

Diffusion Shape Prior for Wrinkle-Accurate Cloth Registration

Jingfan Guo¹Fabian Prada²Donglai Xiang³Javier Romero²Chenglei Wu²Hyun Soo Park¹Takaaki Shiratori²Shunsuke Saito²¹University of Minnesota²Meta Reality Labs³Carnegie Mellon University

Abstract

Registering clothes from 4D scans with vertex-accurate correspondence is challenging, yet important for dynamic appearance modeling and physics parameter estimation from real-world data. However, previous methods either rely on texture information, which is not always reliable, or achieve only coarse-level alignment. In this work, we present a novel approach to enabling accurate surface registration of texture-less clothes with large deformation. Our key idea is to effectively leverage a shape prior learned from pre-captured clothing using diffusion models. We also propose a multi-stage guidance scheme based on learned functional maps, which stabilizes registration for large-scale deformation even when they vary significantly from training data. Using high-fidelity real captured clothes, our experiments show that the proposed approach based on diffusion models generalizes better than surface registration with VAE or PCA-based priors, outperforming both optimization-based and learning-based non-rigid registration methods for both interpolation and extrapolation tests.

1. Introduction

How we dress is important in the perception of identity. The digitization of dynamically deforming clothes is, therefore, one of the core technologies to enable genuine social interaction in virtual environments. This will bring out a myriad of applications including photorealistic telepresence, virtual try-on and visual effects for game and movies. Recently, remarkable progress has been made in computer vision and graphics by modeling photorealistic appearance [49] and plausible geometric deformations [30]. One of the essential building blocks for these approaches is surface registration, which establishes correspondences between a template model and observed 3D reconstruction at each time frame.

Classic methods like Iterative Closest Point (ICP) registration achieve low *surface* error, but suffer from in-plane sliding of the vertices due to the lack of geometric con-

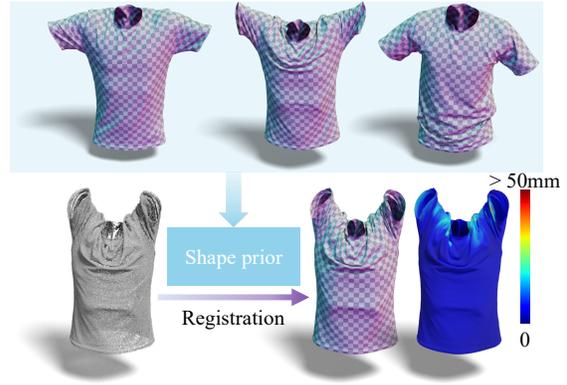


Figure 1. Wrinkle-accurate cloth registration. We learn a strong shape prior from pre-captured 4D data using a diffusion model, and apply it to texture-less registration of the clothing with highly complex deformations.

straints. This problem makes the registration results unsuitable for learning clothes characteristics such as physical parameters or statistical models of deformations. Current registration approaches [49] avoid sliding by relying on texture, i.e., photometric consistency to establish the correspondence between the template and the observed images. Due to the reliance on the texture, the performance of the registration is highly dependent on the uniqueness and contrast of the texture. It is impossible to establish reliable correspondence for regions without salient patterns.

Hand-crafted shape priors like Laplacian [28, 43] and ARAP [42] can be helpful for texture-less cloth registration by regularizing the cloth deformation. They are based on the heuristic that a shape deforms under isometry. However, this may not apply to clothing, because cloth deformation is highly complex, including stretching and bending. The hand-crafted shape priors also require hyperparameter tuning to achieve a balance between registration objective and regularization. To avoid heuristics and parameter tuning, our goal is to learn a shape prior from real-world clothing deformations, which can constrain our solution space within the span of plausible clothing deformations. Existing learning-based approaches like PCA and VAE [2, 14]

have difficulty modeling large deformation and fine details of clothing simultaneously. Noticing the success of diffusion probabilistic models [16, 40, 41] on a wide variety of challenging generative tasks [10, 11, 32, 36], we argue that the representation power of diffusion models is useful for non-rigid surface registration.

In this paper, we present a novel shape prior and illustrate how to use it for surface registration in order to achieve accurate registration of clothing under large motion even without texture information. In particular, we employ a diffusion probabilistic model [16, 40, 41] to learn complex shape distributions of clothing. Given the target 3D point clouds, we employ approximated posterior sampling [10] of the diffusion model with loss functions that optimize the surface alignment. Unfortunately, sampling the posterior directly from scratch results in unstable coarse shape that makes further refinement meaningless. To stably guide the surface registration process, we propose a multi-stage posterior sampling process, where the early stage of the denoising process is guided by a learning-based coarse registration approach [17], and the later stage is only refined with point-to-plane errors. In this way, the registration can avoid local minima while retaining high-fidelity wrinkles with faithful surface deformations.

Real clothing exhibits intricate deformations and interactions with human body parts, which may not be precisely synthesized by physics-based simulation. To evaluate the accuracy of surface registration in real data, we obtain ground-truth correspondence by utilizing a state-of-the-art tracking method based on clothes with a special printed pattern [15]. Experimental results show that our method generalizes to new unseen motion of the garment it was trained on (tested on t-shirts and skirts). In addition, our diffusion-based shape prior significantly outperforms other data-driven shape priors such as PCA and hierarchical VAE as well as state-of-the-art non-rigid registration methods.

Our method has promising applications in 3D cloth modeling. Specifically, after learning the shape prior from a garment with a special printed pattern [15], it can be used to register garments with the same shape but different textures. This enables creating appearance models as in [49], but without requiring the garment to be densely textured.

In summary, our contributions can be summarized as follows:

- a novel diffusion-based shape prior that can effectively encode highly complex clothing geometry.
- a cloth registration approach that leverages the shape prior to achieve accurate cloth registration even in a texture-less setting.
- the first evaluation on ground truth from a wide range of motions and contact of real clothes, quantitatively exposing the accuracy of each registration method in real world scenarios.

2. Related Work

Non-rigid 3D registration is a long-standing problem in the field of computer vision and graphics. In this section we focus on approaches relevant to cloth registration and shape priors. For a more in-depth review of non-rigid 3D registration, please refer to a survey [12].

Optimization-based non-rigid tracking. The goal of non-rigid tracking or registration is to align a stream of unstructured input surfaces with a consistent template mesh, so that the vertices in the template encode the correspondences across different frames. Early work typically uses iterative optimization to find a template deformation which minimizes an energy function including data terms (e.g. target-to-template distance) and regularization terms imposing predefined constraints on the template such as smoothness [43]. Furukawa and Ponce [13] combine rigid local patches with a non-rigid global model for markerless dense 3D tracking. Amberg et al. [1] propose local affine regularization and take optimal greedy steps for non-rigid ICP. Zaharescu et al. [51] propose a 3D feature detector and a 3D feature descriptor for triangle meshes that can be applied to shape matching. Pons-Moll et al. [31] and Xiang et al. [48, 49] use heuristic objectives for non-rigid matching of a template to 4D scans of clothes. SimulCap [50] fits garment templates to a depth video stream by exerting artificial forces in a mass-spring system. Compared with these approaches that rely on hand-crafted regularization constraints, we propose a learning-based approach that effectively leverages garment-specific shape priors directly learned from high-quality ground truth, and therefore achieve more accurate registration.

Learning-based non-rigid tracking. Recent advancement of deep learning has sparked interest in applying deep neural networks to the non-rigid tracking problem. Parameterizing the optimization problem with deep neural networks reduces the number of iterations (typically to a single one) since training data provide better correspondence guesses than the "closest-point" guess used in methods like ICP.

Early learning-based methods use discriminative approaches like regression forest to obtain body correspondences in depth images [44]. The adoption of neural networks further avoids the need of engineering the type of features. A similar approach [19], where the regression forest is replaced by soft classifiers using ResUNet [9], achieves state-of-the-art results in the FAUST dataset [4]. Early work on deep 3D registration like 3D-Coded [14] was based on the PointNet [34] architecture. Numerous representations (3D voxel grids [39], basis point sets [33], zero level sets [3], signed distance functions [6]) have been used to represent the inputs and outputs in these learning-based approaches. Note that deep features can be combined with iterative solvers as in [22], or upgraded to an end-to-end trainable system in Bozic et al. [5]. Most of these methods

are trained in a self-supervised manner using loss functions (e.g. point-to-plane distance) and hand-crafted priors (e.g. smoothness) similar to classical approaches, circumventing the need of accurate ground truth at the cost of fidelity to real deformations. By comparison, our method focuses on utilizing a garment-specific shape prior directly learned from high-quality ground truth to perform more accurate registration.

Shape prior. Deep learning has emerged as a powerful tool to build statistical 3D shape priors directly from data. Such data prior can be useful for various downstream tasks such as animation, reconstruction and tracking. Functional maps [17, 23, 29] are a flexible framework for isometric shape matching, where a shape can be modeled by either deterministic descriptors (e.g., Laplace-Beltrami operator) or learnable descriptors. Different versions of mesh autoencoders (multi-scale [35], an embedded deformation layer [47], and fully convolutional [52]) have been used to model shape variations. Local shape models like PatchNets [46] and DeepLS [7] claim better generalizability. Minimal Neural Atlas [25] models a 3D shape in the parametric domain as a combination of multiple charts, enabling the learning of distortion-minimal parameterization. While these approaches have pushed forward the accuracy of shape generation, implicit surface models lack an explicit parameterization to model in-plane sliding, and approaches like functional maps only provide coarse-level registration when applied to highly non-rigid objects such as clothing.

Learning-based shape priors are also incorporated in clothing or clothed human modeling. Cloth has been modeled independently from the body in a few classic works [21, 30]. TailorNet [30] jointly models the pose, shape and style of clothing, where a high frequency component is responsible for representing fine details. To increase the high frequency details, DeepWrinkles [21] uses conditional GAN to generate high resolution normal maps. Another classic system is CAPE [27], which related the 3D clothes to the underlying body [24] through 3D displacements. All these approaches primarily focus on the synthesis of clothing shapes under different poses, and how to use implicitly learned shape priors for surface registration remains an open question.

Diffusion models. Diffusion models [16, 40, 41] are a class of generative models that can learn the prior from highly complex data distributions by score matching. They have achieved state-of-the-art performance in various image-based generative tasks [10, 11, 36], including a dedicated application to clothing image manipulation [20]. Diffusion models have also been applied to 3D tasks including text-to-3D generation [32], human motion generation [45], point cloud completion [26], and stereo-based human body reconstruction [38]. Specifically, DiffuStereo [38] is closely related to our work, which uses a conditional diffusion model

to refine depth maps for high-quality human body reconstruction. Different from DiffuStereo, which only models small residual deformations in a feed-forward fashion, we propose to estimate both large deformation and detailed deformation in a unified diffusion model by integrating it into an optimization framework.

