape of the garment, and sample around them to form local patches, which are encoded independently to account for local shapes. Specifically, given a point cloud $\mathbf{P}$, we first pass it through a global auto-encoder to form a representative point cloud:

$$\mathbf{C} = D_g(E_g(\mathbf{P}), \theta) \tag{1}$$

where $E_g$ and $D_g$ are the global encoder and decoder, and $\theta$ is the human body parameter. Next, we use K-nearest-neighbor to sample points around the representative ones:

$$\mathbf{P}_i = KNN(\mathbf{P}, \mathbf{c}_i) \tag{2}$$

where $\mathbf{c}_i$ is the i-th element in $\mathbf{C}$, and $\mathbf{P}_i$ is the i-th patch. This step forms local patches around the representative points. Finally, we pass each patch to the shared local auto-encoder, and do a union operation to obtain the reconstructed point cloud:

$$\mathbf{Q}_i = D_l(E_l(\mathbf{P}_i), \mathbf{z}_g, \mathbf{c}_i) \tag{3}$$

$$\mathbf{Q} = \bigcup_i \mathbf{Q}_i \tag{4}$$

where $\mathbf{Q}_i$ and $\mathbf{Q}$ are the reconstructed patches and point cloud, $D_l$ and $E_l$ are the local decoder and encoder, and $\mathbf{z}_g = E_g(\mathbf{P})$ is the global latent code.

### 4.2   Representative Point Set Extraction

Note that Eq. 3 and 4 imply that the representative points $\mathbf{C}$ have to be in the same order as the local latent codes $\mathbf{z}_l$. This is the key reason why traditional methods such as farthest point sampling [47] do not work: its ordering is very sensitive to the input, resulting in an unknown mapping between reconstructed patch centers $\mathbf{C}$ and the local patches $\mathbf{P}_i$ (thus the local latent code $\mathbf{z}_{l_i}$).

To resolve this issue, we encode the entire point cloud and compute the representative points using the decoder itself. Due to the continuous nature of the auto-encoder network, the continuity and consistency regarding similar point clouds are guaranteed, thus ensuring $\mathbf{c}_i$ to be exactly matched with $\mathbf{P}_i$.

### 4.3    Training Losses

During training, we use Chamfer Distance between two point clouds as the loss:

$$d(\mathbf{P}, \mathbf{Q}) = \frac{1}{|\mathbf{P}|} \sum_{\mathbf{p} \in \mathbf{P}} \min_{\mathbf{q} \in \mathbf{Q}} \|\mathbf{p} - \mathbf{q}\| + \frac{1}{|\mathbf{Q}|} \sum_{\mathbf{q} \in \mathbf{Q}} \min_{\mathbf{p} \in \mathbf{P}} \|\mathbf{q} - \mathbf{p}\| \tag{5}$$

In Eq. 6, we apply the Chamfer Distance loss between the representative point set and the point cloud to learn the global shape (first term), and the one between the recovered and the original point clouds, both patch-wisely (second term) and globally (third term) to capture the local details:

$$\mathcal{L}_{AE} = d(\mathbf{P}, \mathbf{C}) + \frac{1}{n} \sum_{i=1}^{n} d(\mathbf{P}_i, \mathbf{Q}_i) + d(\mathbf{P}, \mathbf{Q}) \tag{6}$$

## 5    Material Estimation

With the garment auto-encoder (Sec. 4) at hand, garment material estimation becomes tractable. We design our overall pipeline as shown in Fig. 3. Given the sequence of image frames, we first feed them one by one to a model for estimating the human body and cloth geometry. By predicting the latent vector instead of the exact positions of the point cloud, the single-frame estimation network avoids severe overfitting or producing irrational results, due to the reduction of the degree of freedom by the auto-encoder.

Next, we combine the image features as well as the estimated garment latent code as the temporal signals, which go through a canonical temporal network module (*i.e.* LSTM [25]) to predict the final material type. Since the latent space preserves similarity (*i.e.* positive correlation between distances of latent vectors and distances between the original point clouds), the motion of the estimated latent vector becomes a better indicator of garment motion than image features, which is beneficial to garment material learning. We do not include body features here because the garment material is directly related to the garment motion, which has already taken the human body as the condition (Sec. 5.1). We discuss more details of the network in the following sections.
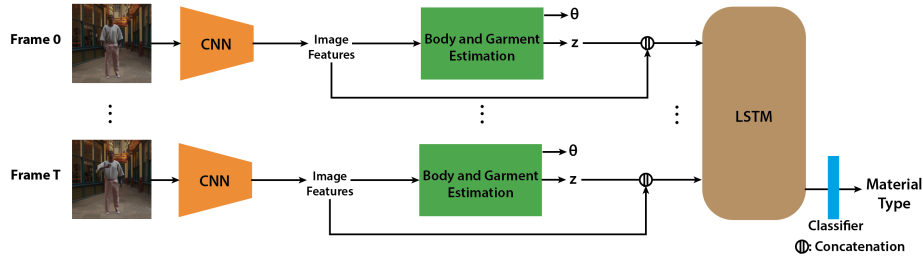


**Fig. 3. Our estimation pipeline**. Each video frame is first processed to obtain the image feature, the human body, and the garment shape. Then the image features are concatenated with the garment latent code as input to LSTM for material recovery.
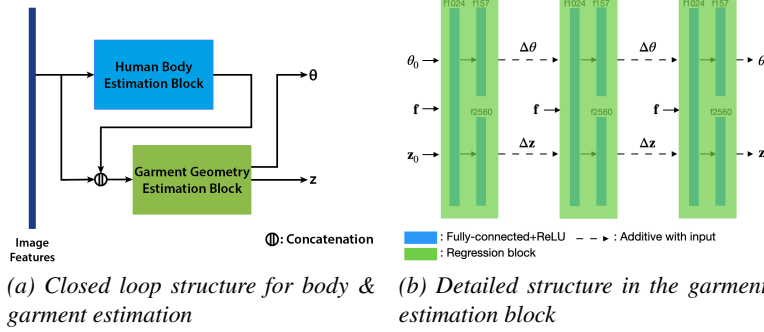
*(a) Closed loop structure for body & garment estimation*

*(b) Detailed structure in the garment estimation block*

**Fig. 4. The network structure for body and garment estimation in each frame**. (a) The garment shape estimation block takes the human body parameters as a prior, but also provides a feedback correction. (b) The garment estimation module consists of three identical, shared-weight blocks, each of which takes image features $\mathbf{f}$ and current predictions of the human body $\theta_0$ and garment $\mathbf{z}_0$, and outputs the corrective values.

### 5.1   Single Frame Closed-Loop Estimation

As shown in Fig. 4, we train a model to estimate the human body and cloth geometry given one single frame. Formally, in each frame, we are given the image features $\mathbf{f}$. We first go through a state-of-the-art body estimation block [33], $HB$, to get a first-hand body estimation, $\hat{\theta} = HB(\mathbf{f})$, where $\hat{\theta} = [\theta, \beta]$ are the human body parameters including pose and shape. In the garment estimation block, we take as input $\mathbf{f}$ together with $\hat{\theta}$ and regress the garment latent code $\mathbf{z}$, consisting of both $\mathbf{z}_g$ and $\mathbf{z}_l$, and body parameters $\theta$. Inside the garment estimation block, we use three shared-parameter small regression blocks, $RB$, to iteratively provide the correction, given the current estimation:

$$\theta_0 = \hat{\theta} \quad \mathbf{z}_0 = \mathbf{0} \tag{7}$$

$$\Delta\theta_i, \Delta\mathbf{z}_i = RB(\theta_{i-1}, \mathbf{z}_{i-1}) \tag{8}$$

$$\theta_i = \theta_{i-1} + \Delta\theta_i \quad \mathbf{z}_i = \mathbf{z}_{i-1} + \Delta\mathbf{z}_i \tag{9}$$

Overall, the garment estimation block forms a closed-loop structure, in which the human body parameters are required to predict the garment, and are later corrected back by the garment prediction as well.

