

3D Body Shapes Estimation from Dressed-Human Silhouettes

Dan Song¹, Ruofeng Tong^{†1}, Jian Chang², Xiaosong Yang², Min Tang¹ and Jian Jun Zhang²

¹Zhejiang University, China
²Bournemouth University, UK

Abstract

Estimation of 3D body shapes from dressed-human photos is an important but challenging problem in virtual fitting. We propose a novel automatic framework to efficiently estimate 3D body shapes under clothes. We construct a database of 3D naked and dressed body pairs, based on which we learn how to predict 3D positions of body landmarks (which further constrain a parametric human body model) automatically according to dressed-human silhouettes. Critical vertices are selected on 3D registered human bodies as landmarks to represent body shapes, so as to avoid the time-consuming vertices correspondences finding process for parametric body reconstruction. Our method can estimate 3D body shapes from dressed-human silhouettes within 4 seconds, while the fastest method reported previously need 1 minute. In addition, our estimation error is within the size tolerance for clothing industry. We dress 6042 naked bodies with 3 sets of common clothes by physically based cloth simulation technique. To the best of our knowledge, We are the first to construct such a database containing 3D naked and dressed body pairs and our database may contribute to the areas of human body shapes estimation and cloth simulation.

Categories and Subject Descriptors (according to ACM CCS): I.3.m [Computer Graphics]: Miscellaneous—Image-based modeling

1. Introduction

Virtual fitting systems provide valuable visual information and size suggestions for users to buy clothes online, offering unique immersive experience. The main difficulties of realistic virtual fitting systems are the efficiency of physically based cloth simulation, the acquisition of 3D customized human model and the determination of physical property of cloth. Virtual fitting systems need to extract relatively accurate 3D body shapes of customers, especially in specific size aspects required for fashion design (e.g., chest/waist/hip size and body length).

3D human body estimation is a hot topic in computer graphics, as it plays an important role in applications such as movies, computer games and virtual fitting. A variety of methods have been proposed to estimate body shapes and they can be classified into non-parametric methods and parametric methods [CTT*].

Non-parametric methods take 3D points as input by scanning a human body in several views. Then a body mesh is acquired through registering or merging scan views and hole filling steps. This kind of method relies on the accuracy of scanners and restricts human to keep still with minimal clothes when scanning. For ordinary online shopping customers, they absolutely cannot bear these restrictions.

Parametric methods deform a 3D template body with a series of parameters. Input information is used as constraints to work out parameters and such information can be 3D or 2D which is more accessible. At early stage, parametric methods reconstructed 3D bodies using minimally-dressed information. Nowadays, more and more researchers utilize parametric human body model to estimate 3D body shapes from dressed-human information, the input of which is more convenient for users to get.

Clothes occlude human bodies and make body shapes estimation challenging. A body space learnt from a database of 3D naked bodies can be used to alleviate clothes effects [HSR*09, WPB*14, NH14]. Due to projecting dressed shapes to the body space, the estimated body is usually fatter than ground truth. Skin areas show more confidence for shapes so that some researchers set higher weight for the exposed-skin part of input information [B-B08, ZCD*15]. Consequently, this kind of method relies on skin areas and still leaves the covered body shapes ambiguous. Zhu et. al. [ZM15] allowed users to interactively estimate some body points under clothes, which depended on users' experience. Chen et. al. [CGZZ13] made the first attempt to model clothes deformation and proposed a parametric dressed body model. Compared with physically based cloth simulation, their clothes is less realistic.

Physically based cloth simulation has been researched for many years with explicit modules (e.g., cloth modeling, numerical time integration and collision handling). Given a 3D naked body and

[†] trf@zju.edu.cn

corresponding clothes, some commercial cloth simulation softwares can composite a realistic 3D dressed model and fossilization. However, it is extremely difficult and complex for the inverse operation, i.e., recovering the naked body from its dressed shape.

To analyze the relationship between a 3D naked body and its dressed shape, we construct a database containing 6042 pairs of 3D naked and dressed bodies for 3 clothes types. Because of the labor intensity of real shapes acquisition, we synthesize 6042 male bodies with a standing pose using 56 real male bodies from MPI database [HSS*09]. We dress the bodies with 3 sets of single-layer clothes using mature cloth simulation technique, which are a set of long-sleeved shirt and long pants (abbr. L&L), a set of short-sleeved shirt and long pants (abbr. S&L) and a set of short-sleeved shirt and short pants (abbr. S&S). Ideally, we aim to acquire same accuracy for "undressing estimation" as "dressing simulation" with the help of our database. To our knowledge, this is the first database that contains 6042 pairs of 3D naked and dressed human bodies, which may benefit the areas of cloth simulation and human body shapes estimation.

Based on our database, we create training samples containing dressed-human silhouettes, initial 3D landmarks and target 3D landmarks. An effective feature descriptor is proposed to combine 3D naked body landmarks with dressed-human silhouettes, and regressors are trained for guiding landmarks movements (from initial landmarks to target landmarks) according to dressed-human silhouettes with training samples. In testing phase, given dressed-human silhouettes and a set of initial landmarks as input, we regress target 3D body landmarks with training results as guidance. The regressed landmarks are used to constrain SCAPE model [ASK*05] for body reconstruction.

Our work provide a tangible solution for ordinary users to access their 3D body data with common clothes holding a standard standing pose. The main application of our work is virtual fitting. Users can reconstruct their own 3D bodies by our method with photos so that they could receive size recommendations and visualize their 3D views with new clothes in virtual fitting room. There are many other potential applications. Take computer games for example, it is exciting to create a virtual customized character with similar shape of the player in real world. It is also possible to combine our method with 3D printing for customized human toys and sculptures.

Our main contributions are summarized as: (a) an automatic framework for efficient body shapes estimation from dressed-human silhouettes, (b) a database containing 6042 pairs of 3D naked and dressed bodies for 3 clothes types and (c) an effective feature descriptor combining 3D naked body landmarks with 2D dressed-human silhouettes.

2. Related Work

2.1. 3D Human Body Reconstruction Methods

3D human body reconstruction methods can be classified to non-parametric methods and parametric methods [CTT*]. A variety of works [TCL*13] [LVG*13] [TZL*12] use non-parametric methods to obtain a 3D mesh which is close to scanned points cloud. If we use non-parametric methods to estimate body shapes, we rely on

3D scanners and restrict human to keep still with minimal clothes when scanning, which is far away from our input requirement of using dressed-human photos.