3. Method

Given the ground-truth 4D scans of cloth in motion, we learn a shape prior using a diffusion model [10, 16] to simultaneously encode large deformation and fine details. The learned shape prior can be used to register the same clothing to noisy 4D scans via multi-stage manifold guidance [10, 11]. In the early stage, our shape prior relies on coarse registration signal to achieve rough alignment. The coarse registration signal can be acquired by markers, visual-based tracking, geometric-based tracking, or any combination of them. In a minimum setting, where markers or visual information are not available, we rely on the geometric information by training SyNoRiM [17] a coarse registration module. In the later stage of manifold guidance, our shape prior further refines the alignment to achieve wrinkle-accurate registration by considering spatial proximity with the input 4D scan.

3.1. Shape Representation

We represent the clothing geometry as a 3D triangle mesh with V vertices $\mathcal{V} \in \mathbb{R}^{V \times 3}$ and F triangles, where the i -th vertex position is denoted as \mathbf{v}_i . The i -th vertex corresponds with at least one point \mathbf{u} on the 2D UV surface. We define the displacement of the mesh from the mean shape with vertices $\bar{\mathcal{V}}$ as a function of UV coordinate, i.e., $\mathcal{U}|_{\mathbf{u}} = \mathbf{v}_i - \bar{\mathbf{v}}_i$, where $\bar{\mathbf{v}}_i$ is the i -th vertex of the mean shape, and $\mathcal{U}|_{\mathbf{u}}$ is the displacement map evaluated at \mathbf{u} . With an abuse of notation, we denote the mapping and its inverse as $\mathcal{U} = \Phi(\mathcal{V})$ and $\mathcal{V} = \Psi(\mathcal{U})$, respectively. Note that this parameterization is not injective, i.e., there exists a vertex that maps to multiple points in \mathcal{U} where these points lie in the boundary (seam) of the unwrapped clothes.

3.2. Diffusion-based Shape Prior

Based on the UV displacement parameterization, we aim to learn a prior distribution of the plausible deformations, which can be used as guidance for cloth registration. We leverage the diffusion model [16] that is made of two processes: forward and reverse. For the forward process, we learn a transition probability from the complete signal to a random noise \mathbf{x}_T by adding noise:

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \beta_t \epsilon, \quad \mathbf{x}_0 = \mathcal{U}, \quad (1)$$

where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is a sample from the Gaussian distribution. We gradually increase the variance schedule β_t as

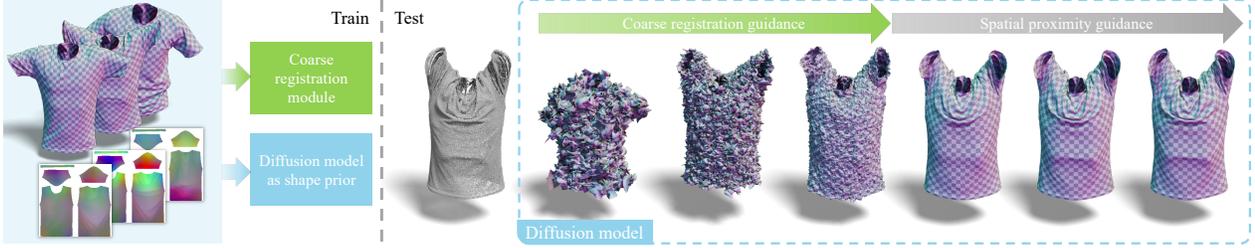


Figure 2. We learn a diffusion-based shape prior from 4D cloth capture, and use it to accurately register the same clothing to noisy scans.



Figure 3. Seam stitching. We enforce the same noise value for corresponding points on the seams in the reverse process.

increasing t , which reduces the impact of \mathbf{x}_{t-1} while increasing that of Gaussian noise. This ensures coarse-to-fine shape learning where the large deformation (low-frequency) are modeled at large t (early denoising stage), and the small deformation (high-frequency) is modeled at small t .

The reverse process, known as ancestral sampling [16], reconstructs the signal from random noise by denoising:

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z} \quad (2)$$

Based on the variance schedule β_t , we define $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$, and $\sigma_t = \sqrt{\tilde{\beta}_t} = \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t}$. The learnable neural network $\boldsymbol{\epsilon}_\theta$ parameterized by θ aims to predict the noise $\boldsymbol{\epsilon}$ from corrupted data \mathbf{x}_t . We train $\boldsymbol{\epsilon}_\theta$ with a weighted variational bound [16] as the objective:

$$L = \mathbb{E}_{t, \mathbf{x}_0, \boldsymbol{\epsilon}} [\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2] \quad (3)$$

Iterating Equation (2) will generate a plausible data sample \mathbf{x}_0 from the learned prior. Figure 4 illustrates the forward and reverse diffusion processes.

Seam stitching. As illustrated in Figure 3, panels of the clothing can be disconnected along seams in UV parameterization. A plausible clothing shape has a smooth surface that maps to smoothly transitioned values across seams. To avoid the clothing being separated apart at the seam, at every time step in the reverse process, we enforce the noise value on the seams to be the same for corresponding points by $\mathbf{x}'_{t-1} = \Phi(\Psi(\mathbf{x}_{t-1}))$. Note that the mapping Ψ from UV to mesh space is not injective, so the ambiguity is solved by averaging the 3D locations of UV positions which refer to the same point.

3.3. Non-rigid Registration via Manifold Guidance

The ancestral sampling in Equation (2) allows generating diverse plausible shape deformations. To ensure that the deformed shape matches the visible surface of the clothing, it is necessary to steer the reverse diffusion process. Noticing that the gradient of log marginal density can be approximated by the learned network $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \simeq -\boldsymbol{\epsilon}_\theta / \sigma_t$, we make use of manifold guidance [10, 11] to sample the near optimal shape deformation that can match the observed points. The manifold guidance maximizes the log-likelihood of every diffusion state, \mathbf{x}_t , given the observation $\mathcal{Y} \in \mathbb{R}^{P \times 3}$ where P is the number of observed 3D points:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathcal{Y}) \simeq -\frac{\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)}{\sigma_t} - \rho \nabla_{\mathbf{x}_t} d(\hat{\mathbf{x}}_0, \mathcal{Y}) \quad (4)$$

where d is a distance measurement in Euclidean space, and ρ is the step size of guidance. The posterior mean $\hat{\mathbf{x}}_0$ can be estimated from \mathbf{x}_t

$$\hat{\mathbf{x}}_0 = \frac{1}{\sqrt{\bar{\alpha}_t}} \mathbf{x}_t - \sqrt{\frac{1 - \bar{\alpha}_t}{\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \quad (5)$$

We take a multi-stage distance measurement d with the decreasing of time step t in the reverse diffusion process

$$d(\hat{\mathbf{x}}_0, \mathcal{Y}) = \begin{cases} \sum_{i=1}^V \varrho(\|\tilde{\mathbf{v}}_i - \hat{\mathbf{v}}_i\|) & t > \tau \\ \sum_{i=1}^P \varrho\left(\left(\mathbf{y}_i - \mathfrak{N}(\hat{\mathcal{V}}, \mathbf{y}_i)\right)^\top \mathbf{n}_{\mathbf{y}_i}\right) & t \leq \tau \end{cases} \quad (6)$$

where $\tilde{\mathbf{v}}_i$ is the target position predicted by the coarse registration module $\mathcal{C}(\mathcal{Y})$, and $\varrho(\cdot)$ is the Huber robust function [18]. The mesh vertices are mapped from the UV displacement map as $\hat{\mathcal{V}} = \Phi(\hat{\mathbf{x}}_0)$ with the i -th vertex denoted as $\hat{\mathbf{v}}_i$. The time step τ is the point where the distance measurement is changed. \mathfrak{N} retrieves the closest point in $\hat{\mathcal{V}}$, and $\mathbf{n}_{\mathbf{y}_i} \in \mathbb{R}^3$ is the normal vector of \mathbf{y}_i .

When $t > \tau$, we use the signals from the coarse registration module \mathcal{C} to guide the large deformation of the clothing shape, until we achieve a rough alignment with the input point cloud. We adopt SyNoRiM [17] as the coarse registration module \mathcal{C} . It learns to predict per-vertex 3D flow from template to target shape given a 3D point cloud as input. The dense correspondences it infers are not strictly necessary since our method works well with sparse ones as shown

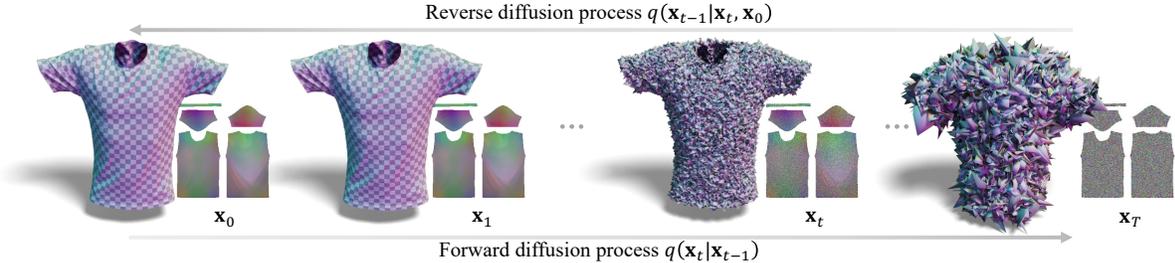


Figure 4. Illustration of the diffusion model for clothing shape. In the forward process, we gradually add noise to the UV displacement map \mathbf{x}_0 to acquire an isotropic Gaussian distribution \mathbf{x}_T . To sample from the learned data distribution, we recover \mathbf{x}_0 by gradually denoising the corrupted UV displacement map.

in Section 4. After $t \leq \tau$, the vertices are guided by point-to-plane errors [8] based on spatial proximity from the point \mathbf{y}_i , which is typically used for non-rigid registration tasks. The point-to-plane distance helps to avoid overestimating the distance when the input point cloud contains holes.

After reaching $t = 0$ in the reverse diffusion process, we repeat the final denoising step with point-to-plane guidance to adjust the inferred vertices to the high-frequency surface details of the point cloud.

4. Experiments

We evaluate our method by comparing with baseline methods quantitatively and qualitatively, demonstrating the effectiveness of the proposed approach. We further discuss the design choice of our method in the ablation study and showcase that our method can be flexible in practical use cases when sparse tracking signals are available. See more experiments in Supp. Mat., including robustness to noisy input and cross-subject generalization.