The key insight of our module design is that the human body and garment shape are highly correlated at different scales and should be jointly learned using shared information. On the global scale, the detailed features of the garments restrict the variance of the human body and reduce ambiguity due to camera projection. On the local scale, the body pose and shape largely defines the valid distribution of the garment wrinkle positions. Our proposed structure is also analogous to iterative optimization and feedback control in other areas, where two objectives serve as prior knowledge of each other and are improved iteratively. This work is the first to introduce this idea for the human and garment joint estimation task.

The loss function for the single-frame estimation is defined as:

$$\mathcal{L}_s = \mathcal{L}_{body} + \mathcal{L}_{AE} \tag{10}$$

$$\mathcal{L}_{body} = \mathcal{L}_{2D} + \mathcal{L}_{3D} + \mathcal{L}_{SMPL} \tag{11}$$

where $\mathcal{L}_{2D}$, $\mathcal{L}_{3D}$, and $\mathcal{L}_{SMPL}$ represent the 2D joint loss, the 3D joint loss, and the body parameter loss defined to supervise the human body estimation [37], and $\mathcal{L}_{AE}$ is the Chamfer distance defined in Eq. 6 to supervise the garment estimation.

### 5.2   Temporal Estimation for Garment Material

Garment material estimation is challenging since the visual difference of different materials is subtle and can easily be overwhelmed by disturbance, *e.g.* various directions or magnitudes of external forces. To tackle this problem, previous works often assume fixed environment settings and cloth shapes [66, 67]. While we follow a similar principle when training the material estimation module, we go one step further that we only assume common human motion for driving the garment instead of the whole external force field. While previous works [66, 67] can only handle videos of a piece of cloth hanging and dragged by the wind, our method possesses a wider applicability regarding the diversity of the garment shapes, sizes and human motions in the input video, which for the first time enables practical usages for garment material cloning.

As shown in Fig. 3, we collect and concatenate the image features and the estimated garment latent vector of each frame as the input signal, and feed the sequence of the signals to LSTM to produce a summary feature. Finally, we pass the summary feature to a fully-connected layer for material type classification. We use the cross-entropy loss for supervision.

Training the entire pipeline from scratch is not ideal because the system is too large and the training could be unstable. Instead, we first trained the single frame body and garment estimation module using single view images. After the convergence on the single frame module, we fixed its parameters and applied it to train the material estimation network. Our experiments demonstrated in Sec. 6 indicate that the multi-scale garment features are not merely useful for detecting the fabric materials; they are the dominant features during the estimation and can boost the test accuracy compared to methods that only use the image features.

## 6   Results

We demonstrate the performance of our model as follows. (1) We compare our work with with the baselines and SOTA methods *quantitatively* (Sec. 6.1) and *qualitatively* (Sec. 6.2). (2) Ablation studies are presented in Sec. 6.2 on the improvements by our network design. (3) A user study on material perception using this work and the similarity between the measurements of real-world fabrics from lab experiments vs. our predicted garment materials from videos is presented in Sec. C, with detail in Appendix. (4) We compare ours with other related learning-based methods (Sec. D) and show application to virtual try-on (Sec. E), with detail in Appendix. (5) Additional latent code interpolation between garments, training data and perceptual study examples can be found in both Appendices and the supplementary video.

**Training process.** We first trained the auto-encoder alone. Next, we trained the single-frame estimation module, with the fixed decoder attached at the end. Finally, we trained the material estimation with other parts fixed. See Appendix for more training details.

### 6.1   Quantitative Analysis

**Ours vs. Image-Only [67].** Due to the difference regarding the input distribution (dressed garment on a human body in our method vs. hanging cloth in theirs), we re-train their model on our datasets for a fair comparison. We study the contribution of image-only features vs. garment-only features, as well as CNN vs. LSTM (that exploits the temporal coherence). Finally, we compare the overall performance difference between ours and [67]. The test classification accuracy is reported in Table 1.

_Findings_: (1) While all three models have learned the relationship between motion and materials and all three outperform random guess, *the garment feature signals are shown to be much more important than the image features.* This finding is not surprising, since the garment shape is directly affected by the material. (2) Combining the two features, as our model does, further improves the test accuracy. A possible reason is that an overall capturing of the garment shape (*e.g.* width and length of the entire piece), which is difficult to retrieve using garment latent codes, could be more easily extracted using image features. (3) By exploiting temporal coherence, unsurprisingly all three versions of the model achieve better accuracy than only using 1 image.

| Method | Mean Accuracy | Temporal Gain | Garment Features Gain |
|---|---|---|---|
| Random guess | 1.85% | - | - |
| Image only, CNN | 5.11% | 40.16% | - |
| Image only, LSTM [67] | 45.27% | | - |
| Garment only, CNN | 11.85% | 53.31% | 6.74% |
| Garment only, LSTM | 65.16% | | 19.89% |
| Image + Garment, CNN | 12.62% | **57.52%** | 7.51% |
| Image + Garment, LSTM **(ours)** | **70.14%** | | **24.87%** |

**Table 1. Comparison on material estimation**: our method achieves ∼1.5x higher accuracy (45.27% vs. 70.14%) in material identification than [67].

**Ours vs. Optimization-based [66, 68]:**  An optimization-simulation framework to obtain the fabric material parameters using wrinkle density of the garment in a single image was proposed in  [66, 68]. In contrast, our method extracts both static image features and spatio-temporal *garment features* across frames. We generate the same set of test scenes as shown in the Appendix. Our model is tested on these sequences under varying lighting and visibility conditions (Appendix, Fig. 11); the average accuracy is reported in Table 2. In this challenging case where the lighting condition and the textures are not seen in the training distribution, our method still achieves comparable accuracy with previous method [66, 68], but it runs **more than 1,000x faster**. *Ours is the first learning-based method to predict fabric materials directly from a video of garments worn on a human body*.

### 6.2   Qualitative Results

We compare ours with the most relevant work of [66] for joint estimation of garment shapes and materials, as shown in Fig. 5. Our method achieves similar reconstruction accuracy and visual quality as [66, 68]. But, [66, 68] uses semantic segmentation, thus suffering from tedious manual processing and long inference time. In contrast, our learning-based method is fully automatic and can compute the prediction in real time. Moreover, our method does not assume the sewing patterns as a prior.

| Method | Accuracy (%) | | | | | | Speed |
|---|---|---|---|---|---|---|---|
| | Mid-day | | | Sunset | | | |
| | T-shirt | Pants | Skirt | T-shirt | Pants | Skirt | |
| [66, 68] | 80.2 | 80.2 | **83.3** | 81.6 | 79.9 | **80.7** | 4-6 hours |
| Ours | **86.5** | **91.6** | 81.6 | **82.4** | **91.6** | 79.6 | **8.7 sec** |

**Table 2. Quantitative comparison with [66, 68]**. Our method achieves comparable or higher accuracy, but runs at least *three orders of magnitude* faster than the state-of-the-art [66, 68].



(a) Input image (b) Results from [68] (c) Our results

**Fig. 5. Qualitative comparison with [66, 68]**: Ours is easier to use and achieves visually comparable reconstruction much faster without priors on garment patterns and topology.

We further compare with several learning methods [7, 24] in the Appendices for reference. Many often use additional information (*e.g.* mesh templates or known garment types) as priors, so direct visual comparison is not meaningful. Nonetheless, our model successfully generalizes to unseen real-world images/videos with comparable visual results, as shown in Appendices. During these experiments, our method is directly applied without any fine-tuning or post-optimization. Although trained using synthetic datasets, our model correctly identifies people and the garments from real-world images, and achieves similar visual results in all examples, when compared with previous works. The network is also capable of the predicting correct sizes of garments relative to the body, due to multiscale auto-encoders.

