

# UNSUPERVISED GEOMETRIC DISENTANGLEMENT FOR SURFACES VIA CFAN-VAE

**N. Joseph Tatro**

Department of Mathematical Sciences  
Rensselaer Polytechnic Institute  
Troy, NY 12180  
tatron@rpi.edu

**Stefan C. Schonsheck**

Department of Applied Mathematics  
University of California Davis  
Davis, CA 95616  
scschonsheck@ucdavis.edu

**Rongjie Lai**

Department of Mathematical Sciences  
Rensselaer Polytechnic Institute  
Troy, NY 12180  
lair@rpi.edu

## ABSTRACT

Geometric disentanglement, the separation of latent codes for intrinsic (i.e. identity) and extrinsic (i.e. pose) geometry, is a prominent task for generative models of non-Euclidean data such as 3D deformable models. It provides greater interpretability of the latent space, and leads to more control in generation. This work introduces a mesh feature, the conformal factor and normal feature (CFAN), for use in mesh convolutional autoencoders. We further propose CFAN-VAE, a novel architecture that disentangles identity and pose using the CFAN feature. Requiring no label information on the identity or pose during training, CFAN-VAE achieves geometric disentanglement in an unsupervised way. Our comprehensive experiments, including reconstruction, interpolation, generation, and identity/pose transfer, demonstrate CFAN-VAE achieves state-of-the-art performance on unsupervised geometric disentanglement.

## 1 INTRODUCTION

With the recent popularity of geometric deep learning (Bronstein et al., 2017), mesh based convolutional autoencoders (MeshVAEs) are now a popular tool for surface generation in the surface processing community (Cheng et al., 2019; Litany et al., 2018; Ma et al., 2019; Ranjan et al., 2018). With these VAEs achieving state-of-the-art performance on tasks such as reconstruction, more attention is being given towards tasks such as latent space interpretability. Geometric disentanglement, where the latent variables controlling *intrinsic* (properties independent of surface embedding) and *extrinsic* (properties dependent on surface embedding) geometry are separated (Aumentado-Armstrong et al., 2019), is an important open problem related to such interpretability. A typical application of a disentangled latent space is the separation of identity and pose in the case of human body generation (Jiang et al., 2019a; Tan et al., 2018). This is useful in fields such as graphics where we may be interested in the procedural generation of certain assets.

We are interested in creating an architecture that explicitly leads to geometric disentanglement in an unsupervised setting. Namely, we do not require labels for mesh identity and pose, allowing the architecture to be applied more broadly. In this work, we:

1. introduce a feature, the conformal factor and normal feature (CFAN), that decouples intrinsic and extrinsic geometry for use in mesh convolutional autoencoders.
2. propose a novel architecture, CFAN-VAE, for unsupervised geometric disentanglement. For a given mesh, we compute the CFAN feature, and encode its components separately into latent vectors representing intrinsic and extrinsic geometry.

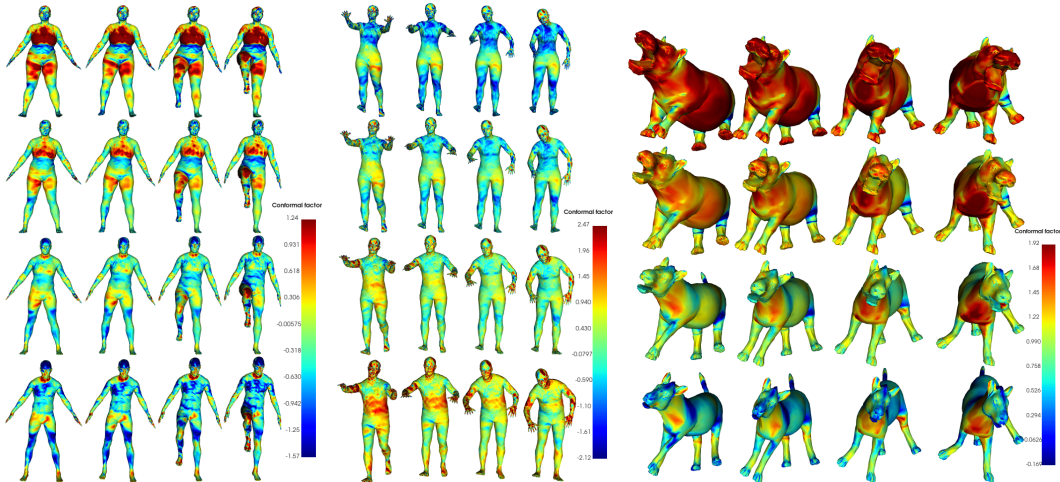


Figure 1: Geometric disentangled interpolations between two meshes from DFAUST, SURREAL, and SMAL datasets. Meshes are generated by CFAN-VAEs. The horizontal/vertical axis display linear interpolation in the normal/conformal latent codes,  $z_n/z_c$ . Color denotes the conformal factor feature of each reconstructed mesh.