4.1. Dataset and Settings

Dataset. We conduct our experiments using the pattern-based cloth registration dataset [15], which provides a template geometry for each clothing type, as well as accurate registrations in the same topology. We use the provided data as ground-truth for both training and evaluation. As the dataset does not release the original point clouds, we instead construct partial 3D reconstruction to be used as input at test time. See Supp. Mat. for details.

Since the body is not included in the dataset, we use the ground-truth registration to compute the global translation and rotation of each frame w.r.t. the mean shape using Procrustes analysis, and then normalize the data by applying inverse of the global transformation on each frame. In real use cases where the body is available, we can apply a similar normalization by estimating the body pose, and taking the pelvis transformation as global transformation.

We use T-shirt on "subject_00" (*T-shirt 1*), T-shirt on "subject_04" (*T-shirt 2*, with a stiffer material than *T-shirt 1*), skirt on "subject_03" (*Skirt 1*), and skirt on "subject_04" (*Skirt 2*, longer than *Skirt 1*) in our experiments. For each

data sequence, we split the frames into training set and test set, which further includes interpolation and extrapolation sets. The interpolation test set is uniformly sampled from the entire sequence, so its data distribution is similar to the training set. The extrapolation test set is a manually selected short sequence consisting of body poses unseen in training set. All learning-based methods use this identical train-test split.

Implementation details. We train the diffusion model with $T = 1,000$ steps, and sample a subset $S = 50$ steps using DDIM [41] with $\eta = 0$ in the reverse process. We use a linear variance schedule that increases from $\beta_1 = 10^{-4}$ to $\beta_T = 0.02$. The network ϵ_θ is implemented as a U-Net [37] that takes in a 256×256 UV displacement map. The invalid pixels on the UV map are masked out during both training and testing. We use SyNoRiM [17] as the coarse registration module in our pipeline. It is trained in a pairwise manner between the mean template shape and each training sample with randomly sub-sampled mesh vertices. Note that our method can seamlessly integrate with not only SyNoRiM but also any coarse registration methods.

Metrics. We quantitatively evaluate the performance of our method and baseline methods using vertex error E_v and bidirectional point-to-plane error E_{pt} and E_{ps} .

$$E_v = \frac{1}{V} \sum_{i=1}^V \|\hat{\mathbf{v}}_i - \mathbf{v}_i\| \quad (7)$$

$$E_{pt} = \frac{1}{V} \sum_{i=1}^V \left| \left(\mathbf{v}_i - \mathfrak{N}(\hat{\mathcal{V}}, \mathbf{v}_i) \right)^\top \mathbf{n}_{\mathbf{v}_i} \right| \quad (8)$$

$$E_{ps} = \frac{1}{V} \sum_{i=1}^V \left| \left(\hat{\mathbf{v}}_i - \mathfrak{N}(\mathcal{V}, \hat{\mathbf{v}}_i) \right)^\top \mathbf{n}_{\hat{\mathbf{v}}_i} \right| \quad (9)$$

Our actual goal is to achieve a low E_v for accurate alignment, while E_{pt}/E_{ps} are also important. E_v directly indicates the accuracy of the registration, while E_{pt}/E_{ps} measures only *surface* alignment. When E_{pt}/E_{ps} is small, E_v can still be large due to in-plane sliding. Similarly, lower E_v with higher E_{pt}/E_{ps} is not preferable due to large deviation from the true surface. Our goal is to achieve low E_v with reasonable E_{pt}/E_{ps} .

	<i>T-shirt 1</i>				<i>T-shirt 2</i>				<i>Skirt 1</i>				<i>Skirt 2</i>			
	Int. set		Ext. set		Int. set		Ext. set		Int. set		Ext. set		Int. set		Ext. set	
	E_v	E_{pt}/E_{ps}	E_v	E_{pt}/E_{ps}	E_v	E_{pt}/E_{ps}	E_v	E_{pt}/E_{ps}	E_v	E_{pt}/E_{ps}	E_v	E_{pt}/E_{ps}	E_v	E_{pt}/E_{ps}	E_v	E_{pt}/E_{ps}
SyNoRiM [17]	5.76	1.37/1.68	11.13	1.38/1.67	6.32	1.39/1.80	10.09	1.41/1.80	18.94	1.57/2.76	24.05	1.61/2.87	20.88	1.68/2.70	22.17	1.71/2.67
Lap. reg. [43]	5.04	0.65/0.64	10.80	0.66/0.64	5.60	0.59/0.62	9.61	0.59/0.62	18.33	0.71/0.73	23.71	0.73/0.73	20.47	0.73/0.72	21.82	0.66/0.67
Lap. precond. [28]	5.00	0.53/0.59	10.64	0.53/0.59	5.53	0.52/0.59	9.44	0.51/0.58	18.09	0.55/0.66	23.40	0.56/0.67	20.42	0.60/0.70	21.84	0.57/0.67
ARAP reg. [42]	4.57	0.52/0.60	10.48	0.51/0.60	5.16	0.53/0.60	9.23	0.51/0.59	17.88	0.59/0.75	23.36	0.60/0.76	19.88	0.65/0.77	21.25	0.58/0.68
3D-CODED [14]	11.02	2.50/3.31	21.19	3.01/5.47	13.24	2.93/5.15	14.76	3.45/5.22	18.44	2.03/5.21	24.30	2.25/6.33	17.78	3.37/6.09	28.25	4.51/9.06
3D-CODED opt.	10.26	2.41/3.06	18.81	2.72/4.48	11.25	2.67/4.09	13.72	2.95/4.07	18.14	1.98/4.83	23.47	2.18/5.61	17.07	3.05/5.09	25.40	3.92/6.72
PCA	10.24	2.62/2.93	14.30	2.95/3.75	8.69	2.99/3.41	12.96	3.58/4.36	15.87	3.36/4.10	21.12	4.10/5.09	17.79	4.87/5.27	24.44	6.22/7.09
Comp. VAE [2]	5.05	0.72/0.78	10.74	0.73/0.79	5.67	0.85/0.93	9.63	0.88/0.95	18.04	0.64/0.72	23.47	0.65/0.74	20.57	0.79/0.83	21.94	0.79/0.86
Ours	3.16	0.57/0.62	9.51	0.61/0.75	3.94	0.57/0.63	8.59	0.61/0.76	15.07	0.65/0.73	21.01	0.73/0.78	16.66	0.71/0.78	19.71	0.67/0.75

Table 1. Quantitative comparison to baseline methods. Error metrics are measured in mm. **Bold** indicates the best E_v . Our goal is to achieve low vertex error E_v with reasonable point-to-plane error E_{pt} and E_{ps} .

4.2. Comparison to Baseline Methods

SyNoRiM [17] is a general-purpose non-rigid registration method. Since SyNoRiM tends to produce over-smoothed results lacking fine details like wrinkles, we further refine SyNoRiM predictions by optimizing point-to-plane distance together with heuristic shape priors (Laplacian [28, 43] and ARAP [42]) as in classical ICP methods. In Table 1, we quantitatively show that our full pipeline consistently outperforms SyNoRiM and its heuristic refinement on the metric of vertex error, which is our main goal. Qualitative results in Figure 5, 6 show that our method consistently produces better registration with lower vertex error and realistic wrinkles.

We also compare our approach to data-driven shape priors 3D-CODED [14], PCA, and compositional VAE [2]. In 3D-CODED, we model the 3D point translation and 3D patch deformation following the original setting. The 3D-CODED results can be further refined by optimizing Chamfer distance (denoted 3D-CODED opt.). For PCA, we model the per-vertex 3D displacement from mean shape, and keep a number of principal components that retains 95% explained variance. At test time, we estimate the PCA coefficients by solving least squares to dense target points given by SyNoRiM. For compositional VAE, we model the UV displacement map similar to our setting. At test time, we initialize the latent code by feeding the SyNoRiM result to the encoder, then optimize the latent code by minimizing point-to-plane distance in Equation (8). Table 1 shows that our shape prior is more effective than baseline data-driven shape priors. Please note that PCA achieves comparable E_v to our method on *Skirt 1* and *Skirt 2*, but it shows significantly worse plane error E_{pt} and E_{ps} . As a linear model, PCA may not be suitable for this inherently nonlinear problem, so it is not flexible enough to achieve accurate surface-level alignment. It cannot fit to large deformations that are far from the mean shape as shown in Figure 5, 6, even though they are from the interpolation set and close to some training samples. Comparing to all baseline methods, Figure 7 shows that our method consistently achieves the best balance of lower E_v and E_{pt}/E_{ps} .

	<i>T-shirt 1</i>			
	Interpolation set		Extrapolation set	
	E_v	E_{pt} / E_{ps}	E_v	E_{pt} / E_{ps}
$\tau = 40$	32.16	0.67 / 1.29	36.46	0.72 / 1.93
$\tau = 30$	3.39	0.57 / 0.62	9.40	0.61 / 0.73
$\tau = 20$	3.16	0.57 / 0.62	9.51	0.61 / 0.75
$\tau = 10$	3.47	0.58 / 0.64	9.82	0.63 / 0.79
$\tau = 0$	3.82	0.59 / 0.67	10.10	0.66 / 0.86

Table 2. The effect of the guidance breakpoint τ . A large τ significantly impair the performance.

4.3. Ablation Study

Guidance breakpoint τ . In Table 2, we quantitatively show the impact of varying the guidance breakpoint τ in Equation (6) on the *T-shirt 1* sequence. The performance with large τ is significantly worse, indicating that the first stage of the manifold guidance plays a key role, and a coarse registration module is necessary in our method. When τ is in a reasonable range, it does not significantly affect the performance, although a too small τ amplifies the influence of error from the coarse registration.

Seam stitching. Given a 2D parameterization with multiple islands for the clothing, it is important to enforce the continuity across the seams. Figure 8 illustrates that the proposed seam stitching strategy prevents generating implausible shape with separated clothing parts.

4.4. Application

In real-world scenarios, clothing usually comes with textures that make visual keypoint tracking possible. Our method is flexible in that it can take advantage of such information when available. To mimic the use case where keypoint tracking is available, we experiment with a synthetic setting where coarse registration is replaced by sparse ground-truth guidance, as it could be provided by perfectly accurate sparse texture tracking. Specifically, we randomly select N_k vertices from the ground-truth mesh, and use them to compute the distance in Equation 6 in the first stage of manifold guidance with $t > \tau$. This replaces the coarse registration module, so coarse registration module is not used under this setting.