**Material Cloning for Virtual Try-On:** we show three application scenarios of our method. In Fig. 6, given an RGB video of a person wearing garments of different fabric materials, our method can identify the underlying material and *clone* it onto other garment models using cloth simulation. Our method is the first to achieve fast and accurate material extraction from videos of dressed garments on a body. We further show the ability of our method to reconstruct the entire human appearance from the input video using one single network. We first estimate the body and the garment geometry frame by frame, and use the temporal information to infer the material. The three parts are combined using cloth simulation to generate the final output. Fig. 7 shows the reconstruction results (also see the supplement video). Our reconstructed garment shapes and wrinkles match those in the input video frames.

### 6.3   Ablation Study

We verify the effectiveness of our network model in the following ablation studies reported in Table 3. We compare our method with baselines that replaces (a) the two-level auto-encoder with AtlasNet [19], and (b) the joint body-garment estimation block with a parallel estimation structure, respectively. The metrics include both the reconstruction accuracy and the material classification accuracy in the final stage. Our method results in notably smaller errors than both baselines in reconstruction and material prediction. See Appendix on the details of this ablation study.
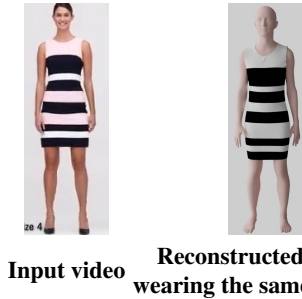
**Input video**        **Reconstructed Avatar wearing the same garment**

**Fig. 6. Material transfer between videos**. Our method can take videos of a person wearing any garments and clone the underlying fabric materials onto a virtual avatar wearing the same garment with the same fabric.
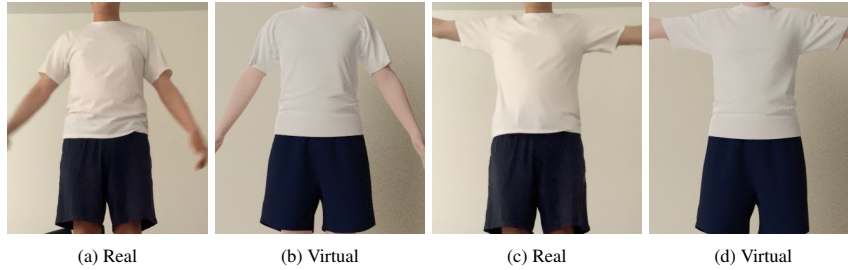


(a) Real                    (b) Virtual                    (c) Real                    (d) Virtual

**Fig. 7. Qualitative results:** our method faithfully recovers the T-shirt materials (a, c) in video so that the wrinkles around the simulated t-shirt sleeves (b, d) appear similar under different poses.

| Method | CD | SD | Accuracy | Method | MPJPE | CD | Accuracy |
|---|---|---|---|---|---|---|---|
| AtlasNet [19] | 0.31 | 8.05 | 54.20% | Separate | 80.89 | 1.55 | 49.15% |
| Ours | **0.12** | **1.03** | **70.14%** | Ours | **55.20** | **0.88** | **70.14%** |
| (a) Auto-encoder | | | | (b) Human and garment estimation | | | |

**Table 3. Ablation study for different parts of our proposed network.** CD stands for errors in Chamfer Distance; SD stands for Sinkhorn Divergence [20]; and MPJPE stands for Mean Per Joint Position Error – all in millimeters (mm). *Separate* predicts the body and the garment separately in parallel branches. Our method results in notably smaller errors (in CD, SD, MPJPE) than all baselines in reconstruction and material prediction, with 30% to 40% higher accuracy.

We further compare our method with a recent work [7] since their task is the most relevant to ours. As shown in Table 4, our model, without any fine-tuning or domain adaptation, achieves the smallest error – whether the ground-truth human body model is provided or not during reconstruction. Since our model is trained on a wider range of body poses and garment types than theirs while achieving better accuracy, it has shown to offer good generality to unseen inputs. See Appendix for more detail.

## 6.4   Lab Experiments and User Study

In this experiment, we test the prediction accuracy of our method using real-world materials. We used five real-world materials measured from lab experiments [61], which are sampled from sweater, t-shirt, tablecloth, jeans, and blanket, respectively. The measured values are then compared with the one predicted from our method, reported in Table 5. Our method achieves a relatively small error between 9.5% to 16.7%. See Appendix for more detail.

| Methods | [2] | [7] | | Ours | |
|---|---|---|---|---|---|
| | GT Pose | GT Pose | Full Pred. | GT Pose | Full Pred. |
| Pants (mm) | 5.44 | 5.57 | 10.16 | **1.58** | **3.08** |
| Short Pants (mm) | 8.23 | 5.97 | 10.00 | **4.92** | **5.69** |
| T-shirt (mm) | 5.80 | 5.63 | 11.97 | **1.67** | **3.08** |
| Shirt (mm) | 5.71 | 6.33 | 9.05 | **2.29** | **3.75** |
| Coat (mm) | 5.85 | 5.66 | 9.09 | **2.84** | **3.65** |

**Table 4. Test errors on the Multi-Garment Net datasets [2,7]**. Our method achieves the lowest errors (by up to **3x–4x**) across all garment types, without any fine-tuning or reference body.

| Material Name [61] | Stretching Ratio (GT/Prediction) | Bending Ratio (GT/Prediction) | Mean Relative Error |
|---|---|---|---|
| gray-interlock | 1.01/1 | 1.6/2 | 10.5% |
| navy-sparkle-sweat | 0.56/0.5 | 1.7/2 | 12.8% |
| white-dots-on-blk | 15.8/20 | 3.5/4 | 16.7% |
| 11oz-black-denim | 3.6/3 | 3.0/3 | 8.3% |
| pink-ribbon-brown | 2.93/3 | 12/10 | 9.5% |

**Table 5. Lab experiment results.** Our material estimation achieves relatively small errors compared to lab measurements on all real-world materials tested.

| Method | Input | Dependencies | Generality | Dresses support | Separate mesh | Material Estimation |
|---|---|---|---|---|---|---|
| MGN [7] | Semantic seg. + 2D joints | Garment correspondences | One model per garment | No | Yes | No |
| DeepCap [24] | Foreground seg. | Template mesh | One model per garment | Yes (w/ known template) | No | No |
| [66] | Semantic seg. | Template mesh | One model per garment | Yes (w/ known template) | Yes | Yes |
| DeepFashion3D [74] | RGB frame (garment only) | None | One model for all | Yes (limited topologies) | Yes | No |
| Tailornet [43] | Body parameters | Garment correspondences | One model per garment | Yes (limited topologies) | Yes | No |
| BCNet [28] | RGB frame | Garment correspondences | One model per garment | Yes (limited topologies) | Yes | No |
| SIZER [57] | Body scan | Garment labels | One model per garment | No | Yes | No |
| ARCH [26] | RGB frame (foreground only) | None | One model for all | Yes (water-tight) | No | No |
| Ours | RGB frame | None | One model for all | Yes | Yes | Yes |

**Table 6. Comparison with previous works**. Our method can handle the largest set of garments, using the fewest possible information (*i.e.* widely available RGB images only), in one stand-alone network.



**Fig. 8. Material transfer examples.** Our method can accurately estimate the material parameters from input videos (a, e) and replicate the same 'feel' in other animations (b, c, d) and (f, g, h), respectively, creating significantly different visual effects for the same pose.