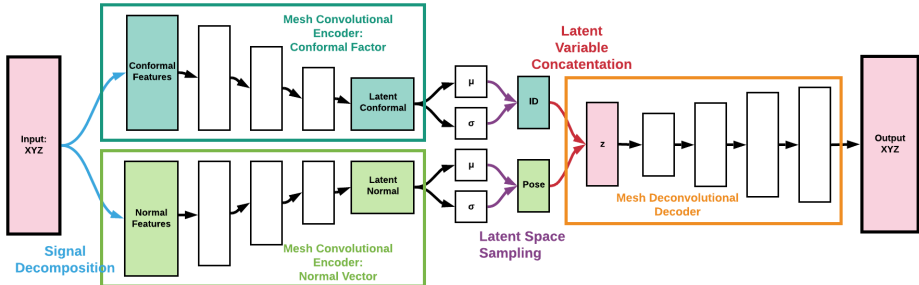


Figure 2: Diagram of the CFAN-VAE architecture. 3D mesh coordinates are transformed to the CFAN feature and separately encoded to promote disentanglement.

**Related Work** Geometric disentanglement is an important feature leading to a more interpretable latent space and thus more control in generation. For instance, (Jiang et al., 2019a;b; Tan et al., 2018) learn a latent space separating identity and pose, advocating for the use of the *as-consistent-as-possible deformation representation* (ACAPDR). These works require a reference pose for all identities during training. Recently, several works proposed networks using non-Euclidean convolution for this task (Cosmo et al., 2020; Levinson et al., 2019; Zhou et al., 2020). These networks either utilize supervision or loss functions that require identification of isometric pairs in training data. In practice, the latter needs supervised identity to scale on large datasets. Training datasets lacking isometries may be desirable in applications such as differential privacy (Abadi et al., 2016).

There is limited work to achieve unsupervised geometric disentanglement that *completely* eliminates reliance on both identity and pose labeling. To our knowledge, the only previous work concerning unsupervised disentanglement for 3D data is (Aumentado-Armstrong et al., 2019), in which the authors introduce GDVAE, a VAE built on a PointNet (Charles et al., 2017). This network requires the computation of Laplace-Beltrami eigenvalues on each surface of the training set. We compare our results to this work as it is the state-of-the-art in unsupervised geometric disentanglement. Later, we will find advantages of our approach include the use of a signal, unlike Laplace-Beltrami eigenvalues, that is not globally sensitive to local deformation. Additionally, as our model is not trained on PointNet embeddings, we are able to use significantly less parameters (on the order of millions), for our architecture.

## 2 CFAN-VAE

Our novel method for geometric disentanglement is motivated by the Fundamental Theorem of Surfaces (Do Carmo, 2016). We describe surfaces using conformal factors and surface normal vectors. This CFAN feature is easily leveraged by our introduced CFAN-VAE for geometric disentanglement.

**CFAN Feature** The Fundamental Theorem of Surfaces states that a surface can be uniquely reconstructed up to rigid motion given feasible metric tensor and surface normals (Do Carmo, 2016). Then we can use metric tensors and normal vector fields to characterize surfaces. In this work, we consider genus zero surfaces. It is well known that all genus zero surfaces are conformally equivalent (Jost & Jost, 2008). Namely, given two genus-zero surfaces,  $(\mathcal{M}_1, g_1)$  and  $(\mathcal{M}_2, g_2)$ , there exists a diffeomorphism,  $\phi : (\mathcal{M}_1, g_1) \rightarrow (\mathcal{M}_2, g_2)$  where  $\phi^*(g_2) = \exp(\lambda)g_1$ .

Here the function,  $\lambda$ , is known as the conformal factor and defines a conformal deformation from  $\mathcal{M}_1$  to  $\mathcal{M}_2$ . Then any pair of conformally equivalent surfaces can be deformed into one another (up to isomorphism) by choosing the correct conformal factor. The surface embedding is defined up to translation by the normal field. These two components define the CFAN feature, with the conformal factor and normal field representing intrinsic and extrinsic geometry respectively.