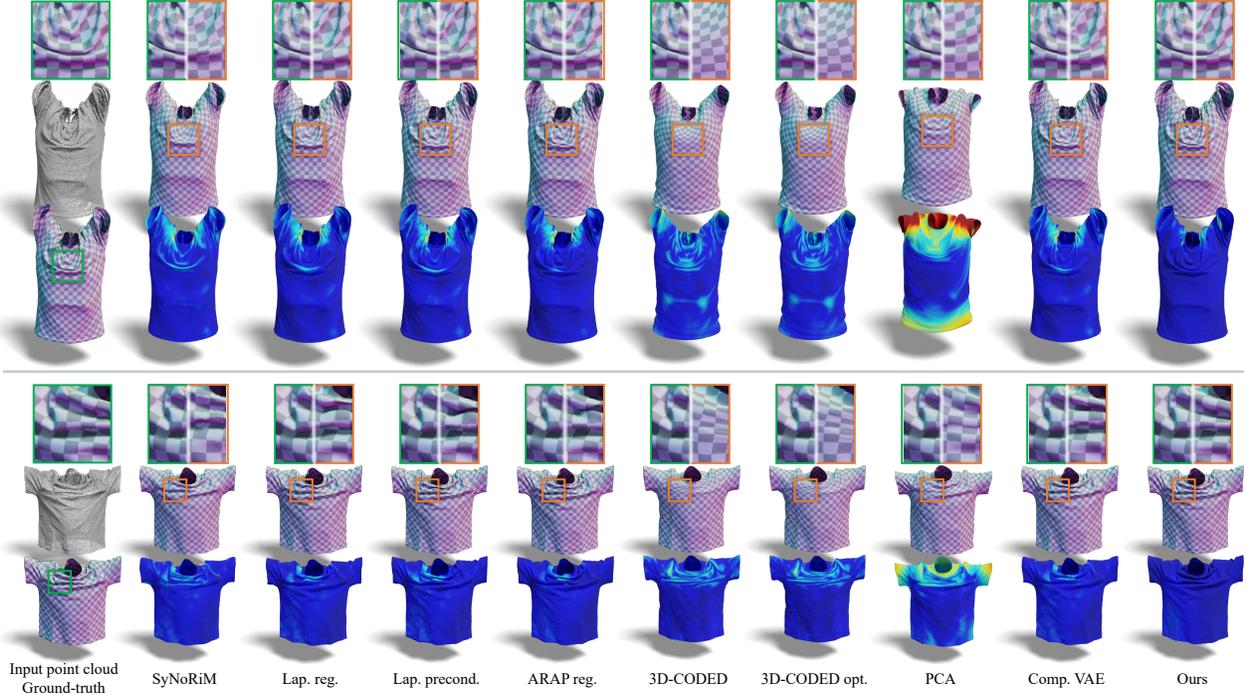


Figure 5. Comparison to baseline methods on *T-shirt 1* and *T-shirt 2*. In each example, the middle-left is the input point cloud, the bottom-left is the ground-truth, the top-left is zoom-in view of ground-truth. The rest are the results of different methods, where the top row shows side-by-side comparison to ground-truth, the middle row shows the geometry with normal rendering, while the bottom row shows vertex error E_v in color ($0mm$ $> 50mm$).

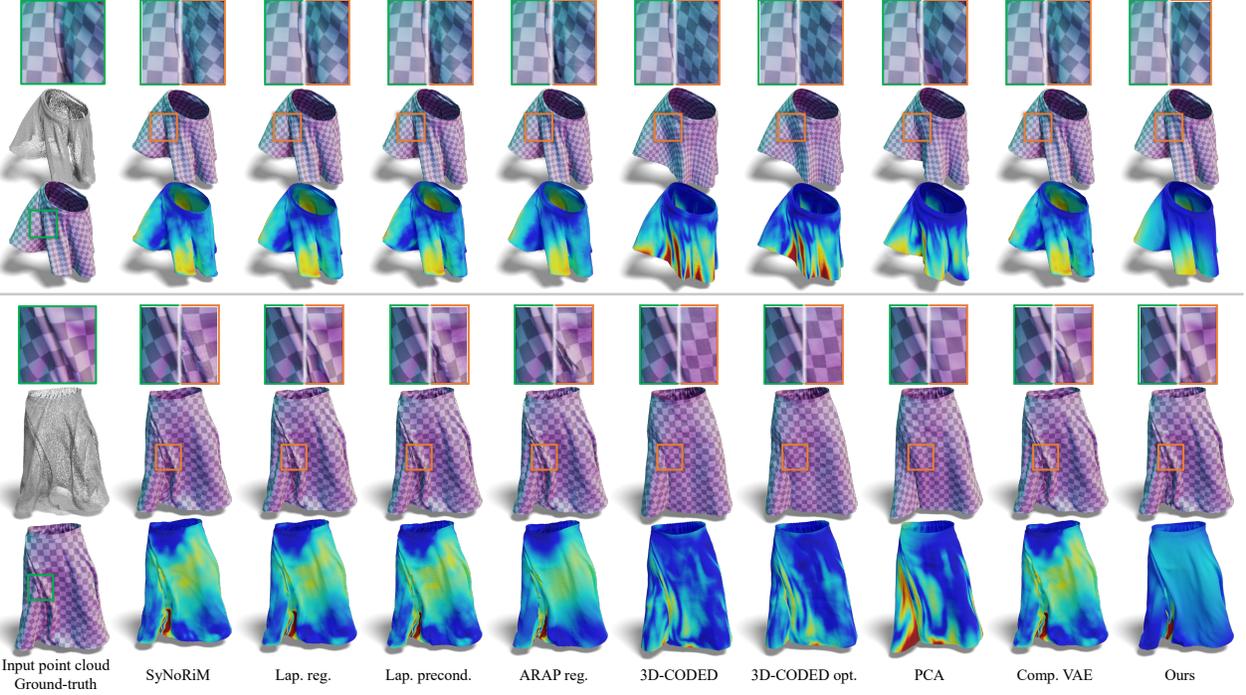


Figure 6. Comparison to baseline methods on *Skirt 1* and *Skirt 2*. See Figure 5 for explanation, and Supp. Mat. for more results.

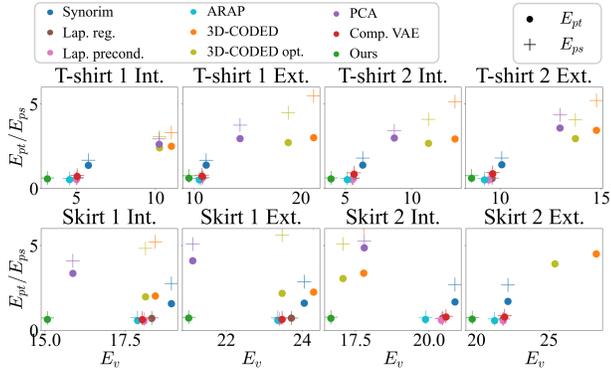


Figure 7. Error plot of E_v vs. E_{pt}/E_{ps} . Closer to origin (lower E_v and E_{pt}/E_{ps}) is a preferred solution.

	N_k	T-shirt 1			
		Interpolation set		Extrapolation set	
		E_v	E_{pt} / E_{ps}	E_v	E_{pt} / E_{ps}
PCA	50	10.65	2.67 / 3.08	13.86	2.98 / 3.87
	100	10.22	2.60 / 2.94	13.36	2.91 / 3.72
	200	10.28	2.61 / 2.97	13.55	2.95 / 3.78
Comp. VAE	50	30.08	1.67 / 6.00	37.16	1.45 / 8.23
	100	28.41	1.24 / 5.94	34.55	1.07 / 7.93
	200	25.41	1.03 / 5.67	30.54	0.87 / 7.44
Ours	50	3.59	0.58 / 0.65	6.10	0.62 / 0.79
	100	2.62	0.57 / 0.62	4.89	0.61 / 0.74
	200	2.39	0.57 / 0.62	4.53	0.61 / 0.74

Table 3. Sparse ground-truth guidance. Our method works with sparse tracking signals. PCA performs similarly to the original setting, while the performance of compositional VAE significantly decreases.



Figure 8. The effect of seam stitching strategy. Without seam stitching, the generated clothing shape may have separated parts, as continuity is not enforced across seams.

From Table 3 and Figure 9, we can see that our method performs well with very sparse keypoint tracking signals, and the accuracy improves when the number of sparse keypoints increases. As a compact linear model, PCA performs similarly to the original setting, but increasing the number of keypoint does not help. The compositional VAE fails to learn a meaningful latent space for plausible clothing shapes, because large deformation and fine wrinkles are coupled together. It may require a significant amount of tracking signals to find a latent code corresponding to a plausible clothing shape.

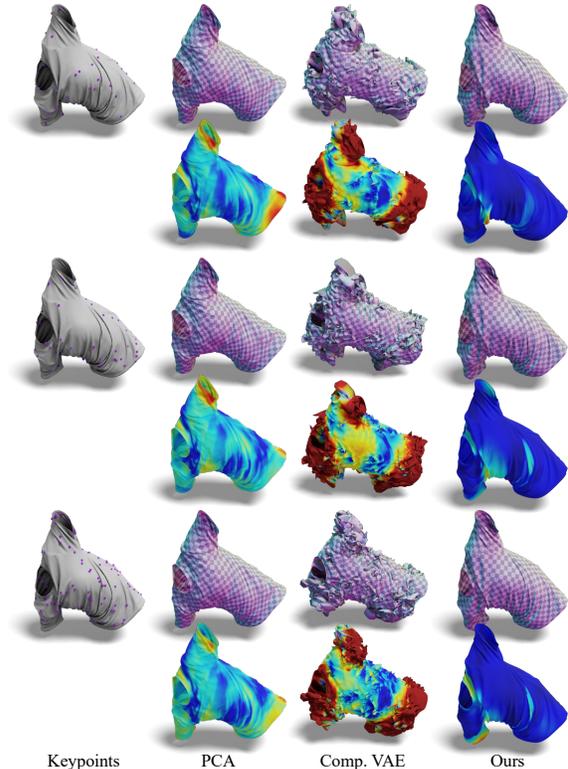


Figure 9. Sparse ground-truth guidance. Keypoints used for tracking are shown in purple in the left column. Our method works well with very sparse tracking signal. PCA cannot reproduce the correct pose or wrinkles, and Compositional VAE fails to generate a plausible clothing shape.