Many parametric methods have been proposed for 3D human body reconstruction. Allen et. al. [ACP03] came up with a statistical model to learn a shape space for a similar pose. Similar ideas were used for human body shapes estimation from one or more images [SYW06] [CC09] [BSWX13] and body measurements [WS13]. To allow pose variation, Anguelov et. al. [ASK*05] proposed SCAPE model which considered body deformation as the combination of pose deformation and shape deformation. SCAPE model successfully models human body variations and attracts lots of researchers.

Some researchers used SCAPE model to represent human bodies for animating realistic clothing [GRWB12]. Some researchers used it to estimate body shapes from a single image or painting of minimally-dressed people [GWBB09]. Balan et. al. [BB08] adopted it to estimate body shapes with 4 images of normally-dressed people from different views. Such estimated mesh can be utilized to modify the input image [ZFL*10] or video [JTST10]. Weiss et. al. [WHB11] adopted it to estimate human body with noisy Kinect data. We use SCAPE model for our body reconstruction with several 3D landmarks.

2.2. 3D Body Shapes Estimation Under Clothes

Hasler et. al. [HSR*09], Wuhler et. al. [WPB*14] and Neophytou et. al. [NH14] regarded clothes as noises. They trained their models with databases of minimally-dressed bodies to learn human body spaces which did not contain clothes. They used dressed-human information (3D dense points) as constraints to deform a template mesh and got a coarse body mesh which was affected by clothes. Then the coarse mesh was represented in the learnt body space to alleviate clothes effects. These methods work better for tight clothes and the reconstructed bodies are usually fatter than ground truth. With RGBD data acquired by Kinect, Zeng et. al. [ZCD*15] used the RGB image to detect skin areas as tight constraints for body shapes estimation. These methods should find correspondences between input 3D dense points and target mesh vertices, which is time-consuming.

Compared with 3D information obtained by scanners, images are more accessible. Balan et. al. [BB08] took dressed-human silhouettes from 4 views as input. They detected exposed-skin parts to decide weights for input information and then constrained SCAPE model with input information. The reconstruction energy was represented as pixels differences between input silhouettes and projections of target mesh in 4 views. This representation and computation are complex so that they used a gradient-free direct search simplex method to optimize the energy. Fitting takes approximately 40 minutes for a single model.

Chen et. al. [CGZZ13] made the first attempt to consider clothes and they extended SCAPE model to a dressed human shape model. They used deformation transfer [SP04] technique to construct a database of naked and dressed body pairs. For each type of clothes, only one naked and dressed body pair in database was generated

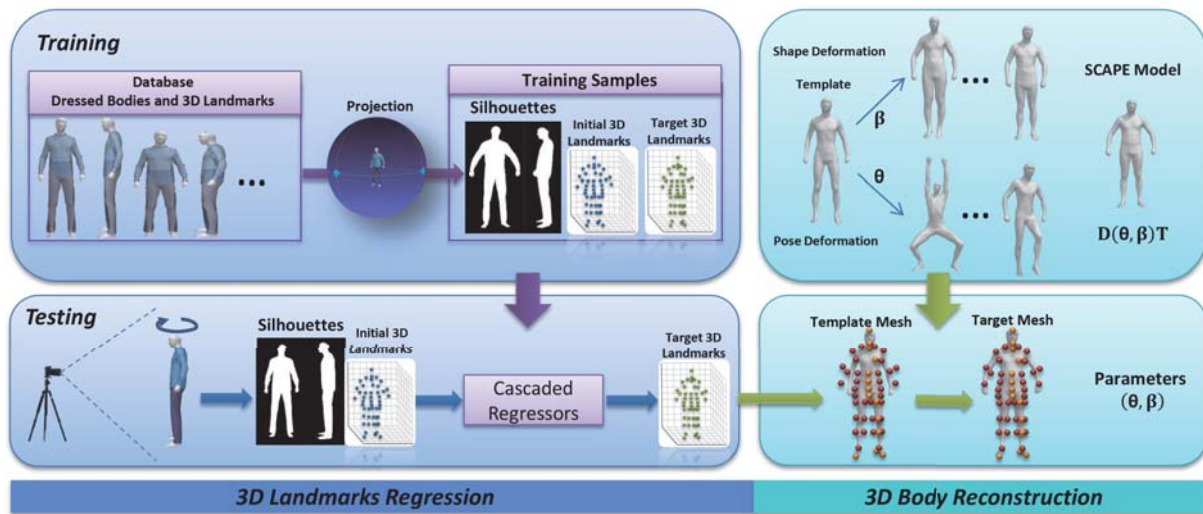


Figure 1: Overview. Our novel automatic framework experiences two stages: 3D landmarks regression and 3D body reconstruction. We create training samples using our database which contains 3D naked and dressed body pairs. Regressors are trained to learn the predictive relationship between dressed-human silhouettes and target 3D naked body landmarks. The regressed landmarks are used as constraints to optimize SCAPE model which is trained with a database of 3D naked bodies.

by an animation software (POSER) with high quality while the other pairs were obtained by deformation transfer. They assumed that clothes deformation was only related to body shape deformation for a specific clothes type and learnt clothes-related coefficients for the model additionally. This also leads to the non-linear optimization problem which costs some time. They made users to manually specify initial 2D joints and then found the correspondences between 3D dressed body vertices and 2D dressed body contour points with a HMM method [KSvdP09], the computation of which is not fast.

Zhu et. al. [ZM15] predicted body shapes under clothes using orthogonal-view dressed-human photos. They allowed users to interactively estimate some body points for photos and then searched naked body silhouettes in a database according to those points. Unlike previous work, they did not directly build the energy between 2D silhouettes and 3D bodies. To acquire efficiency, they defined 2D features for silhouettes and 3D features for bodies and learnt the relationship between them using a database of 3D naked bodies and corresponding naked body silhouettes. Finally, they optimized body mesh using 3D features.

3. Approach Overview

There are two obstacles for "undressing estimation" and we use 3D naked body landmarks to overcome them. One obstacle that hinders the effectiveness of previous methods is the way to remove clothes. We construct a database of naked and dressed body pairs and propose a data-driven method to predict 3D naked body landmarks from dressed-human silhouettes to solve this obstacle. Previous work should find correspondences between the vertices of target mesh and points from input information for body reconstruction, which is another obstacle blocking the efficiency. In our work,

the landmarks indices of target mesh are pre-defined and we use predicted naked body landmarks to constrain target mesh, avoiding time-consuming correspondences mapping process. For clear exposition, in the following contents we use clothes type 1 (L&L) as an example to show our method when there is no relevance with clothes types.