Given a triangle mesh  $(\mathbf{P}, \mathbf{T})$ , where  $\mathbf{P} \in \mathbb{R}^{n \times 3}$  is the set of vertices and  $\mathbf{T} \in \mathbb{R}^{k \times 3}$  is the corresponding set of faces, we define the discretized CFAN feature. Each face  $\tau$  in  $\mathbf{T}$  has a corresponding exterior face normal  $\mathbf{n}_\tau$ . We compute the weighted average of these face normals around the first ring structure of each vertex to define a pointwise normal. In the continuous setting, the local area on the surface is given by  $\sqrt{\det g_2}$ , so it follows that the logarithm of local area is an affine function of the conformal factor,  $\lambda + \frac{1}{2} \log \det g_1$ . We have equivalence if  $g_1$  is taken to be the trivial metric. With this in mind, we approximate a discrete conformal factor by taking the logarithm of a third of the area of the first ring structure about each point. Formally we define the CFAN feature,  $(c_i, \mathbf{n}_i)$  as:

$$c_i = \log \sum_{\substack{\tau \in \mathbf{T}; \\ i \in \tau}} \frac{\text{Area}(\tau)}{3}, \quad \mathbf{n}_i = \frac{\sum_{\substack{\tau \in \mathbf{T}; \\ i \in \tau}} \text{Area}(\tau) \mathbf{n}_\tau}{\|\sum_{\substack{\tau \in \mathbf{T}; \\ i \in \tau}} \text{Area}(\tau)\|} \quad (1)$$

**Network Architecture** Based on the CFAN feature, we propose a simple architecture, CFAN-VAE, as shown in Figure 2 to achieve unsupervised geometric disentanglement. The intuition is to encode the conformal factor and the normal features separately. The 3D coordinates  $\mathbf{p}$  of an input mesh under a fixed triangulation are first converted to the CFAN feature by transformations  $\phi_c$  and  $\phi_n$  provided by equation 1, yielding  $\mathbf{c} = \phi_c(\mathbf{p})$  and  $\mathbf{n} = \phi_n(\mathbf{p})$ .

Then, the feature components are separately encoded by  $E_c$  and  $E_n$  to create two different latent variables, the conformal latent variable  $\mathbf{z}_c$  and the normal latent variable  $\mathbf{z}_n$  in the disentangled latent space  $\mathbf{Z}_{c,n} = \mathbf{Z}_c \times \mathbf{Z}_n$ . That is,  $\mathbf{z}_c = E_c(\mathbf{c})$ ,  $\mathbf{z}_n = E_n(\mathbf{n})$ , and  $\hat{\mathbf{p}} = D(\mathbf{z}_c, \mathbf{z}_n)$ .

As a result,  $\mathbf{z}_c$  corresponds to intrinsic geometry, controlling surface identity, and  $\mathbf{z}_n$  corresponds to extrinsic geometry, controlling surface pose. After that, the CFAN latent variable  $\mathbf{z}_{c,n} = [\mathbf{z}_c, \mathbf{z}_n] \in \mathbf{Z}_{c,n}$  is decoded by  $D$  to obtain the reconstructed 3D coordinates  $\hat{\mathbf{p}}$ .

Our model is built on geometric convolutions given by PTC layers (Schonsheck et al., 2018). The mesh signal is decimated during encoding and refined during decoding by a factor of four. Sampling is precomputed by the geodesic Farthest Point Sampling (FPS) method (Moenning & Dodgson, 2003). After the final mesh convolution in the encoder, dense layers map the signal to the variational statistics,  $\mu := [\mu_c, \mu_n]$  and  $\sigma := [\sigma_c, \sigma_n]$ . A dense layer also maps  $\mathbf{z}_{c,n}$  to the first mesh signal in the decoder. Details on the loss function used for training can be found in Section A.1.

## 3 EXPERIMENTS AND RESULTS

We consider three datasets in our experiments. This includes DFAUST, a real-world motion-capture dataset of 10 people performing several poses (Bogo et al., 2017). We randomly split the dataset into a training/validation/test set of 35,720/500/5,000 meshes respectively. In addition, we generated a synthetic dataset of animals referred to as SMAL following (Zuffi et al., 2017), based on

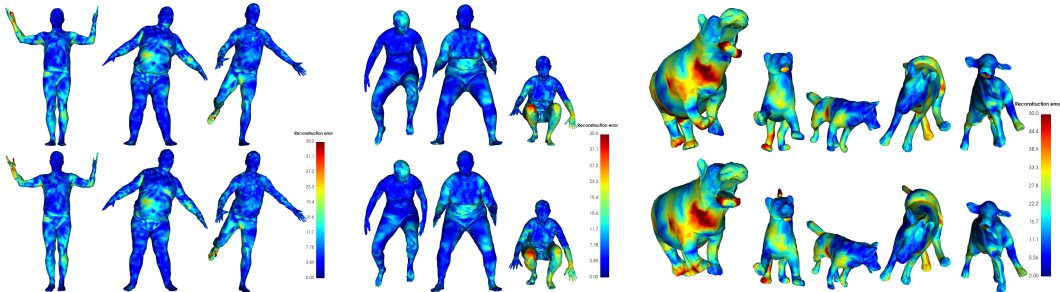


Figure 3: Reconstructions from the CFAN-VAE models in Figure 1. The mesh color represents the pointwise Euclidean error after translation. The top row are the ground truth meshes, and the bottom row are the reconstructions.