5. Conclusion

We have presented a diffusion-based shape prior for highly deformable clothing geometry, and how the prior can be incorporated into fine-grained non-rigid registration tasks. Our approach, for the first time, achieves the adoption of diffusion models into 3D cloth modeling by leveraging UV parameterization. Our experiments using real data show the versatility of the proposed multi-stage manifold guidance, demonstrating superior performance with multiple clothing types and diverse motions. As our approach is simple and general, we believe it can open a new venue for various 3D optimization problems that benefit from strong 3D shape priors.

Limitations and future work. As our early stage guidance relies on an off-the-shelf coarse registration module, large error introduced by this module cannot be fully removed in the following refinement stage. Eliminating the need of the coarse registration or building a more robust shape prior is an interesting venue for future work. Also, UV parameterization limits the capability of cross-garment generalization, and leveraging a single UV parameterization may be non-trivial for more complex clothing, which could be addressed by extending diffusion models to other 3D representations.

References

- [1] Brian Amberg, Sami Romdhani, and Thomas Vetter. Optimal step nonrigid icp algorithms for surface registration. In *CVPR*, pages 1–8, 2007. 2
- [2] Timur Bagautdinov, Chenglei Wu, Jason Saragih, Pascal Fua, and Yaser Sheikh. Modeling facial geometry using compositional vaes. In *CVPR*, pages 3877–3886, 2018. 1, 6, 3
- [3] Bharat Lal Bhatnagar, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Loopreg: Self-supervised learning of implicit surface correspondences, pose and shape for 3d human mesh registration. In *NeurIPS*, pages 12909–12922, 2020. 2
- [4] Federica Bogo, Javier Romero, Gerard Pons-Moll, and Michael J Black. Dynamic FAUST: Registering human bodies in motion. In *CVPR*, pages 6233–6242, 2017. 2
- [5] Aljaz Bozic, Pablo Palafox, Michael Zollhöfer, Angela Dai, Justus Thies, and Matthias Nießner. Neural non-rigid tracking. In *NeurIPS*, pages 18727–18737, 2020. 2
- [6] Aljaz Bozic, Pablo Palafox, Michael Zollhofer, Justus Thies, Angela Dai, and Matthias Nießner. Neural deformation graphs for globally-consistent non-rigid reconstruction. In *CVPR*, pages 1450–1459, 2021. 2
- [7] Rohan Chabra, Jan E Lenssen, Eddy Ilg, Tanner Schmidt, Julian Straub, Steven Lovegrove, and Richard Newcombe. Deep local shapes: Learning local sdf priors for detailed 3d reconstruction. In *ECCV*, pages 608–625, 2020. 3
- [8] Yang Chen and Gérard Medioni. Object modelling by registration of multiple range images. *Image and vision computing*, pages 145–155, 1992. 5
- [9] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *CVPR*, pages 3075–3084, 2019. 2
- [10] Hyungjin Chung, Jeongsol Kim, Michael T Mccann, Marc L Klasky, and Jong Chul Ye. Diffusion posterior sampling for general noisy inverse problems. In *ICLR*, 2022. 2, 3, 4
- [11] Hyungjin Chung, Byeongsu Sim, Dohoon Ryu, and Jong Chul Ye. Improving diffusion models for inverse problems using manifold constraints. In *NeurIPS*, 2022. 2, 3, 4
- [12] Bailin Deng, Yuxin Yao, Roberto M Dyke, and Juyong Zhang. A survey of non-rigid 3d registration. In *Comput. Graph. Forum*, pages 559–589, 2022. 2
- [13] Yasutaka Furukawa and Jean Ponce. Dense 3d motion capture from synchronized video streams. In *CVPR*, pages 1–8, 2008. 2
- [14] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. 3d-coded: 3d correspondences by deep deformation. In *ECCV*, pages 230–246, 2018. 1, 2, 6, 3
- [15] Oshri Halimi, Tuur Stuyck, Donglai Xiang, Timur Bagautdinov, He Wen, Ron Kimmel, Takaaki Shiratori, Chenglei Wu, Yaser Sheikh, and Fabian Prada. Pattern-based cloth registration and sparse-view animation. *ACM TOG*, pages 1–17, 2022. 2, 5
- [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, pages 6840–6851, 2020. 2, 3, 4
- [17] Jiahui Huang, Tolga Birdal, Zan Gojcic, Leonidas J Guibas, and Shi-Min Hu. Multiway non-rigid point cloud registration via learned functional map synchronization. *IEEE TPAMI*, 2022. 2, 3, 4, 5, 6
- [18] Peter J Huber. Robust estimation of a location parameter. In *Breakthroughs in statistics: Methodology and distribution*, pages 492–518, 1992. 4
- [19] Hyomin Kim, Jungeon Kim, Jaewon Kam, Jaesik Park, and Seungyong Lee. Deep virtual markers for articulated 3d shapes. In *ICCV*, pages 11615–11625, 2021. 2
- [20] Chaerin Kong, DongHyeon Jeon, Ohjoon Kwon, and Nojun Kwak. Leveraging off-the-shelf diffusion model for multi-attribute fashion image manipulation. In *WACV*, pages 848–857, 2023. 3
- [21] Zorah Lahner, Daniel Cremers, and Tony Tung. Deepwrinkles: Accurate and realistic clothing modeling. In *ECCV*, pages 667–684, 2018. 3
- [22] Yang Li, Aljaz Bozic, Tianwei Zhang, Yanli Ji, Tatsuya Harada, and Matthias Nießner. Learning to optimize non-rigid tracking. In *CVPR*, pages 4910–4918, 2020. 2
- [23] Or Litany, Tal Remez, Emanuele Rodola, Alex Bronstein, and Michael Bronstein. Deep functional maps: Structured prediction for dense shape correspondence. In *ICCV*, pages 5659–5667, 2017. 3
- [24] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *ACM TOG*, pages 1–16, 2015. 3
- [25] Weng Fei Low and Gim Hee Lee. Minimal neural atlas: Parameterizing complex surfaces with minimal charts and distortion. In *ECCV*, 2022. 3
- [26] Zhaoyang Lyu, Zhifeng Kong, XU Xudong, Liang Pan, and Dahua Lin. A conditional point diffusion-refinement paradigm for 3d point cloud completion. In *ICLR*, 2022. 3
- [27] Qianli Ma, Jinlong Yang, Anurag Ranjan, Sergi Pujades, Gerard Pons-Moll, Siyu Tang, and Michael J Black. Learning to dress 3d people in generative clothing. In *CVPR*, pages 6469–6478, 2020. 3
- [28] Baptiste Nicolet, Alec Jacobson, and Wenzel Jakob. Large steps in inverse rendering of geometry. *ACM TOG*, pages 1–13, 2021. 1, 6, 3
- [29] Maks Ovsjanikov, Mirela Ben-Chen, Justin Solomon, Adrian Butscher, and Leonidas Guibas. Functional maps: a flexible representation of maps between shapes. *ACM TOG*, pages 1–11, 2012. 3
- [30] Chaitanya Patel, Zhouyingcheng Liao, and Gerard Pons-Moll. Tailornet: Predicting clothing in 3d as a function of human pose, shape and garment style. In *CVPR*, pages 7365–7375, 2020. 1, 3
- [31] Gerard Pons-Moll, Sergi Pujades, Sonny Hu, and Michael J Black. Clothcap: Seamless 4d clothing capture and retargeting. *ACM TOG*, pages 1–15, 2017. 2
- [32] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 2, 3
- [33] Sergey Prokudin, Christoph Lassner, and Javier Romero. Efficient learning on point clouds with basis point sets. In *ICCV*, pages 4332–4341, 2019. 2

- [34] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, pages 652–660, 2017. 2
- [35] Anurag Ranjan, Timo Bolkart, Soubhik Sanyal, and Michael J Black. Generating 3d faces using convolutional mesh autoencoders. In *ECCV*, pages 704–720, 2018. 3
- [36] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022. 2, 3
- [37] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, pages 234–241, 2015. 5
- [38] Ruizhi Shao, Zerong Zheng, Hongwen Zhang, Jingxiang Sun, and Yebin Liu. Diffustereo: High quality human reconstruction via diffusion-based stereo using sparse cameras. In *ECCV*, pages 702–720, 2022. 3
- [39] Soshi Shimada, Vladislav Golyanik, Edgar Tretschk, Didier Stricker, and Christian Theobalt. Dispvoxnets: Non-rigid point set alignment with supervised learning proxies. In *3DV*, pages 27–36, 2019. 2
- [40] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, pages 2256–2265, 2015. 2, 3
- [41] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021. 2, 3, 5
- [42] Olga Sorkine and Marc Alexa. As-rigid-as-possible surface modeling. In *SGP*, pages 109–116, 2007. 1, 6, 3
- [43] Olga Sorkine, Daniel Cohen-Or, Yaron Lipman, Marc Alexa, Christian Rössl, and H-P Seidel. Laplacian surface editing. In *SGP*, pages 175–184, 2004. 1, 2, 6, 3
- [44] Jonathan Taylor, Jamie Shotton, Toby Sharp, and Andrew Fitzgibbon. The vitruvian manifold: Inferring dense correspondences for one-shot human pose estimation. In *CVPR*, pages 103–110, 2012. 2
- [45] Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Amit H Bermano, and Daniel Cohen-Or. Human motion diffusion model. *arXiv preprint arXiv:2209.14916*, 2022. 3
- [46] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Carsten Stoll, and Christian Theobalt. Patchnets: Patch-based generalizable deep implicit 3d shape representations. In *ECCV*, pages 293–309, 2020. 3
- [47] Edgar Tretschk, Ayush Tewari, Michael Zollhöfer, Vladislav Golyanik, and Christian Theobalt. Demea: Deep mesh autoencoders for non-rigidly deforming objects. In *ECCV*, pages 601–617, 2020. 3
- [48] Donglai Xiang, Fabian Prada, Timur Bagautdinov, Weipeng Xu, Yuan Dong, He Wen, Jessica Hodgins, and Chenglei Wu. Modeling clothing as a separate layer for an animatable human avatar. *ACM TOG*, pages 1–15, 2021. 2
- [49] Donglai Xiang, Timur Bagautdinov, Tuur Stuyck, Fabian Prada, Javier Romero, Weipeng Xu, Shunsuke Saito, Jingfan Guo, Breannan Smith, Takaaki Shiratori, et al. Dressing avatars: Deep photorealistic appearance for physically simulated clothing. *ACM TOG*, 2022. 1, 2
- [50] Tao Yu, Zerong Zheng, Yuan Zhong, Jianhui Zhao, Qionghai Dai, Gerard Pons-Moll, and Yebin Liu. Simulcap: Single-view human performance capture with cloth simulation. In *CVPR*, pages 5504–5514, 2019. 2
- [51] Andrei Zaharescu, Edmond Boyer, Kiran Varanasi, and Radu Horaud. Surface feature detection and description with applications to mesh matching. In *CVPR*, pages 373–380, 2009. 2
- [52] Yi Zhou, Chenglei Wu, Zimo Li, Chen Cao, Yuting Ye, Jason Saragih, Hao Li, and Yaser Sheikh. Fully convolutional mesh autoencoder using efficient spatially varying kernels. In *NeurIPS*, pages 9251–9262, 2020. 3