**Perceptual Validation:** To further validate and quantify the material similarity, we conduct a user study to examine how close our estimation results are to the ground-truth data in human perception. Our results show that the average similarity ratings for five tested materials vs. the ground-truth data are all larger than 5, ranging from 5.7 to 8.5, with an overall mean value of 7.1. These indicate that our method indeed can recover

fabric materials with only minor perceptible differences to the real-world materials. Furthermore, we also conducted on material perceptions under different environmental conditions. Please see Appendix for details on these perceptual studies.

### 6.5   More Comparison with Previous Works

In Table 6, we extensively compare our work with previous ones regarding different assumptions, functionalities, and abilities. We define 'one model per garment' in 'generality' as that the method needs to create extra templates or registrations to the body, or need to retrain part of the network in order to predict a different garment type. Although DeepFashion3D [74] and ARCH [26] also have generality to different topologies to some extent, there are still limitations in their pipeline. The output from Deep-Fashion3D has to be continuous in one body part, meaning that they cannot support all topologies (*e.g.* dresses with holes). ARCH does support different garments on the body, but the output is a water-tight mesh together with the body, which is not always convenient for certain applications like virtual try-on. In contrast, our method naturally supports all kinds of topologies, and predicts the body and the garment in separate meshes. Additional comparison results with DeepCap [24] and MGN [7] can be found in the Appendix D and Fig. 13.

**Virtual Try-on:**   Visual results of our work on application to virtual try-on are shown in Fig. 8 and Fig. 12. More animations are shown in the supplementary video.

## 7   Conclusion

In this paper, we introduced a learning model for garment material estimation using RGB videos. We do not assume other inputs (e.g. segmentation, 3D scans, multi-views, etc.) or any prior knowledge on the garment shape/topology, design patterns/templates, or correspondences. We extract the multi-scale features to effectively represent the dynamic geometry structure of garments, which can be combined with image features to estimate fabric materials by learning their temporal patterns, while improving the human body reconstruction using a feedback loop. This approach is perhaps the first to introduce a unified parametric model for all garment types, and it can thereby support garments of different topologies without the need to retrain different models. Experiments show that our method achieves much higher accuracy up to 70.14% in estimating fabric materials than prior works, while offering capabilities in recovering garment types and topologies with generality and simplicity for an unification of multiple correlated tasks.

**Limitations:**   We assume that garment motion is captured as videos of adequate image resolution under sufficient lighting to show fabric movement, wrinkles and folds. The current implementation does not support multi-layer, folded garments, or detection of different materials at once. These issues can likely be addressed by adding more structural prior to encode multi-layer clothing, introduction of curvature representation for multi-fold features, and a point cloud segmentation module. The accuracy of fabric material estimation perhaps can probably be further enhanced by integrating neural rendering [56] with this work.

# References

1. Hdri haven. `https://hdrihaven.com/` (2020)
2. Alldieck, T., Magnor, M.A., Bhatnagar, B.L., Theobalt, C., Pons-Moll, G.: Learning to reconstruct people in clothing from a single RGB camera. In: IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019. pp. 1175–1186. Computer Vision Foundation / IEEE (2019). https://doi.org/10.1109/CVPR.2019.00127, `http://openaccess.thecvf.com/content\_CVPR\_2019/html/Alldieck\_Learning\_to\_Reconstruct\_People\_in\_Clothing\_From\_a\_Single\_RGB\_CVPR\_2019\_paper.html`
3. Alldieck, T., Magnor, M.A., Xu, W., Theobalt, C., Pons-Moll, G.: Video based reconstruction of 3d people models. In: 2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018. pp. 8387–8397. IEEE Computer Society (2018). https://doi.org/10.1109/CVPR.2018.00875, `http://openaccess.thecvf.com/content\_cvpr\_2018/html/Alldieck\_Video\_Based\_Reconstruction\_CVPR\_2018\_paper.html`
4. Alldieck, T., Pons-Moll, G., Theobalt, C., Magnor, M.A.: Tex2shape: Detailed full human body geometry from a single image. In: 2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019. pp. 2293–2303. IEEE (2019). https://doi.org/10.1109/ICCV.2019.00238, `https://doi.org/10.1109/ICCV.2019.00238`
5. Bednarik, J., Parashar, S., Gundogdu, E., Salzmann, M., Fua, P.: Shape reconstruction by learning differentiable surface representations. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4716–4725 (2020)
6. Bhat, K.S., Twigg, C.D., Hodgins, J.K., Khosla, P., Popovic, Z., Seitz, S.M.: Estimating cloth simulation parameters from video (2003)
7. Bhatnagar, B.L., Tiwari, G., Theobalt, C., Pons-Moll, G.: Multi-garment net: Learning to dress 3d people from images. In: 2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019. pp. 5419–5429. IEEE (2019). https://doi.org/10.1109/ICCV.2019.00552, `https://doi.org/10.1109/ICCV.2019.00552`
8. Bi, W., Jin, P., Nienborg, H., Xiao, B.: Estimating mechanical properties of cloth from videos using dense motion trajectories: Human psychophysics and machine learning. Journal of vision **18**(5), 12–12 (2018)
9. Bi, W., Xiao, B.: Perceptual constancy of mechanical properties of cloth under variation of external forces. In: Proceedings of the ACM symposium on applied perception. pp. 19–23 (2016)
10. Bickel, B., Bächer, M., Otaduy, M.A., Lee, H.R., Pfister, H., Gross, M., Matusik, W.: Design and fabrication of materials with desired deformation behavior. ACM Transactions on Graphics (TOG) **29**(4), 1–10 (2010)
11. Bouman, K.L., Xiao, B., Battaglia, P., Freeman, W.T.: Estimating the material properties of fabric from video. In: Proceedings of the IEEE international conference on computer vision. pp. 1984–1991 (2013)
12. Bradley, D., Popa, T., Sheffer, A., Heidrich, W., Boubekeur, T.: Markerless garment capture. In: ACM SIGGRAPH 2008 papers, pp. 1–9 (2008)
13. Cao, Z., Simon, T., Wei, S.E., Sheikh, Y.: Realtime multi-person 2d pose estimation using part affinity fields. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 7291–7299 (2017)