As figure 1 shows, our automatic framework experiences two stages: 3D landmarks regression and 3D body reconstruction. To analyze the relationship between naked body and its dressed shape, we construct a database of naked and dressed body pairs. With the help of our database, we create training samples consisting of dressed-human silhouettes, initial 3D landmarks and target 3D landmarks. The target landmarks for each body are pre-defined (Figure 3) and we obtain initial landmarks by randomly using other sample's target landmarks. Database construction and training samples preparation are introduced in section 5. We propose an effective feature descriptor to combine 3D naked body landmarks with dressed-human silhouettes, and train regressors for guiding landmarks movements (from initial landmarks to target landmarks) according to dressed-human silhouettes with training samples. Our regression framework is explained in section 6.1 and feature descriptor is illustrated in section 6.2. In the testing phase, with training results as guidance and a set of initial landmarks, we regress target 3D naked body landmarks. The regressed landmarks are used to constrain SCAPE model [ASK*05] for our body reconstruction, which is introduced in section 4.

4. Parametric Body Reconstruction

We adopt SCAPE model [ASK*05] for our parametric body reconstruction. SCAPE model decouples human body deformation into pose deformation and shape deformation which are separate-

ly controlled by pose parameter θ and shape parameter β . Given parameters θ and β , the vertices positions $\mathbf{Y} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_V\}$ of target mesh are solved by minimizing the least square error:

$$\arg \min_{\mathbf{y}_1, \dots, \mathbf{y}_V} \sum_{k=1}^K \sum_{d=2}^3 \|\mathbf{R}_{p[k]}(\theta) \mathbf{S}_k(\beta) \mathbf{Q}_k(\theta) \hat{\mathbf{v}}_{d,k} - (\mathbf{y}_{d,k} - \mathbf{y}_{1,k})\|^2 \quad (1)$$

where K denotes the total number of triangles and V is the total number of vertices. $\mathbf{y}_{1,k}$, $\mathbf{y}_{2,k}$ and $\mathbf{y}_{3,k}$ are three vertices of a triangle k . $\hat{\mathbf{v}}_{d,k}$ is an edge of template mesh, and $\mathbf{y}_{d,k} - \mathbf{y}_{1,k}$ represents the corresponding edge of target mesh. The human body is divided into 17 partitions and $p[k]$ means the partition p that triangle k locates at. $\mathbf{R}_{p[k]}(\theta)$ is a 3×3 matrix with a 3-dimensional parameter θ , which represents the rigid rotation for partition p . $\mathbf{Q}_k(\theta)$ is a 3×3 matrix that shows the non-rigid deformation (e.g., muscle bulging) induced by pose variation. $\mathbf{S}_k(\beta)$ is a 3×3 matrix explaining the shape variation between different individuals.

The formulation of $\mathbf{R}_{p[k]}$, \mathbf{S}_k and \mathbf{Q}_k can be found in [ASK*05]. We train SCAPE model using MPI database [HSS*09] which consists of pose database and shape database. Pose database contains one individual with 35 different poses and it is used to train the relationship between \mathbf{Q}_k and θ . Shape database contains 56 individuals with one standard standing pose and we utilize it to train the relationship between \mathbf{S}_k and β .

According to formula (1), Cheng et al. [CTT*] derived the linear representation of \mathbf{Y} with reference to θ and β respectively. When fixing shape parameter β and rotation deformation $\mathbf{R}_{p[k]}$, \mathbf{Y} is represented as equation (2). \mathbf{c} and \mathbf{d} are determined by shape parameter β and rigid rotation $\mathbf{R}_{p[k]}$. Similarly, \mathbf{Y} is represented as equation (3) when fixing pose parameter θ . \mathbf{f} and \mathbf{g} are decided by pose parameter θ . Both the detailed explanation of \mathbf{c} , \mathbf{d} , \mathbf{f} and \mathbf{g} and the derivation of equation (2) and (3) can be found in [CTT*].

$$\mathbf{Y} = \mathbf{c} \cdot \theta + \mathbf{d} \quad (2)$$

$$\mathbf{Y} = \mathbf{f} \cdot \beta + \mathbf{g} \quad (3)$$

As we address in section 3, we use body landmarks to constrain SCAPE model for our body reconstruction to leave out the time-consuming correspondences mapping process. We compute parameters θ and β by alternately optimizing two energies:

$$E(\theta) = \sum_{l=1}^L \|\mathbf{y}_l(\theta) - \mathbf{P}_l\|^2 + w_\theta \sum_{p_1, p_2} \|\theta_{p_1} - \theta_{p_2}\|^2 \quad (4)$$

$$E(\beta) = \sum_{l=1}^L \|\mathbf{y}_l(\beta) - \mathbf{P}_l\|^2 + w_\beta \left(\frac{1}{2} \beta^T \Lambda \beta \right) \quad (5)$$

L is the number of body landmarks. \mathbf{P}_l denotes the position of body landmark l and \mathbf{y}_l is the corresponding vertex of target mesh. p_1 and p_2 are two adjacent partitions of body. The second term of equation (4) is a quadratic smoothness term to keep the adjacent body parts changing continuously. In equation (5), $\Lambda = \text{diag}(1/\sigma_1^2, 1/\sigma_2^2, \dots, 1/\sigma_B^2)$. σ_i^2 is an eigenvalue from SCAPE shape parameters and B is the dimension of β . The second term in equation (5) is to regularize β . w_θ and w_β are weight coefficients