Diffusion Shape Prior for Wrinkle-Accurate Cloth Registration

Supplementary Material

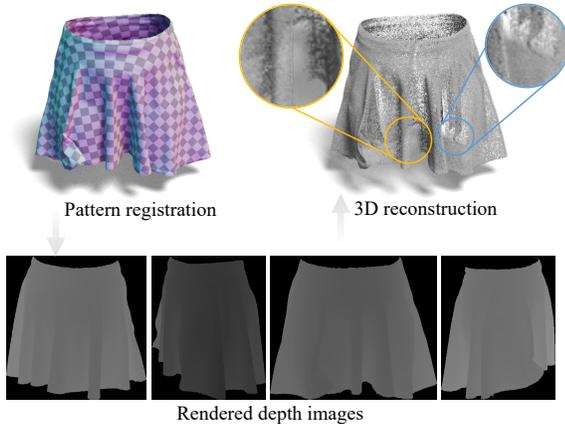


Figure 10. We generate 3D point clouds to mimic the typical 3D reconstruction acquired by a 3D capture system. We render multi-view depth images using the ground-truth registration, then we fuse the multi-view depth images and sub-sample the 3D point cloud.

	<i>T-shirt 1</i>	<i>T-shirt 2</i>	<i>Skirt 1</i>	<i>Skirt 2</i>
Training	3575	4463	1699	2734
Int. set	188	234	90	144
Ext. set	200	150	150	150
Total	3963	4847	1939	3028

Table 4. Frame count in each sequence.

A. More experiment details

Dataset. As illustrated in Figure 10, we render seven depth images for the ground-truth meshes from side and top views, then fuse them into a 3D point cloud with proper occlusion reasoning.

For each garment, there is a long sequence where the actor performs various movements (clips). All clips in the sequence form the full dataset. For extrapolation, we select a clip where the actor performs a rare movement with unique cloth deformation, not present in the rest of the sequence. For the rest of the clips, every 20th frame is selected as interpolation set, while the remaining frames are the training set. The number of frames in each set is shown in Table 4.

Running time. We conduct the experiments on an NVIDIA Tesla V100 GPU with 32GB memory. For each garment, we train the diffusion model for 100k iterations, which takes 20 hours. The inference time is 53.57 seconds/frame.

SyNoRiM refinement in baseline comparison. We use the Adam (UniformAdam for "Lap. precondition.") with step size 10^{-3} for 200 iterations in all the refinement experiments. The weights for the regularization terms are $\lambda_{\text{Lap. reg.}} = 10^3$,

		<i>T-shirt 1</i>			
		Int. set		Ext. set	
		E_v	E_{pt} / E_{ps}	E_v	E_{pt} / E_{ps}
No noise		3.16	0.57 / 0.62	9.51	0.61 / 0.75
Gaussian noise	$\sigma = 1$	3.77	0.59 / 0.65	9.60	0.64 / 0.79
	$\sigma = 2$	5.55	0.64 / 0.73	10.23	0.70 / 0.89
	$\sigma = 3$	7.55	0.73 / 0.89	11.11	0.81 / 1.08
	$\sigma = 4$	9.48	0.85 / 1.12	12.09	0.92 / 1.32
	$\sigma = 5$	11.30	0.96 / 1.37	13.21	1.03 / 1.60
Laplace noise	$b = 3$	10.74	0.82 / 1.10	12.66	0.89 / 1.31
	$b = 4$	13.53	0.96 / 1.43	14.51	1.03 / 1.68
	$b = 5$	17.73	1.11 / 1.93	16.96	1.17 / 2.14

Table 5. The quantitative evaluation of our method on noisy input point cloud. Error metrics are measured in mm.

$$\lambda_{\text{Lap. precondition.}} = 10 \text{ (Eq. (14) in [27])}, \text{ and } \lambda_{\text{ARAP reg.}} = 10^{-3}.$$

B. Registration to noisy point cloud

We evaluate the robustness of our registration method to noisy measurement. Specifically, we perturb the input point cloud of the *T-shirt 1* sequence by adding Gaussian noise and Laplace noise to the 3D point locations

$$\mathbf{y}_i^G = \mathbf{y}_i + \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \quad (10)$$

$$\mathbf{y}_i^L = \mathbf{y}_i + \mathcal{L}(\mathbf{0}, b\mathbf{I}) \quad (11)$$

where \mathbf{y}_i is the i -th point of the input point cloud \mathcal{Y} . \mathbf{y}_i^G is the point cloud perturbed by Gaussian noise, with σ being the standard deviation of the Gaussian noise. \mathbf{y}_i^L is the point cloud perturbed by Laplace noise, with b being the scale of the Laplace noise. Comparing to Gaussian noise, Laplace noise has longer tail noise distribution that can mimic outliers. In our experiments, we add Gaussian noise to the input point cloud with $\sigma = 1, 2, 3, 4, 5$ mm, and Laplace noise with $b = 3, 4, 5$ mm, as shown in the top row of Figure 13 and Figure 14.

We take the noisy point cloud as input, and quantitatively evaluate the registration result of our method in Table 5. We also compare our method to baseline methods in Figure 11, Figure 13 and Figure 14. The performance of our method drops with the increase of noise level, but it consistently outperforms the baseline methods. Please note in Figure 11 (c), PCA shows comparable E_v to our method when $b = 5$. However, PCA consistently shows worse plane error E_{pt} and E_{ps} as discussed in the main manuscript. This can be qualitatively verified in Figure 14 as well.

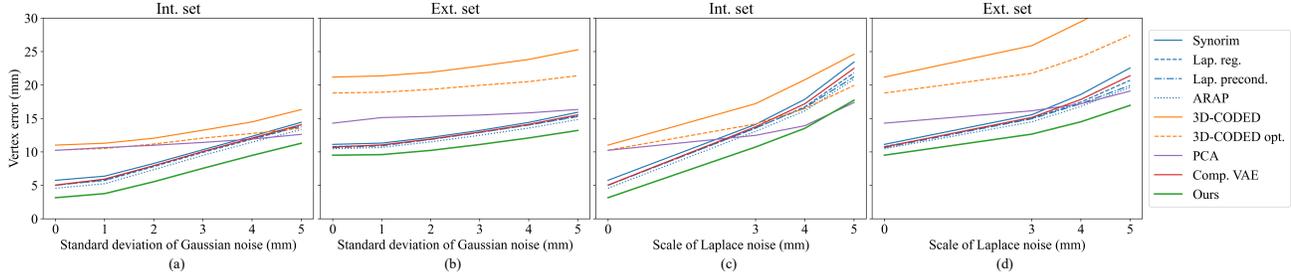


Figure 11. Vertex error of the registration results given noisy point cloud as inputs. (a) and (b) Adding Gaussian noise to input point cloud, then testing on interpolation set and extrapolation set, respectively. (c) and (d) Adding Laplace noise to input point cloud, then testing on interpolation set and extrapolation set, respectively



Figure 12. The difference between predicting data \mathbf{x}_0 and predicting noise ϵ . Data prediction cannot represent fine details like wrinkles.

C. Data prediction vs. noise prediction

Conceptually, it is equivalent to use a data prediction network that predicts \mathbf{x}_0 , or a noise prediction network that predicts ϵ in the diffusion model. In practice, however, we find that noise prediction is more effective. As shown in Figure 12, a data prediction network has difficulty modeling high-frequency signals like wrinkles. It cannot enforce continuity across the seams even if the seam stitching strategy is applied.

D. Cross-subject generalization

The same clothing worn by different subjects may deform differently because of the variation of body shapes. In Table 6 and Figure 15, we empirically show that the proposed method can generalize to unseen subjects. Specifically, *T-shirt 1* is worn by 4 different actors in 4 long sequences "subject_00", "subject_01", "subject_02", and "subject_03". In this experiment, we take *T-shirt 1* on all subjects as the full dataset, and divide it into 4 parts, each contains one subject. We conduct 4-fold cross-validation, where each time we use 3 subjects for training, and the 4th subject for validation. In each validation, the training set consists of

all frames of the 3 subjects, while the validation set only contains every 20th frame of the 4th subject. In Table 6, we report the quantitative evaluation of the 4-fold cross-validation, showing that the proposed method consistently outperforms the baseline methods on the cross-subject generalization task.

Please note that *T-shirt 1* on "subject_01", "subject_02", and "subject_03" are only used in the cross-subject generalization experiment, while other experiments on *T-shirt 1* are done with "subject_00".

E. Additional qualitative results

We show more qualitative results on *T-shirt 1*, *T-shirt 2*, *Skirt 1* and *Skirt 2* sequences in Figure 16, Figure 17, Figure 18 and Figure 19, respectively.

Validation subject	subject_00		subject_01		subject_02		subject_03	
	E_v	E_{pt}/E_{ps}	E_v	E_{pt}/E_{ps}	E_v	E_{pt}/E_{ps}	E_v	E_{pt}/E_{ps}
SyNoRiM [17]	10.62	1.38 / 1.66	13.59	1.54 / 2.16	8.50	1.33 / 1.57	10.67	1.40 / 1.80
Lap. reg. [43]	10.32	0.66 / 0.64	12.99	0.64 / 0.67	8.13	0.61 / 0.62	10.21	0.65 / 0.65
Lap. precondition. [28]	10.09	0.53 / 0.59	12.97	0.55 / 0.63	8.10	0.52 / 0.58	10.20	0.56 / 0.62
ARAP reg. [42]	10.04	0.52 / 0.60	12.57	0.54 / 0.65	7.87	0.51 / 0.58	9.94	0.54 / 0.65
3D-CODED [14]	36.94	3.93 / 15.77	18.17	4.09 / 6.05	23.91	4.03 / 6.76	15.80	3.25 / 4.58
3D-CODED opt.	33.03	3.13 / 13.13	16.04	3.00 / 4.17	20.52	3.35 / 5.34	14.14	2.77 / 3.72
PCA	11.65	3.02 / 3.44	12.11	2.72 / 3.28	10.19	2.87 / 3.27	10.41	2.38 / 2.71
Comp. VAE [2]	10.35	1.21 / 1.23	13.12	1.06 / 1.19	8.21	1.05 / 1.08	10.27	0.97 / 1.04
Ours	9.85	0.61 / 0.69	11.77	0.60 / 0.75	7.58	0.56 / 0.68	9.69	0.60 / 0.72

Table 6. Quantitative results for cross-subject generalization. For each validation, the models are trained on the other 3 subjects. Error metrics are measured in mm. **Bold** indicates the best E_v .