14. Casati, R., Daviet, G., Bertails-Descoubes, F.: Inverse elastic cloth design with contact and friction. Ph.D. thesis, Inria Grenoble Rhône-Alpes, Université de Grenoble (2016)
15. Chen, X., Zhou, B., Lu, F.X., Wang, L., Bi, L., Tan, P.: Garment modeling with a depth camera. ACM Trans. Graph. **34**(6), 203–1 (2015)
16. Clyde, D., Teran, J., Tamstorf, R.: Modeling and data-driven parameter estimation for woven fabrics. In: Proceedings of the ACM SIGGRAPH/Eurographics Symposium on Computer Animation. pp. 1–11 (2017)
17. CMU: Carnegie-mellon mocap database. created with funding from nsf eia- 0196217 (2003), http://mocap.cs.cmu.edu/
18. Daněřek, R., Dibra, E., Öztireli, C., Ziegler, R., Gross, M.: Deepgarment: 3d garment shape estimation from a single image. In: Computer Graphics Forum. vol. 36, pp. 269–280. Wiley Online Library (2017)
19. Deprelle, T., Groueix, T., Fisher, M., Kim, V., Russell, B., Aubry, M.: Learning elementary structures for 3d shape generation and matching. In: Advances in Neural Information Processing Systems. pp. 7433–7443 (2019)
20. Feydy, J., Séjourné, T., Vialard, F.X., Amari, S.I., Trouvé, A., Peyré, G.: Interpolating between optimal transport and mmd using sinkhorn divergences. arXiv preprint arXiv:1810.08278 (2018)
21. Gong, W., Zhang, X., Gonzàlez, J., Sobral, A., Bouwmans, T., Tu, C., Zahzah, E.h.: Human pose estimation from monocular images: A comprehensive survey. Sensors **16**(12), 1966 (2016)
22. Guarnera, G.C., Hall, P., Chesnais, A., Glencross, M.: Woven fabric model creation from a single image. ACM Transactions on Graphics (TOG) **36**(5), 1–13 (2017)
23. Gundogdu, E., Constantin, V., Seifoddini, A., Dang, M., Salzmann, M., Fua, P.: Garnet: A two-stream network for fast and accurate 3d cloth draping. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 8739–8748 (2019)
24. Habermann, M., Xu, W., , Zollhoefer, M., Pons-Moll, G., Theobalt, C.: Deepcap: Monocular human performance capture using weak supervision. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE (jun 2020)
25. Hochreiter, S., Schmidhuber, J.: Long short-term memory. Neural computation **9**(8), 1735–1780 (1997)
26. Huang, Z., Xu, Y., Lassner, C., Li, H., Tung, T.: Arch: Animatable reconstruction of clothed humans. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3093–3102 (2020)
27. Jeong, M.H., Han, D.H., Ko, H.S.: Garment capture from a photograph. Computer Animation and Virtual Worlds **26**(3-4), 291–300 (2015)
28. Jiang, B., Zhang, J., Hong, Y., Luo, J., Liu, L., Bao, H.: Bcnet: Learning body and cloth shape from a single image. arXiv preprint arXiv:2004.00214 (2020)
29. Kanazawa, A., Black, M.J., Jacobs, D.W., Malik, J.: End-to-end recovery of human shape and pose. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 7122–7131 (2018)
30. Kazhdan, M., Hoppe, H.: Screened poisson surface reconstruction. ACM Transactions on Graphics (ToG) **32**(3), 1–13 (2013)
31. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
32. Klokov, R., Lempitsky, V.: Escape from cells: Deep kd-networks for the recognition of 3d point cloud models. In: The IEEE International Conference on Computer Vision (ICCV) (Oct 2017)
33. Kolotouros, N., Pavlakos, G., Black, M.J., Daniilidis, K.: Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2252–2261 (2019)

34. Lahner, Z., Cremers, D., Tung, T.: Deepwrinkles: Accurate and realistic clothing modeling. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 667–684 (2018)
35. Li, J., Chen, B.M., Hee Lee, G.: So-net: Self-organizing network for point cloud analysis. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
36. Li, Y., Bu, R., Sun, M., Wu, W., Di, X., Chen, B.: Pointcnn: Convolution on x-transformed points. In: Advances in neural information processing systems. pp. 820–830 (2018)
37. Liang, J., Lin, M.C.: Shape-aware human pose and shape reconstruction using multi-view images. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 4352–4362 (2019)
38. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: Smpl: A skinned multi-person linear model. ACM transactions on graphics (TOG) **34**(6), 1–16 (2015)
39. Ma, Q., Saito, S., Yang, J., Tang, S., Black, M.J.: Scale: Modeling clothed humans with a surface codec of articulated local elements. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16082–16093 (2021)
40. Mehta, D., Sridhar, S., Sotnychenko, O., Rhodin, H., Shafiei, M., Seidel, H.P., Xu, W., Casas, D., Theobalt, C.: Vnect: Real-time 3d human pose estimation with a single rgb camera. ACM Transactions on Graphics (TOG) **36**(4), 1–14 (2017)
41. Miguel, E., Bradley, D., Thomaszewski, B., Bickel, B., Matusik, W., Otaduy, M.A., Marschner, S.: Data-driven estimation of cloth simulation models. In: Computer Graphics Forum. vol. 31, pp. 519–528. Wiley Online Library (2012)
42. Miguel, E., Tamstorf, R., Bradley, D., Schvartzman, S.C., Thomaszewski, B., Bickel, B., Matusik, W., Marschner, S., Otaduy, M.A.: Modeling and estimation of internal friction in cloth. ACM Transactions on Graphics (TOG) **32**(6), 1–10 (2013)
43. Patel, C., Liao, Z., Pons-Moll, G.: Tailornet: Predicting clothing in 3d as a function of human pose, shape and garment style. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE (jun 2020)
44. Pumarola, A., Sanchez-Riera, J., Choi, G., Sanfeliu, A., Moreno-Noguer, F.: 3dpeople: Modeling the geometry of dressed humans. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2242–2251 (2019)
45. Qi, C.R., Liu, W., Wu, C., Su, H., Guibas, L.J.: Frustum pointnets for 3d object detection from rgb-d data. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
46. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 652–660 (2017)
47. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In: Advances in neural information processing systems. pp. 5099–5108 (2017)
48. Rasheed, A.H., Romero, V., Bertails-Descoubes, F., Wuhrer, S., Franco, J.S., Lazarus, A.: Learning to measure the static friction coefficient in cloth contact. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9912–9921 (2020)
49. Ravi, N., Reizenstein, J., Novotny, D., Gordon, T., Lo, W.Y., Johnson, J., Gkioxari, G.: Pytorch3d. https://github.com/facebookresearch/pytorch3d (2020)
50. Saito, S., Huang, Z., Natsume, R., Morishima, S., Kanazawa, A., Li, H.: Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2304–2314 (2019)
51. Saltelli, A.: Sensitivity analysis for importance assessment. Risk analysis **22**(3), 579–590 (2002)
52. Santesteban, I., Otaduy, M.A., Casas, D.: Learning-based animation of clothing for virtual try-on. In: Computer Graphics Forum. vol. 38, pp. 355–366. Wiley Online Library (2019)

53. Shen, Y., Liang, J., Lin, M.C.: Gan-based garment generation using sewing pattern images. In: Proceedings of the European Conference on Computer Vision (ECCV). vol. 1, p. 3 (2020)
54. Smith, D., Loper, M., Hu, X., Mavroidis, P., Romero, J.: Facsimile: Fast and accurate scans from an image in less than a second. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 5330–5339 (2019)
55. Tan, Q., Pan, Z., Gao, L., Manocha, D.: Realtime simulation of thin-shell deformable materials using cnn-based mesh embedding. IEEE Robotics and Automation Letters **5**(2), 2325–2332 (2020)
56. Thies, J., Zollhöfer, M., Nießner, M.: Deferred neural rendering: Image synthesis using neural textures. ACM Transactions on Graphics (TOG) **38**(4), 1–12 (2019)
57. Tiwari, G., Bhatnagar, B.L., Tung, T., Pons-Moll, G.: Sizer: A dataset and model for parsing 3d clothing and learning size sensitive 3d clothing. arXiv preprint arXiv:2007.11610 (2020)
58. Varol, G., Ceylan, D., Russell, B., Yang, J., Yumer, E., Laptev, I., Schmid, C.: Bodynet: Volumetric inference of 3d human body shapes. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 20–36 (2018)
59. Varol, G., Romero, J., Martin, X., Mahmood, N., Black, M.J., Laptev, I., Schmid, C.: Learning from synthetic humans. In: CVPR (2017)
60. Vidaurre, R., Casas, D., Garces, E., Lopez-Moreno, J.: Brdf estimation of complex materials with nested learning. In: 2019 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 1347–1356. IEEE (2019)
61. Wang, H., O'Brien, J.F., Ramamoorthi, R.: Data-driven elastic models for cloth: modeling and measurement. ACM transactions on graphics (TOG) **30**(4), 1–12 (2011)
62. Wang, T.Y., Ceylan, D., Popovic, J., Mitra, N.J.: Learning a shared shape space for multimodal garment design. arXiv preprint arXiv:1806.11335 (2018)
63. Wang, Y., Sun, Y., Liu, Z., Sarma, S.E., Bronstein, M.M., Solomon, J.M.: Dynamic graph cnn for learning on point clouds. ACM Transactions on Graphics (TOG) **38**(5), 1–12 (2019)
64. Wei, S.E., Ramakrishna, V., Kanade, T., Sheikh, Y.: Convolutional pose machines. In: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition. pp. 4724–4732 (2016)
65. Xu, Y., Zhu, S.C., Tung, T.: Denserac: Joint 3d pose and shape estimation by dense render-and-compare. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 7760–7770 (2019)
66. Yang, S., Ambert, T., Pan, Z., Wang, K., Yu, L., Berg, T., Lin, M.C.: Detailed garment recovery from a single-view image. arXiv preprint arXiv:1608.01250 (2016)
67. Yang, S., Liang, J., Lin, M.C.: Learning-based cloth material recovery from video. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 4383–4393 (2017)
68. Yang, S., Pan, Z., Amert, T., Wang, K., Yu, L., Berg, T., Lin, M.C.: Physics-inspired garment recovery from a single-view image. ACM Transactions on Graphics (TOG) **37**(5), 1–14 (2018)
69. Yu, T., Zheng, Z., Zhong, Y., Zhao, J., Dai, Q., Pons-Moll, G., Liu, Y.: Simulcap: Single-view human performance capture with cloth simulation. In: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5499–5509. IEEE (2019)
70. Zakharkin, I., Mazur, K., Grigorev, A., Lempitsky, V.: Point-based modeling of human clothing. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 14718–14727 (2021)
71. Zheng, Z., Yu, T., Wei, Y., Dai, Q., Liu, Y.: Deephuman: 3d human reconstruction from a single image. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 7739–7749 (2019)
72. Zhou, B., Chen, X., Fu, Q., Guo, K., Tan, P.: Garment modeling from a single image. In: Computer graphics forum. vol. 32, pp. 85–91. Wiley Online Library (2013)