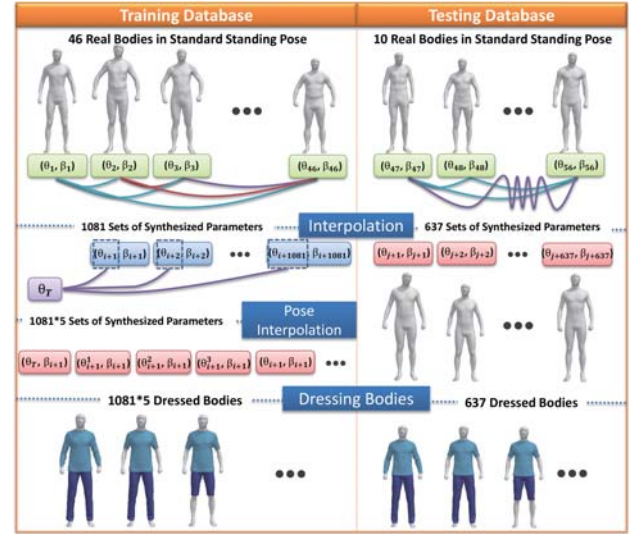


Figure 2: Database construction. Our database consists of 6042 (1081 × 5 + 637) synthesized naked bodies and corresponding dressed bodies with 3 sets of clothes in suitable size. 56 male bodies with standard standing pose in MPI database are represented as parameters of SCAPE model. These parameters are used for interpolation to generate more sets of parameters which determine our synthesized bodies. We dress the synthesized bodies using a physically based cloth simulation software with 3 common clothes types (L&L, S&L and S&S).

and we set $w_\theta = 0.143$ and $w_\beta = 0.002$ for our implementation. The value of w_θ , w_β and the number of iterations are validated in Appendix II.

5. Database Construction and Training Samples Preparation

Because of the difficulty in transformation from dressed body model to naked body model, we use a data-driven approach to indicate naked body information from dressed-human silhouettes. To the best of our knowledge, this is the first database containing thousands of 3D naked and dressed body pairs. Due to the massive manual work of dressing simulation, we currently only use male bodies with standard standing pose in MPI database [HSS*09] for our database construction. We use our database to validate our automatic body shapes estimation framework. The framework can be applied to female situation directly when provided with a database of female naked and dressed body pairs.

Figure 2 illustrates the process of our database construction. Parameters (θ, β) for each body are acquired by representing the mesh with SCAPE model, and one 3D body mesh is determined by one set of parameters. We interpolate one set of parameters from two known sets of body parameters with a random weight between 0 and 1. 1081 sets of parameters are synthesized using 46 sets of parameters for training database. Similarly, 637 synthesized sets of parameters are obtained using 10 sets of parameters to construct testing database. For each one of 1081 body shapes in training database, we assign 5 poses for him by inter-

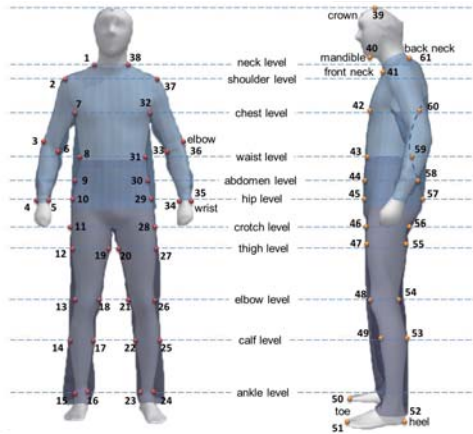


Figure 3: Landmarks. Landmarks on naked body under clothes. Red vertices are used for front view while yellow ones are used for side view.

polating 5 pose parameters between a template pose and its current pose. Therefore, the distribution of human body poses in our training database are around a template pose while preserving small pose variations. However, we do not interpolate poses for testing database to avoid artificial pose overlapping with training database. Finally, we design a T-shirt, a long-sleeved shirt, shorts and long pants using physically based cloth simulation software (Marvelous Designer: www.marvelousdesigner.com) and dress 6042 ($1081 \times 5 + 637$) bodies with three sets of clothes (L&L, S&L and S&S) in suitable size.

Based on the observation of landmarks selection in previous work [SI03, ZMK13], we set several critical landmarks (figure 3) to better represent body shapes, especially for virtual fitting. These landmarks are used for parametric body reconstruction and to some extent, they have intrinsic relationships with dressed-human silhouettes. We select 61 body landmarks on a template mesh and transfer them to all the meshes in our database which have same number of vertices and same topology as the template mesh, avoiding manual landmarks annotation for every body mesh.

The constructed database is used to prepare training samples for landmarks regression. As shown in figure 1, a training sample contains initial 3D landmarks, target 3D landmarks and dressed-human silhouettes. We project dressed bodies in front and side views to prepare dressed-human silhouettes. For each body, we have its target 3D landmarks. We obtain initial 3D landmarks by randomly using other samples' target landmarks. We set a vertex (nearly located at the center of body) as an anchor point to locate human body. For two bodies used for extracting initial landmarks and target landmarks, the difference of their anchor points' positions in $x/y/z$ direction is a random value between $-5cm$ and $5cm$. To enhance training samples, we assign E sets of initial positions to each of 5405 (1081×5) target positions in our database. Thus, the number of training samples N equals to $5405 \times E$. Here we set $E = 30$, and effects caused by different enhance extents of training samples are shown in Appendix I.A.

6. 3D Landmarks Regression

3D landmarks regression process plays an important role in our efficient automatic framework. For one hand, 3D landmarks are used to bridge naked body shapes and dressed-human silhouettes, which is a novel idea for "undressing estimation". For another hand, 3D landmarks leave out time-consuming correspondences mapping for body reconstruction, guaranteeing the efficiency for our work.

We adopt a boosting tree regression method which has been proven effective in face alignment [DWP10, CWWS14, CWLZ13] and body landmarks detection [CTT*]. Dollar et. al. [DWP10] and Cao X. et. al. [CWWS14] used this method to regress 2D facial landmarks from a facial image. Cao C. et. al. [CWLZ13] adopted it to predict 3D facial landmarks. Cheng et. al. [CTT*] regressed 2D body landmarks from a depth image (of people with minimal clothes) using this method. We aim to regress 3D body landmarks from dressed-human silhouettes.

6.1. Regression Framework

As formula (6) shows, we try to learn a regression function R using training samples, in order to get target landmarks positions according to a set of initial landmarks and dressed human silhouettes.

$$\arg \min_R \sum_{j=1}^N \| \mathbf{P}_T^j - (\mathbf{P}_I^j + R(\mathbf{I}^j, \mathbf{P}_I^j)) \|^2 \quad (6)$$

N is the number of training samples and \mathbf{I}^j represents dressed-human silhouettes. \mathbf{P}_T^j denotes the target landmarks positions and \mathbf{P}_I^j is the initial landmarks positions. However, it is extremely complicated to learn such a function. Fortunately, boosting tree regression method combines a series of weak regressors to achieve this goal. Although a simple weak regressor may be not accurate enough, the cascaded regressors become powerful [CWWS14, CTT*].