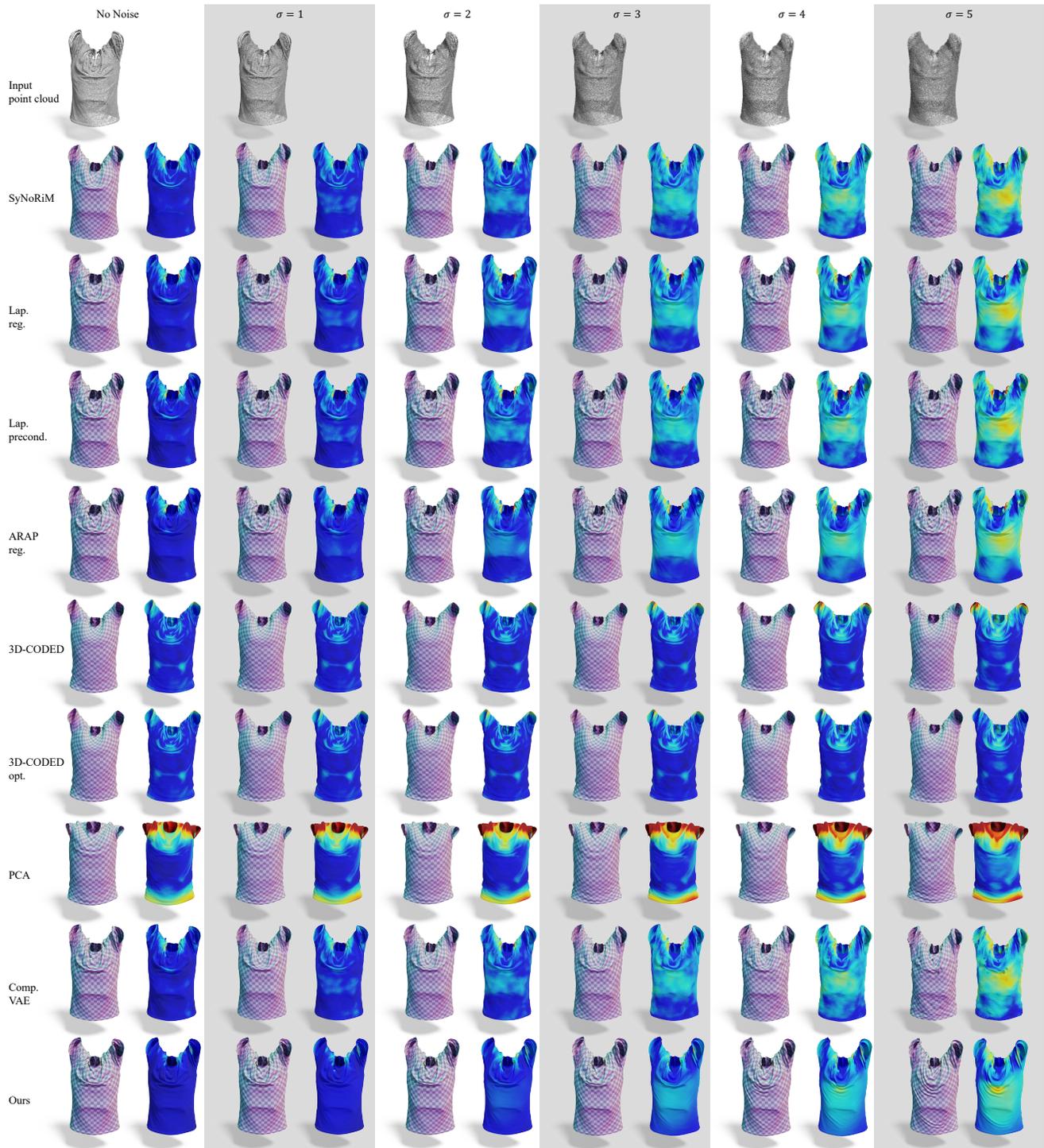


Figure 13. Registration results of baseline methods and our method on noisy point cloud, where σ is the standard deviation of the Gaussian noise measured in mm. Vertex error E_v is shown in color ($0mm$ $> 50mm$).

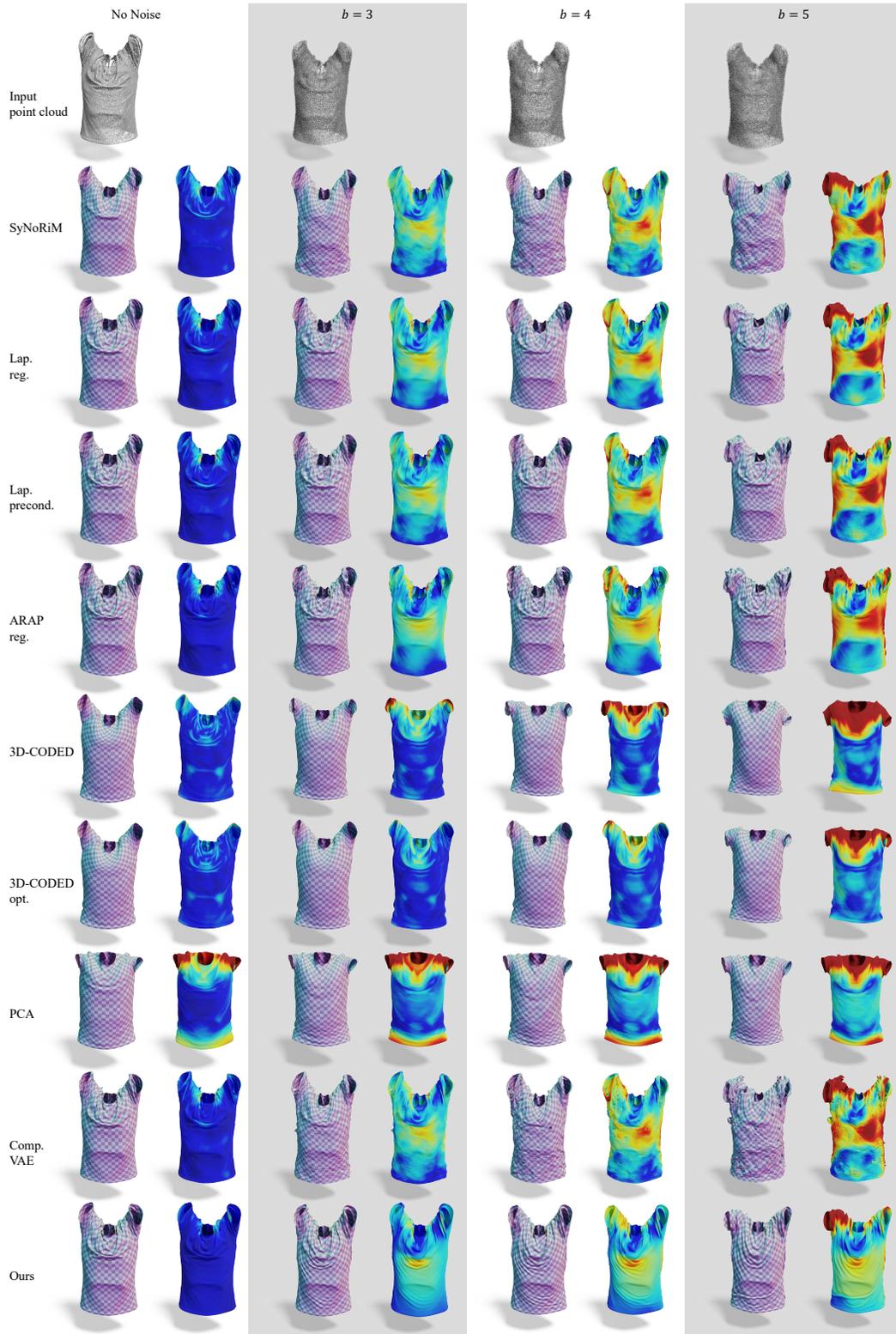


Figure 14. Registration results of baseline methods and our method on noisy point cloud, where b is the scale of the Laplace noise measured in mm. Vertex error E_v is shown in color ($0mm$ $> 50mm$).



Figure 15. Qualitative comparison for cross-subject generalization. Each example is from a different subject. For each example, the model is trained on the other 3 subject, so the test subject is always unseen during training. See Figure 16 for explanation of the figure.



Figure 16. Comparison to baseline methods on *T-shirt 1*. In each example, the middle-left is the input point cloud, the bottom-left is the ground-truth, the top-left is normal map of ground-truth. The rest are the results of different methods, where the top row shows normal map, the middle row shows the geometry with normal rendering, while the bottom row shows vertex error E_v in color ($0mm$ $> 50mm$).