73. Zhou, Y., Tuzel, O.: Voxelnet: End-to-end learning for point cloud based 3d object detection. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
74. Zhu, H., Cao, Y., Jin, H., Chen, W., Du, D., Wang, Z., Cui, S., Han, X.: Deep fashion3d: A dataset and benchmark for 3d garment reconstruction from single images. arXiv preprint arXiv:2003.12753 (2020)

**Ethics Statement:** Our training dataset is synthetically generated using unidentifiable human body textures drawn from the SURREAL dataset [59]. We used several real-world images and videos from unidentifiable models for qualitative testing. Our user study has been approved by the IRB. The results are only used to develop a statistical understanding of the perceptual performance of our method.

**Reproducibility:**  Detail on implementation is given below, with data and code to be released upon publication.

## Appendix A    Data Preparation and Training

In order to train our model, a large number of examples that contain ground truth human body parameters, garment meshes, and the corresponding material parameters are needed. These are very challenging to capture in real world. To supplement a very limited number of such real-world videos, we create a large dataset of videos generated with controlled variables with the corresponding ground-truth values for validation. We vary different conditions to generate this dataset:

**Human motion.** We sample common human motion sequences and shape parameters in the CMU Mocap dataset [17], including walking, sitting, boxing and climbing stairs.

**Garment meshes.** We use the garment dataset from [53].

**Material space.** We sample different materials uniformly in the discretized space mentioned in Sec. 3.

**Human and garment textures.** We use random human textures from SURREAL [59] and random garment textures from online images.

**Lighting.** We employ diverse outdoor and indoor environment maps downloaded from [1].

In total, we create a dataset of 250,000 images (10 motions * 100 garments * 250 frames) for single-frame human and garment geometry estimation, and 140,000 video sequences with each consisting of 25 images sampled at a frame rate of 5. Some examples of our training datasets are shown in the Appendices in the supplementary document (Fig. 9). We will make our code and datasets publicly available with this paper.

### A.1    Training Details

The split percentages for the train/validation/test sets are 85%/5%/10%. During training, we use Adam [31] to minimize the loss. The learning rate is $10^{-3}$ for the autoencoder, and $10^{-4}$ for others. They are scaled down to $10^{-4}$ and $10^{-5}$ respectively after 10 epochs. All hyper-parameters are chosen empirically. As shown in Fig. 2, we use a representative point set of size 256, and generate 128 points for each patch. The sampling size for the ground-truth point set (N in Fig. 2) is 16,384.

We trained and tested our model on a machine with 8 CPUs (Intel Xeon, 3.60GHz) for data loading and 2 GPUs for computing (GeForce GTX 1080). We trained our autoencoder, single-frame module, and the temporal module for 20 epochs, respectively.

## A.2 Recovery from Point Clouds to Garment Meshes

The point cloud representation of a garment mesh does not explicitly store the connectivity information, so it is necessary to apply certain prior when recovering meshes from point clouds. One straightforward way is to connect each point with its neighbors, determined by a distance threshold. However, this method does not guarantee the resulting mesh to be a manifold. Instead, we use manifold surfaces to approximate the results.

The overall pipeline to recover point clouds to meshes is split into several steps. We first employ Screen Poisson algorithm [30] on the point cloud with estimated normals using neighboring points to recover the overall shape and topology of the garment. Next, we remove the reconstructed vertices that are too far away from the point cloud, since the first step tends to generate water-tight meshes. We focus on removing large clusters that forms holes for neck, arms, and legs. After upsampling on the resulting mesh, we deform the mesh by minimizing the distance $d_M(\mathbf{P}, \mathbf{F})$ between the point cloud $\mathbf{P}$ and the mesh $\mathbf{F}$ as defined below:

$$d_M(\mathbf{P}, \mathbf{F}) = \frac{1}{|\mathbf{P}|} \sum_{\mathbf{p} \in \mathbf{P}} \min_{\mathbf{f} \in \mathbf{F}} \hat{d}(\mathbf{p}, \mathbf{f}) + \frac{1}{|\mathbf{F}|} \sum_{\mathbf{f} \in \mathbf{F}} \min_{\mathbf{p} \in \mathbf{P}} \hat{d}(\mathbf{p}, \mathbf{f}) \tag{12}$$

where $\hat{d}$ is the point-to-face distance, $\mathbf{f} \in \mathbf{F}$ is a face in the garment mesh, $\mathbf{F}$ is the set of all faces of the mesh. We use other regularization terms similar to the mesh deformation demo in Pytorch3D [49].



**Fig. 9. Training data examples**. Our dataset includes various garment topologies with rich body poses, textures, and background environments. Some examples are shown here.

## Appendix B     Ablation Studies

In the following ablation studies, we verify the effectiveness of our network model. We use the test errors to compare our model with the baselines, which include previous methods or their combinations. During the test, all other conditions are held the same, except for the network structure itself.

**Garment auto-encoder.** We use a two-level auto-encoder structure to address the highly dynamic geometry of garments, as discussed in Sec. 4. We demonstrate here its effectiveness compared to the baseline, which simply consists of a PointNet [46] encoder and an AtlasNet [19] decoder. There are two theoretical advantages of our method against the baseline. First, the two-level structure partitions the entire point cloud by patches that only overlap by a small fraction. This approach makes the local features more focused on its local shape rather than on a more expanded surface. Next, the two-level network also offers more capability to express detailed features and prevents overly-smooth reconstructions.

In Table 3(a), we report our test errors compared to the baseline. In addition to Chamfer Distance as we used in the training, we also use Sinkhorn Divergence [20], which is a fast approximator for computing Earth Moving Distance between two distributions. While Chamfer Distance indicates the average distance between two point clouds, Sinkhorn Divergence captures the density difference across space. The numbers in the table show that our method not only reproduces point clouds closer to the ground-truth but also has more evenly-distributed points due to our patch-wise partition during training. It also affects the performance of the material classification accuracy, which is our ultimate goal in this work. It is indicated in the table that a single-layer auto-encoder is not sufficient for the detailed geometry of dressed garments, resulting in much lower classification accuracy for material estimation. The results align with our theoretical analysis.