As equation (7) shows, boosting tree method regress 3D body landmarks in an incremental manner. \mathbf{P}_i represents positions of landmarks in the i^{th} stage, \mathbf{I} denotes dressed-human silhouettes and R_i is the regressor in the i^{th} stage. The goal of regressor R_i is illustrated in formula (8), which is similar to (6).

$$\mathbf{P}_i = \mathbf{P}_{i-1} + R_i(\mathbf{I}, \mathbf{P}_{i-1}) \quad i = 1, 2, \dots, m \quad (7)$$

$$\arg \min_{R_i} \sum_{j=1}^N \| \mathbf{P}_T^j - (\mathbf{P}_{i-1}^j + R_i(\mathbf{I}^j, \mathbf{P}_{i-1}^j)) \|^2 \quad (8)$$

We aim to represent regressor R_i as a piecewise function. Suppose training sample j is classified into class Ω_c in i^{th} stage, its R_i is computed as equation (9). $\delta \mathbf{P}_i^j$ equals to $\mathbf{P}_T^j - \mathbf{P}_{i-1}^j$ and $|\Omega_c|$ represents the total number of training samples in class Ω_c . We should notice that the landmarks movements of samples in a subset can be approximated as the average movements if their target landmarks movements are similar to each other [CTT*]. Therefore, the problem converts to how to classify training samples in each stage, which is explained in the next section.

$$R_i = \frac{\sum_{j \in \Omega_c} \delta \mathbf{P}_i^j}{|\Omega_c|} \quad (9)$$

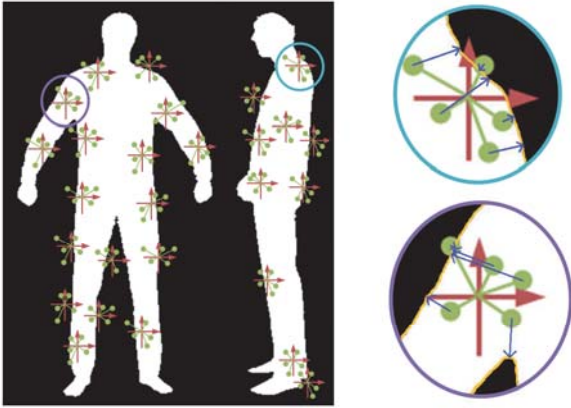


Figure 4: Sampling points and feature descriptor. 2D landmarks are marked as the centers of red crosses and green points represent sampling points. For clear display, this figure does not show all points. The feature descriptor is defined as displacements from sampling points to their nearest contour points (blue arrows).

6.2. Feature Descriptor and Classification

Since both dressed-human silhouettes and target landmarks have intrinsic relationship with ground truth body shapes, we assume that similar target landmarks movements (from current landmarks to target landmarks) mean similar relationship between current landmarks and dressed-human silhouettes. We propose a novel feature descriptor to describe the relative relationship between 3D landmarks and 2D dressed-human silhouettes. Previous work [DWP10] [CWLZ13] [CWWS14] [RCWS14] [CTT*] used image intensity to define feature descriptor while image intensity is so weak for our black-and-white silhouettes.

Figure 4 shows our feature descriptor. Firstly, we project 3D landmarks to get 2D image points. For simple illustration, we call those points 2D landmarks in the following content. The camera configuration for projection is the same as that used in projecting 3D dressed bodies to get silhouettes. When we use real photos for testing, we estimate the camera configuration, which will be illustrated in section 7.2. Secondly, G sampling points are acquired by sampling around 2D landmarks with a Gaussian distribution whose mean value is 20 pixels and standard deviation value is 2 pixels for a 800×600 image. Finally, we define our feature descriptor as displacements from sampling points to their nearest dressed body contour points. In our implementation, we set 8 sampling points around a 2D landmark. After processing every 300 regressors, we generate another Gaussian sampling. All the sampling configurations are recorded to guide testing phase.

The displacement contains values in both x and y direction, so each training sample has a $2G$ -dimensional feature descriptor. Using the $2G$ -dimensional vector for classification is a straightforward but awkward idea for our regression task. Instead, we adopt a random fern algorithm [OCLF10] to acquire feature reduction and select F out of $2G$ to constitute a F -bit label as feature for classification. The values of F bits of its $2G$ -dimensional feature descriptor are compared with corresponding preset thresholds. If the value is

less than its threshold, the corresponding fern is set to 0. Otherwise, it is set to 1. Consequently, each training sample gets a F -bit binary label and all samples are classified into 2^F classifications in each stage. F is set to 4 in our implementation, and the value of F is discussed in our Appendix I.B.

We propose a variance-ranked method to select F bits. The i th bit for N training samples forms a set D^i ($i = 1, 2, \dots, 2G$) which contains N elements. Then we compute the variance of each set and choose the top F maximal ones. We reduce the variance after using it to make full use of all bits. Comparison with other feature selection method can be found in our Appendix I.D.

6.3. Regression Configurations and Training Results

We stop cascading regressors until the decrease of average landmarks error between two adjacent regressors is less than a pre-set threshold (i.e. $0.0005mm$), and the total number of regressors is recorded. For each regressor, we also record the Gaussian sampling configuration, selected F bits and their corresponding thresholds. Most importantly, the landmarks movements for each classification are recorded. These configurations and training results are used to guide 3D landmarks regression for testing phase.

6.4. Testing Phase

For each regressor in the testing phase, feature is computed and the testing sample is classified into one classification. Then landmarks positions are updated using the landmarks movements for that classification. After processing the same number of regressors (as recorded in training phase), we acquire target 3D landmarks.

7. Experiment Results

Our program is run on a 64-bit desktop machine with 4.0GHz Intel(R) Core(TM) i7-4790K CPU and 32GB RAM. With single-core and single-thread programming, 3D landmarks regression only takes average of 0.028 second using two silhouettes with resolution 800×600 . 3D body reconstruction costs 3.583 seconds on average.

To demonstrate the accuracy of our method, we show the statistical landmarks regression error and body reconstruction error when high-quality silhouettes are provided for three clothes types (section 7.1). To demonstrate the practicability of our method, we use real human photos as input to estimate 3D body shapes, the silhouettes of which may contain noises (section 7.2). Section 7.3 highlights our advantages when compared with previous "undressing estimation" works.