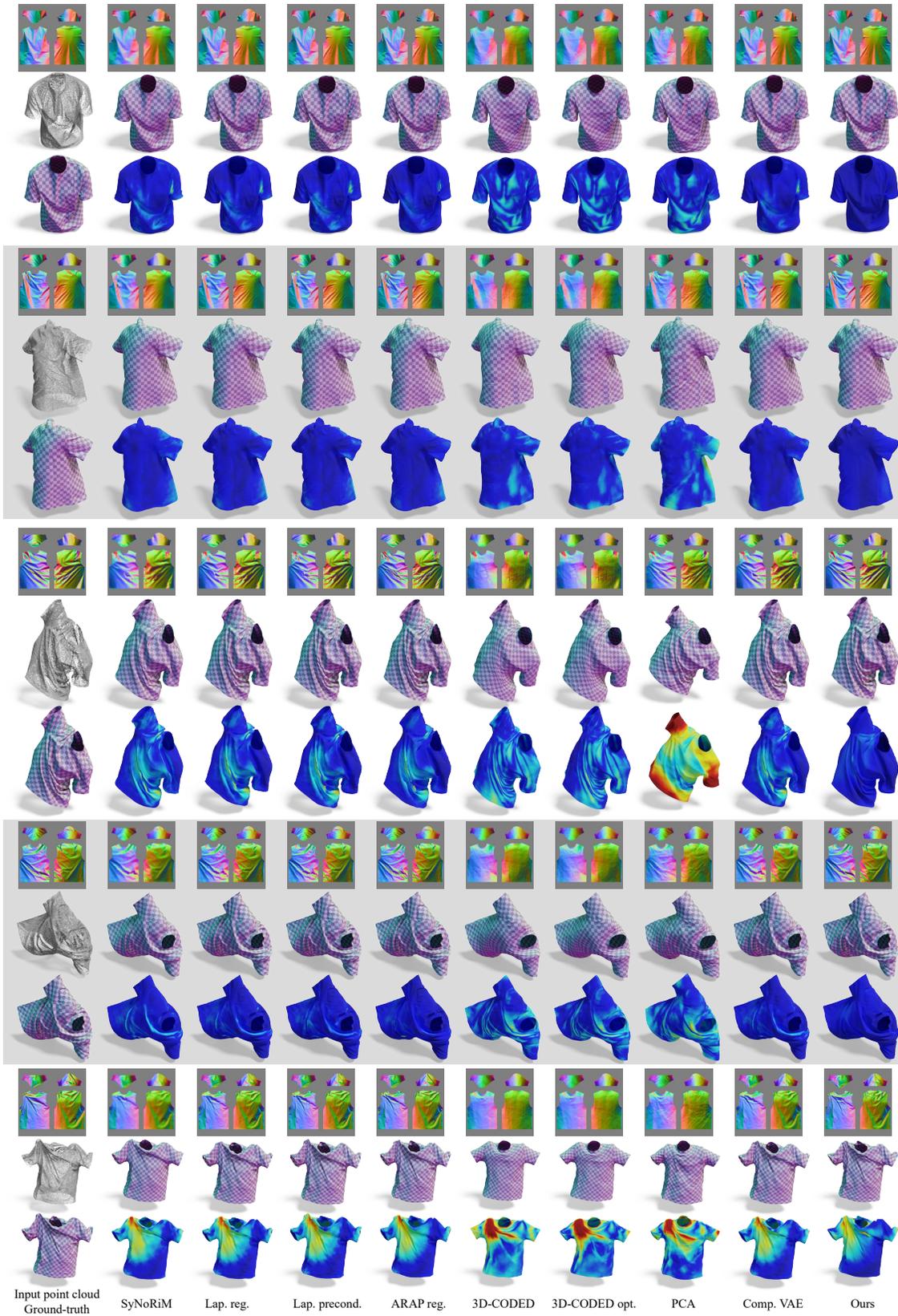


Figure 17. Comparison to baseline methods on *T-shirt 2*. See Figure 16 for explanation.

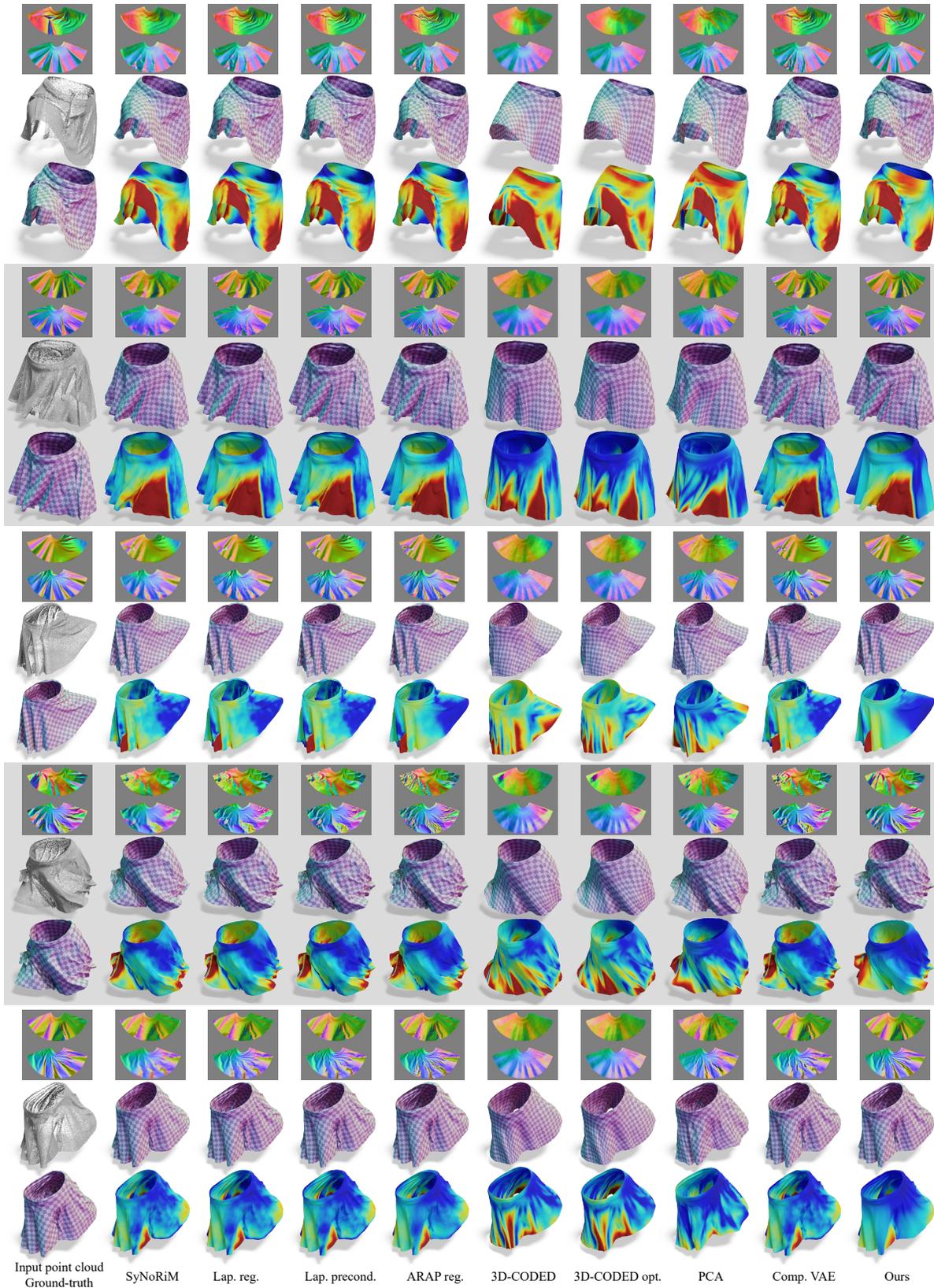


Figure 18. Comparison to baseline methods on *Skirt 1*. See Figure 16 for explanation.

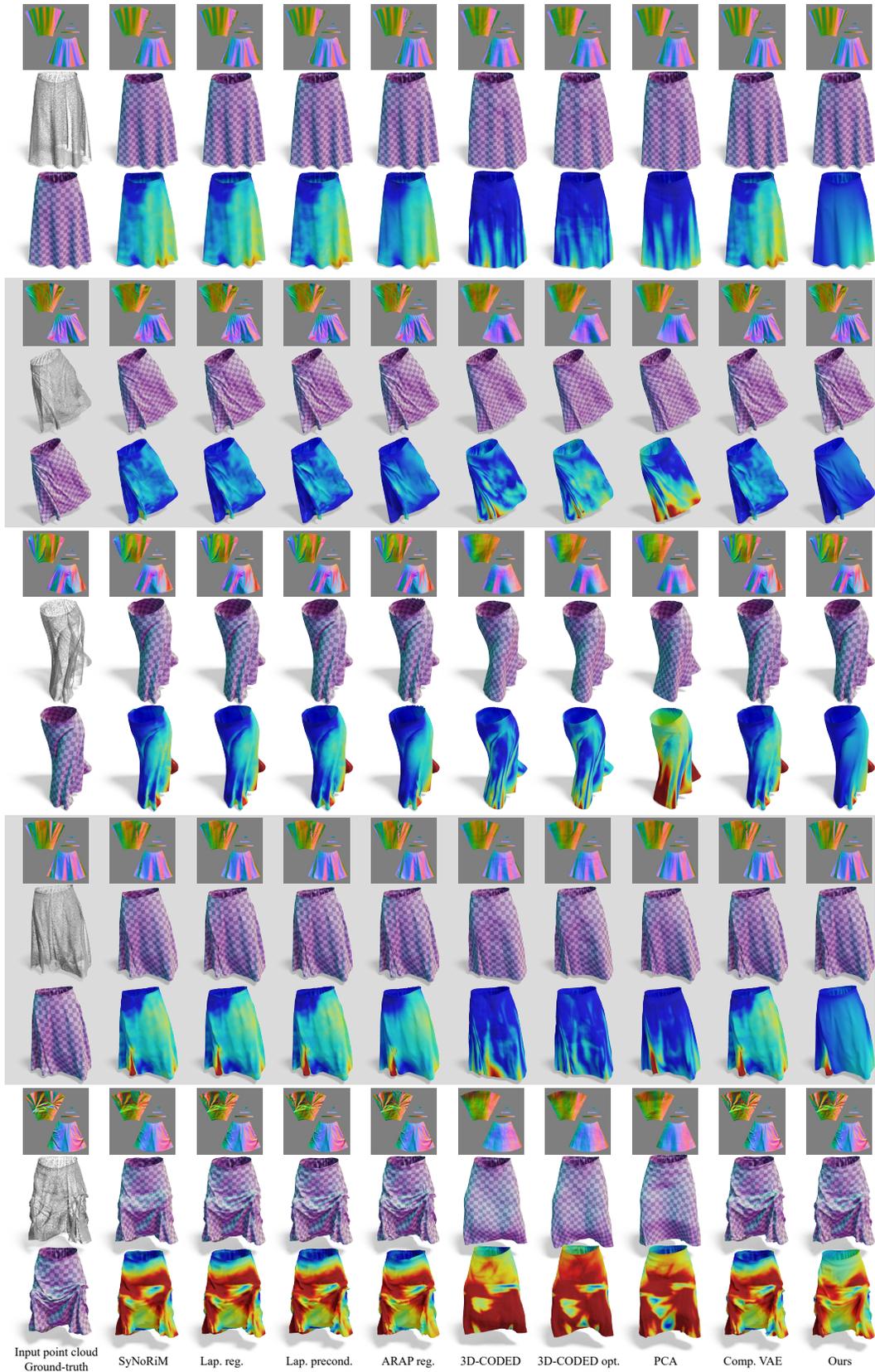


Figure 19. Comparison to baseline methods on *Skirt 2*. See Figure 16 for explanation.