**Garment geometry estimation.** In Sec. 5.1, we proposed a closed-loop feedback structure to improve the estimation of both the human body and garment shape. We conducted an ablation study to show the difference introduced by this structure. Our baseline is a two-branch estimation block that predicts the human body and the garment in parallel but sharing the image features as input. Although it does not have the feedback correction, it already benefits from multi-tasking which can extract useful common image features. We use Mean Per Joint Position Error (MPJPE) for human estimation metric, and Chamfer Distance for garment estimation.

As shown in Table 3(b), by introducing a joint learning mechanism by conditioning the garment estimation with the body prediction followed by a corrective feedback loop, our model results in a much smaller error on human body estimation and garment estimation. More importantly, the accuracy of the body and garment estimation block has a large impact on the material estimation accuracy as well. As shown in the table, the material classification accuracy is greatly improved when using our proposed joint learning method for the body and garment estimation.

We further compare our method with a recent work [7] since their task is the most relevant to ours. We use the public dataset from Multi-Garment-Net [7], which consists of 95 scans of people. Nearly half of them are not suitable due to incorrect or incomplete garment labeling. We tested our model in the dataset without any fine-tuning and

used their reported numbers for comparison. We use the Chamfer Distance defined in Eqn. 5 as our test metric. As shown in Table 4, our model, without any fine-tuning or domain adaptation, achieves the smallest error – whether the ground-truth human body model is provided or not during reconstruction. Since our model is trained on a wider range of body poses and garment types than theirs while achieving better accuracy, it has shown to offer good generality to unseen inputs.
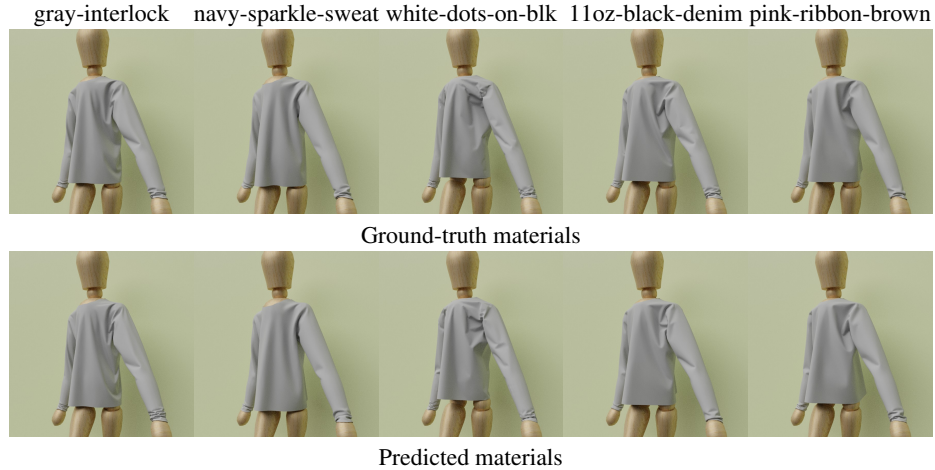
| gray-interlock | navy-sparkle-sweat | white-dots-on-blk | 11oz-black-denim | pink-ribbon-brown |



Ground-truth materials



Predicted materials

**Fig. 10. User study examples**. Our predictions received similarity scores from 5.7 to 8.5 in a 0-10 range, where 10 indicates exactly the same material and 0 completely different.

## Appendix C   Lab Experiments and User Study

In this experiment, we test the prediction accuracy of our method using real-world materials. We used five real-world materials measured from lab experiments [61], which are sampled from sweater, t-shirt, tablecloth, jeans, and blanket, respectively. The materials are used to create videos using the same pipeline as we generate the training data. The videos are then input to our network model for material estimation, which is used to generate the resulting videos with other conditions being held the same (see Fig. 10 for sample frames).

We first quantify our material recovery using the sensitivity analysis described in [67]. For each of the materials, we measure its stiffness ratio according to the deformations under a fixed amount of external forces. The measured values are then compared with the one predicted from our method, reported in Table 5. Our method achieves a relatively small error between 9.5% to 16.7%.

To further validate and quantify the material similarity, we conduct a user study to examine how close our estimation results are to the ground-truth data in human perception. In the study, we place two videos side by side for each material (example frames

shown in Fig. 10); then ask each participant to rate the level of material similarity: from 0 (totally different) to 10 (identical). There are 25 subjects in our study group: 17 male and 8 female, with age ranging from 20 to 40. To further calibrate the subjective score range, we use a pair of videos generated from the *same material* but with *different mesh resolutions* and inform the participants that this example has a similarity score of 5 for calibration. Our results show that the average similarity ratings for five tested materials vs. the ground-truth data are all larger than 5, ranging from 5.7 to 8.5, with an overall mean value of 7.1. These indicate that our method indeed can recover fabric materials with only minor perceptible differences to the real-world materials.

Besides the main results, we also conducted several other studies to investigate how people perceive the garment materials in different environmental conditions:

**Garment color and texture.** We changed the garment colors and textures; then asked participants the same questions. We found that by varying the garment colors, either brighter or darker, does not affect the similarity scores much – within a maximum difference of 0.4. On the other hand, changing the textures results in more perceptible effect – with similarity score differences of 0.5 to 0.9.

**Lighting.** We varied the lighting conditions in the rendered results to understand how shading on the garments affects material perception. The results show that the similarity scores are decreased by 1.1-1.3 when one of the videos has a different lighting angle than the other. These noticeable difference indicate the effects of lighting and shading on how humans perceive wrinkles and folds.

**Stiffness range.** We took the material called 'gray-interlock', consisting of 60% cotton and 40% polyester, and multiplied its material parameters by 1, 2, 5, and 10, respectively. The participants were asked to distinguish which of the two sampled materials is stiffer. Our findings indicate that there is little perceptible visual difference between lower stiffness values till when the garment stiffness is increased to a certain threshold. This finding also reconfirms our design rationale of the material space discretization using sensitivity analysis in Sec. 3.

## Appendix D   More Comparison with Previous Works

In Table 6, we extensively compare our work with previous ones regarding different assumptions, functionalities, and abilities. We define 'one model per garment' in 'generality' as that the method needs to create extra templates or registrations to the body, or need to retrain part of the network in order to predict a different garment type. Although DeepFashion3D [74] and ARCH [26] also have generality to different topologies to some extent, there are still limitations in their pipeline. The output from DeepFashion3D has to be continuous in one body part, meaning that they cannot support all topologies (*e.g.* dresses with holes). ARCH does support different garments on the body, but the output is a water-tight mesh together with the body, which is not always convenient for certain applications like virtual try-on. In contrast, our method naturally supports all kinds of topologies, and predicts the body and the garment in separate meshes.

**Additional Results.** We present more comparisons with other related learning-based methods for garment shape and pose recovery. Please note that none of these methods

(a) T-shirt and pants (day)     (b) T-shirt and skirt (day)

(c) T-shirt and pants (night)     (d) T-shirt and skirt (night)

**Fig. 11. Sample test images from [66]** for comparison in Table 2.

was designed to *recover the fabric materials*, the main focus and motivation of this work for simulation-based virtual try-on, as shown in Fig. 12.



(a) Input video clip (b) Estimated human body and garment (c) Simulation result

**Fig. 12. Virtual try-on example**. Our network model can clone a person's appearance from the physical world (in a video) to the virtual world, enabling simulation under different motions with accurate estimations of the body shape, garment geometry, and fabric materials.