7.1. Testing Database Trials

Our testing database consists of 637 naked and dressed body pairs, which is not overlapped with training database. We project those dressed bodies in front and side views to get silhouettes that are free of noises. In our implementation, the difference of anchor points' positions in x/y/z direction of the two bodies for generating silhouettes and initial landmarks is a random value between $-5cm$ and

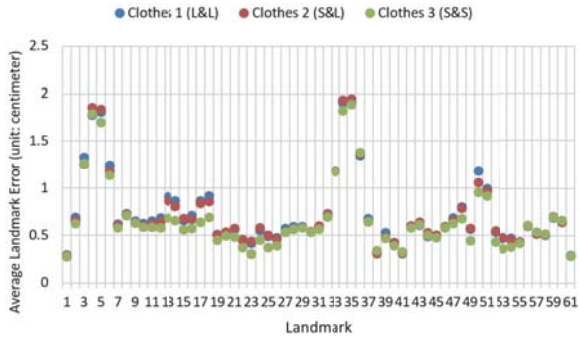


Figure 5: Average error for 61 landmarks. The location of each landmark is illustrated in figure 3. Each point shows the mean value of landmark error in Euclidean distance for 637 samples.

5cm. Therefore, we do not need to match 2D silhouettes and initial landmarks.

We evaluate the regression method through computing 3D landmarks error. For each landmark, we compute the average error of 637 testing samples separately for 3 clothes types and show the error for 61 landmarks in figure 5. The error for landmarks located at legs (indexed 12-27, 47-55) decreases for clothes type 3 (S&S) condition because of more convincing skin-exposed information. Those 8 landmarks with large error locate at arms (indexed 3-6, 33-36), where pose ambiguity exists when we only have front-view and side-view silhouettes. However, we aim to estimate 3D body shapes more than accurate poses so that we tolerate the large error for those points. Excluding those 8 landmarks, the average error for other 53 landmarks in Euclidean distance are shown in table 1.

Table 1: Average landmarks error

Clothes 1 (L&L)	Clothes 2 (S&L)	Clothes 3 (S&S)
0.599 cm	0.590 cm	0.539 cm

For 637 testing samples, let the registered estimated bodies and ground truth bodies own the same pose parameter and we compute vertex error to evaluate shape divergence. The average error of all vertices is illustrated in table 2 and figure 6 shows the average error for each vertex. We further evaluate shape estimation in body measurements (height, chest size, waist size and hip size) and the cumulative error distribution is displayed in figure 7. The curves for 3 clothes types in each body measurement are very close to each other. About 90% of testing samples' height/hip size/chest size/waist size error is less than 1cm/1.2cm/2.2cm/3cm, which satisfies the size tolerance for ordinary clothes [ZM15].

Table 2: Average vertex error

Clothes 1 (L&L)	Clothes 2 (S&L)	Clothes 3 (S&S)
0.515 cm	0.512 cm	0.502 cm

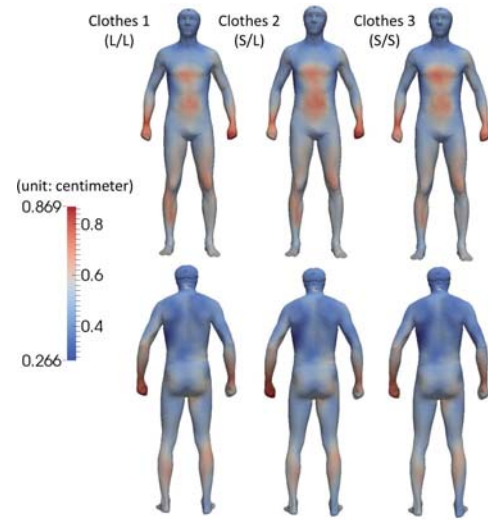


Figure 6: Average vertex error. The colors of points show the average vertex error of 637 samples.

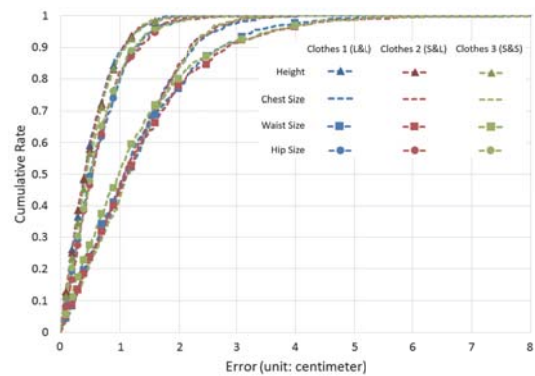


Figure 7: Cumulative error distribution in body measurements. Cumulative error distribution in height, chest size, waist size and hip size aspects.

7.2. Real Human Photos Testing

We take photos of 9 persons wearing casual clothes (similar to the ones in our database) holding the standard standing pose with only 1 camera at a known fixed distance, and figure 8 shows our environment for taking photos. Dressed-human silhouettes are obtained with the help of Photoshop. After being captured the front view, the user turns to his right side holding the same pose. The pose clues provided by orthogonal-view silhouettes do not contradict with each other with small pose differences. We use camera calibration toolbox for Matlab (http://www.vision.caltech.edu/bouguetj/calib_doc/) to estimate camera intrinsic parameters and extrinsic parameters. We can compute the 3D position of anchor point for initial landmarks, when we have its 2D position in image, its distance to camera and camera parameters. Thus, we do not need to match the initial landmarks and input silhouettes.

Table 3 shows the mean error between estimated bodies in body



Figure 8: Testing photos. Examples of testing photos in 2 views for 3 clothes types.

sizes aspects and our measured sizes for 9 persons. At the same time, we acquired the range data of those persons with minimal clothes using Kinect Fusion and reconstruct 3D human bodies with Cheng's method [CTT*]. Figure 9 visually compares our reconstructed bodies for 3 clothes types with theirs.

Table 3: Average body measurements error

CT	Error	Height	Chest Size	Waist Size	Hip Size
1 (L&L)		1.69 cm	1.87 cm	2.01 cm	1.90 cm
2 (S&L)		1.72 cm	1.80 cm	1.85 cm	1.87 cm
3 (S&S)		1.83 cm	1.83 cm	1.85 cm	1.79 cm

7.3. Comparison with Previous Work

We use dressed-human silhouettes as input, which is more convenient for users. Some works [HSR*09, WPB*14, NH14, ZCD*15] taking 3D scanning points as input required 3D scanners. Zhu et. al. [ZM15] made users to interactively estimate some body points under clothes.