We tested our model without any fine-tuning on the images provided in MGN [7] and DeepCap [24], as shown in Fig 13. Although our model is never exposed to real-world images, it successfully predicts the human body pose and garment shape correctly. Note that our model has no prior knowledge to any information about the garment shape, while MGN assumes that there is a one-to-one approximation mapping between body vertices and garment vertices. DeepCap has an initial ground-truth mesh
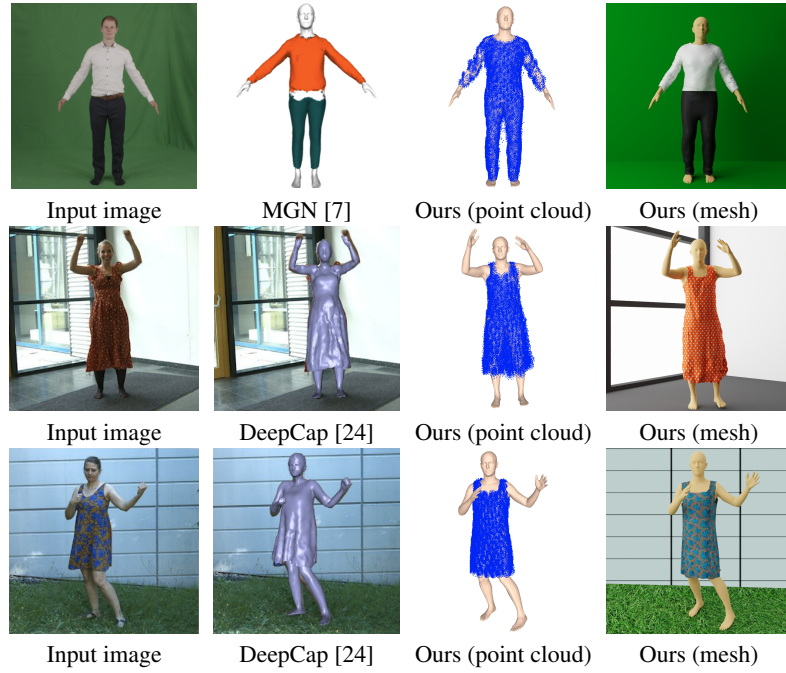
| Input image | MGN [7] | Ours (point cloud) | Ours (mesh) |
| Input image | DeepCap [24] | Ours (point cloud) | Ours (mesh) |
| Input image | DeepCap [24] | Ours (point cloud) | Ours (mesh) |

**Fig. 13. Qualitative Results**. Our model can achieve similar visual results with previous work without any knowledge of the target garment or any assumption of the topology (See Table 6).

that is to be optimized. Our model can thereby easily generalize to unseen garments, which is not possible using these two works.

# Appendix E    Application: Virtual Try-On

To further showcase the strength of our network, we apply it to a virtual try-on application. An online video clip showing a person wearing a dress is taken as input to our network (Fig. 12a). The body and the dress are estimated in each frame, and the fabric material is inferred using the garment motion and the image features. As shown in Fig. 12b, our method successfully infers the correct human body and the garment.

We then simulate the garment in a different body motion, which is the key functionality in virtual try-on systems. The simulated results (Fig. 12c) show that the garment motion provides similar visual impression with the input dress (mostly from the wrinkle motions of the dress). This example shows that our method can effectively extract the correct type of fabric material and transfer the given fabric material in a video to a simulation-based virtual try-on system. More animation results can be found in Fig. 8 and the supplementary video.
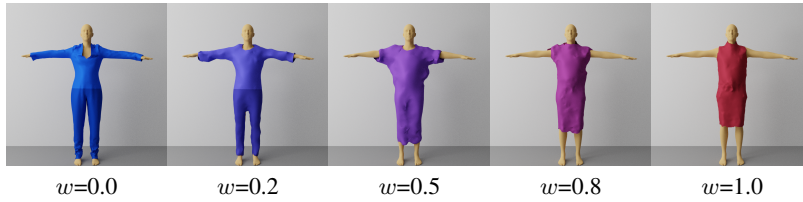
| $w=0.0$ | $w=0.2$ | $w=0.5$ | $w=0.8$ | $w=1.0$ |

**Fig. 14. Interpolation between different garments**. $w$ is the interpolation weight. We extract the latent code of 2 different garments: a long-sleeve top & pants (leftmost) and a short-sleeve dress (rightmost). We use the decoder to produce new garments of linearly-interpolated latent codes, enabling smooth transitioning between topologically different garments not achieved before.

## Appendix F    Latent Code Interpolation

To showcase the expressiveness of our learned latent space, we conducted an experiment generating new garments by interpolation. We encode two different garments to obtain their latent codes and linearly interpolate in between. We generate the new point clouds using the interpolated results accordingly. The visual results are shown in Fig. 14. Our results show that the interpolations represent a smooth transition between the two original garments, creating new garment styles that are not seen before during training. Note that modeling long dresses or garments with different mesh structures are not achieved in previous works, which either use displacement maps [7, 43] (thus not able to model dresses) or mesh-CNN for encoding local features [23] (thus not applicable to different mesh topologies). Our method is the first to propose *a feature space that unifies garments of different topologies with different body poses and shapes*, which is the key component to accurate garment material estimation, as demonstrated in Table 1.

We provide more interpolation results in Fig. 15. As described in Sec. F, our algorithm enables two garment meshes to be interpolated using their latent code to generate new garments. The first row shows an example where the dress is shorten at the bottom and extended at the sleeves gradually. We also show interpolation results in the point cloud form in Row 2-4. The second row is another example of transformation from dress to pants, and the last two examples demonstrate the ability to interpolate the garment with different body poses. In these two cases, the garment type is the same, but the body pose is changing. Although the poses in between are never seen in the training set, the interpolated garment point clouds follow the pose transition, showing very little interpenetrations with the body. Note that when the human legs are moving, the garment correctly deforms with the pose, close enough to be consistent with the body motion, while still being collision-free. See the supplementary video for the interpolation animation.
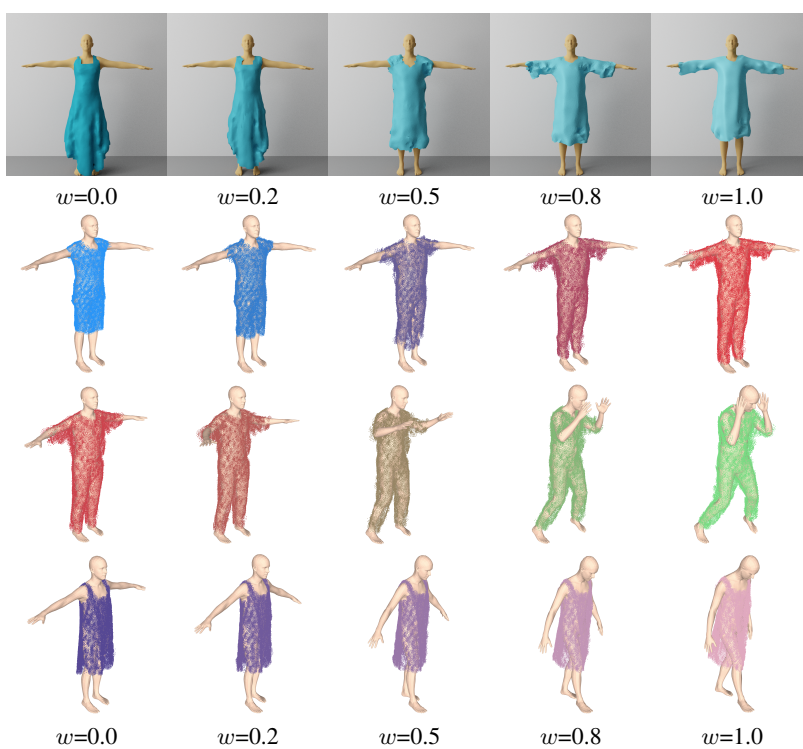
$w$=0.0          $w$=0.2          $w$=0.5          $w$=0.8          $w$=1.0

$w$=0.0          $w$=0.2          $w$=0.5          $w$=0.8          $w$=1.0

**Fig. 15. Interpolation results**. Our method can smoothly interpolate garments between different topologies and body poses.