The details of 3D body reconstruction of previous work are explained in section 2.2. We do not run their programs using the same hardware. Balan et. al. [BB08] used a gradient-free direct search simplex method for optimization which costs 40 minutes with a 2GHz CPU. Hasler et. al. [HSR*09] spent 11.5 minutes on optimization. Chen et. al. [CGZZ13] spent 1 minute with a 3GHz CPU and 2GB RAM. We adopt an efficient regression method to acquire 3D landmarks positions within 0.03 second. We avoid time-consuming vertices correspondence finding process because those landmarks are preset on registered bodies. Our body shape optimization target is minimizing the Euclidean distance between regressed 3D landmarks and corresponding vertices of a template mesh and the process takes less than 4 seconds.

Only a few related works provided statistical error for real human photos testing. Since the input requirements for their works and ours are different, we cannot compare these methods using the same testing data set. Balan et. al. [BB08] tested 6 persons wearing 6~10 kinds of ordinary clothes with 11 poses. Wuhrer et. al. [WPB*14] had 18 scans of 5 persons in casual office clothes with up to 5 poses. For average height/chest size/waist sizes error, Balan et. al. achieved about 1.03/4.65/4.73 centimeters while

Wuhrer et. al. got 2.52/14/15.8 centimeters. The input requirement for our method is illustrated in the first paragraph of section 7.2 and the statistical error is shown in table 3.

8. Conclusions, Limitations and Future Work

We propose a novel automatic framework to efficiently estimate 3D body shapes from dressed-human silhouettes. We build a database containing 6042 pairs of 3D naked and dressed bodies for 3 clothes types (L&L, S&L and S&S), which may benefit the areas of human body estimation and cloth simulation. Critical vertices are selected as landmarks to represent body shapes, which leaves out the time-consuming vertices correspondences finding process for body reconstruction and guarantees efficiency of our method. We explore a novel landmark-indexed feature descriptor to combine 3D body landmarks with 2D dressed-human silhouettes. Based on our constructed database, we learn a regression function to predict 3D landmarks positions according to 2D dressed-human silhouettes with our effective feature. 3D bodies are acquired by constraining SCAPE model with regressed landmarks.

Experiments show that our approach achieves good reconstruction results in body measurements, satisfying the size tolerance of clothing industry. 3D body shapes are estimated within 4 seconds automatically while the fastest method reported previously need 1 minute. We also validate key implementation configurations for our method and show the robustness of our feature descriptor in the appendix. Our work makes it more convenient for ordinary users to access their 3D body shapes with some kinds of common clothes, which will accelerate the evolution of virtual fitting industry in the future.

In the following contents, we show the limitations and future work of our work. We use silhouettes as input of our core algorithm, because clothes texture is of less use for "undressing estimation". Our method is affected by the quality of silhouettes obtained from photos, so we may explore excellent image segmentation technique in the future.

Our database depends on cloth simulation technique, which currently performs realistic simulation on some fabric properties and clothes types. With the development of cloth simulation, we could extend our method to more clothes types. Because of the tremendous manual efforts of dressing simulation for database construction, we now only have male bodies and 3 sets of clothes in our database. We would like to explore automatic dressing simulation technique and extend our work to female and more clothes types situations.

The current feature descriptor which combines 3D landmarks with silhouettes restricts that we should know the camera configuration. The estimation of camera configuration also brings error. In the future, We would explore a new feature descriptor and a new parametric human body model that are more applicable for this problem.

Acknowledgements:

The research is supported in part by NSFC (61572424) and the People Programme (Marie Curie Actions) of the European Unions Seventh Framework Programme FP7 (2007-2013) under

REA grant agreement No.612627-"AniNex". It is partially supported by the grant of the "Sino-UK Higher Education Research Partnership for Ph.D. Studies" Project funded by the Department of Business, Innovation and Skills of the British Government and Ministry of Education of P.R. China. Min Tang is supported in part by NSFC (61572423), Zhejiang Provincial NSFC (LZ16F020003), and the Doctoral Fund of Ministry of Education of China (20130101110133).

References

- [ACP03] ALLEN B., CURLESS B., POPOVIC Z.: The space of human body shapes: reconstruction and parameterization from range scans. *ACM Transactions on Graphics* 22, 3 (JUL 2003), 587–594. 2
- [ASK*05] ANGUELOV D., SRINIVASAN P., KOLLER D., THRUN S., RODGERS J., DAVIS J.: SCAPE: Shape Completion and Animation of People. *ACM Transactions on Graphics* 24, 3 (JUL 2005), 408–416. 2, 3, 4
- [BB08] BALAN A. O., BLACK M. J.: The Naked Truth: Estimating Body Shape Under Clothing. In *Computer Vision - ECCV 2008, PT II, Proceedings* (2008), vol. 5303 of *Lecture Notes in Computer Science*, SPRINGER-VERLAG BERLIN, pp. 15–29. 1, 2, 8
- [BSWX13] BOISVERT J., SHU C., WUHRER S., XI P.: Three-dimensional human shape inference from silhouettes: reconstruction and validation. *Machine Vision and Applications* 24, 1 (JAN 2013), 145–157. 2
- [CC09] CHEN Y., CIPOLLA R.: Learning shape priors for single view reconstruction. In *IEEE 12th International Conference on Computer Vision Workshops (ICCV Workshops)* (2009), pp. 1425–1432. 2
- [CGZZ13] CHEN X., GUO Y., ZHOU B., ZHAO Q.: Deformable model for estimating clothed and naked human shapes from a single image. *Visual Computer* 29, 11 (NOV 2013), 1187–1196. 1, 2, 8
- [CTT*] CHENG K.-L., TONG R.-F., TANG M., SARKIS M., QIAN J.-Y.: Parametric human body reconstruction based on sparse key points. *IEEE Transactions on Visualization and Computer Graphics*. doi:10.1109/TVCG.2015.2511751. 1, 2, 4, 5, 6, 8, 10
- [CWLZ13] CAO C., WENG Y., LIN S., ZHOU K.: 3D Shape Regression for Real-time Facial Animation. *ACM Transactions on Graphics* 32, 4 (JUL 2013). 5, 6
- [CWWS14] CAO X., WEI Y., WEN F., SUN J.: Face Alignment by Explicit Shape Regression. *International Journal of Computer Vision* 107, 2, SI (APR 2014), 177–190. 5, 6
- [DWP10] DOLLAR P., WELINDER P., PERONA P.: Cascaded Pose Regression. In *2010 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2010), IEEE Conference on Computer Vision and Pattern Recognition, IEEE Computer soc, pp. 1078–1085. 5, 6
- [GRWB12] GUAN P., REISS L. AND HIRSHBERG D., WEISS A., BLACK M. J.: Drape: Dressing any person. *ACM Trans. Graphics (Proc. SIGGRAPH)* 31, 4 (jul 2012). 2
- [GWBB09] GUAN P., WEISS A., BALAN A. O., BLACK M. J.: Estimating Human Shape and Pose from a Single Image. In *2009 IEEE 12th International Conference on Computer Vision (ICCV)* (2009), IEEE International Conference on Computer Vision, IEEE, pp. 1381–1388. 2
- [HSR*09] HASLER N., STOLL C., ROSENHAHN B., THORMAEHLEN T., SEIDEL H.-P.: Estimating body shape of dressed humans. *Computers & Graphics-uk* 33, 3, SI (JUN 2009), 211–216. 1, 2, 8
- [HSS*09] HASLER N., STOLL C., SUNKEL M., ROSENHAHN B., SEIDEL H. P.: A Statistical Model of Human Pose and Body Shape. *Computer Graphics Forum* 28, 2 (2009), 337–346. 2, 4
- [JTST10] JAIN A., THORMAEHLEN T., SEIDEL H.-P., THEOBALT C.: MovieReshape: Tracking and Reshaping of Humans in Videos. *ACM Transactions on Graphics* 29, 6 (DEC 2010). 2
- [KSvdP09] KRAEVOY V., SHEFFER A., VAN DE PANNE M.: Modeling from contour drawings. In *Proceedings of the 6th Eurographics Symposium on Sketch-Based Interfaces and Modeling* (2009), pp. 37–44. 3
- [LVG*13] LI H., VOUGA E., GUDYM A., LUO L., BARRON J. T., GUSEV G.: 3D Self-Portraits. *ACM Transactions on Graphics* 32, 6 (NOV 2013). 2
- [NH14] NEOPHYTOU A., HILTON A.: A layered model of human body and garment deformation. In *2014 2nd International Conference on 3D Vision (3DV)* (2014), vol. 1, pp. 171–178. 1, 2, 8
- [OCLF10] OEZUYSAL M., CALONDER M., LEPETIT V., FUA P.: Fast Keypoint Recognition Using Random Ferns. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 32, 3 (MAR 2010), 448–461. 6
- [RCWS14] REN S., CAO X., WEI Y., SUN J.: Face Alignment at 3000 FPS via Regressing Local Binary Features. In *2014 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2014), IEEE Conference on Computer Vision and Pattern Recognition, IEEE, pp. 1685–1692. 6
- [SI03] SIMMONS K. P., ISTOOK C. L.: Body measurement techniques: Comparing 3d body-scanning and anthropometric methods for apparel applications. *Journal of Fashion Marketing and Management: An International Journal* 7, 3 (2003), 306–332. 5
- [SP04] SUMNER R., POPOVIC J.: Deformation transfer for triangle meshes. *ACM Transactions on Graphics* 23, 3 (AUG 2004), 399–405. 2
- [SYW06] SEO H., YEO Y. I., WOHN K.: 3D body reconstruction from photos based on range scan. In *Technologies for E-learning and Digital Entertainment, Proceedings* (2006), vol. 3942 of *Lecture Notes in Computer Science*, Springer-verlag Berlin, pp. 849–860. 2
- [TCL*13] TAM G. K. L., CHENG Z.-Q., LAI Y.-K., LANGBEIN F. C., LIU Y., MARSHALL D., MARTIN R. R., SUN X.-F., ROSIN P. L.: Registration of 3D Point Clouds and Meshes: A Survey from Rigid to Non-rigid. *IEEE Transactions on Visualization and Computer Graphics* 19, 7 (JUL 2013), 1199–1217. 2
- [TZL*12] TONG J., ZHOU J., LIU L., PAN Z., YAN H.: Scanning 3D Full Human Bodies Using Kinects. *IEEE Transactions on Visualization and Computer Graphics* 18, 4 (APR 2012), 643–650. 2
- [WHB11] WEISS A., HIRSHBERG D., BLACK M. J.: Home 3D Body Scans from Noisy Image and Range Data. In *2011 IEEE International Conference on Computer Vision (ICCV)* (2011), IEEE, pp. 1951–1958. 2
- [WPB*14] WUHRER S., PISHCHULIN L., BRUNTON A., SHU C., LANG J.: Estimation of human body shape and posture under clothing. *Computer Vision and Image Understanding* 127 (OCT 2014), 31–42. 1, 2, 8
- [WS13] WUHRER S., SHU C.: Estimating 3D human shapes from measurements. *Machine Vision and Applications* 24, 6 (AUG 2013), 1133–1147. 2
- [ZCD*15] ZENG M., CAO L., DONG H., LIN K., WANG M., TONG J.: Estimation of human body shape and cloth field in front of a kinect. *Neurocomputing* 151, 2 (MAR 5 2015), 626–631. 1, 2, 8
- [ZFL*10] ZHOU S., FU H., LIU L., COHEN-OR D., HAN X.: Parametric Reshaping of Human Bodies in Images. *ACM Transactions on Graphics* 29, 4 (JUL 2010). 2
- [ZM15] ZHU S., MOK P.: Predicting realistic and precise human body models under clothing based on orthogonal-view photos. *Procedia Manufacturing* 3 (2015), 3812–3819. 1, 3, 7, 8
- [ZMK13] ZHU S., MOK P. Y., KWOK Y. L.: An efficient human model customization method based on orthogonal-view monocular photos. *Computer-aided Design* 45, 11 (NOV 2013), 1314–1332. 5



Figure 9: Comparison for body reconstruction results. The first two rows show the reconstruction results for one person and the last two rows show the reconstruction results for another person. The first image in the first row is the range data (of a person with minimal clothes) acquired by Kinect Fusion, which is the input for Cheng's method [CTT*]. The following two images are their reconstruction result in front and back views. Later we show 3 sets of reconstruction results when we use the front-view and side-view silhouettes of the same person with 3 types of clothes separately as input.