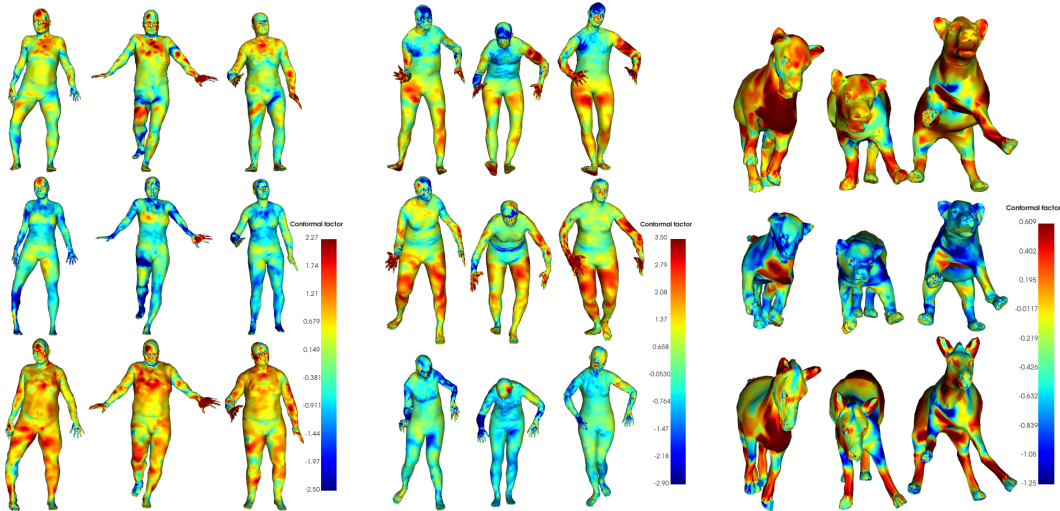


Figure 4: Surfaces generated by sampling the CFAN-VAE latent space. Horizontal and vertical axes correspond to  $z_n$  and  $z_c$  respectively. CFAN-VAE reliably generates meshes in a specific pose or with nearly the same identity, showing geometric disentanglement. Color denotes conformal factor.

the SMPL model introduced in (Loper et al., 2015). Similarly, we formed a synthetic dataset of humans referred to as SURREAL from (Varol et al., 2017). These datasets are created following (Aumentado-Armstrong et al., 2019). For SURREAL, we add small normal noise with  $\sigma = 0.2$  to the intrinsic shape parameters  $\beta$  to eliminate any isometry in the dataset. These two datasets provide challenging environments where identity-supervised methods fail.

Details on the network and training hyperparameters are provided in Section A.2. For comparisons to GDVAE, we use the available dataset-specific pretrained models available in (Aumentado-Armstrong et al., 2019). We verify that CFAN-VAE is able to accurately reconstruct surfaces. Figure 3 displays high quality reconstructions for each dataset using CFAN-VAEs for DFAUST, SURREAL, and SMAL.

**Surface Interpolation** Figure 1 displays examples of disentangled interpolations for each dataset using CFAN-VAE. In this figure, the vertical and horizontal axis represents conformal and normal interpolations respectively. The color map denotes the pointwise normalized conformal factor of the reconstructions. These results verify that CFAN-VAE produces meaningful geometric disentanglement between intrinsic and extrinsic information. Thus, it enjoys flexibility to generate high quality meshes preserving either identity or pose. Notice, the colormaps and position are almost identical for fixed  $z_c$  and  $z_n$  respectively. We include additional examples in Figure 5 in the appendix.

We compare the performance of CFAN-VAE to GDVAE on transfer tasks. For SURREAL and SMAL, we can create ground-truth mesh solutions to the transfer task by exchanging the appropriate

Table 1: Chamfer distance of the mesh generated in pose and identity transfer tasks to the target mesh. The generated mesh is decoded from swapping the latent codes of the source meshes,  $(z_c^{(1)}, z_n^{(1)})$  and  $(z_c^{(2)}, z_n^{(2)})$ . This distance is the mean over 256 transfer pairs with standard error.

Networks	SURREAL		SMAL	
	Error(mm)	#Param.(M)	Error(mm)	#Param.(M)
Distance between sources	$76.0 \pm 1.8$	—	$313.9 \pm 9.7$	—
GDVAE <sup>1</sup>	$31.6 \pm 0.6$	21.30	$96.3 \pm 1.9$	21.30
CFAN-VAE	<b><math>26.8 \pm 0.8</math></b>	1.75	<b><math>35.5 \pm 0.8</math></b>	2.15

SMPL parameters of the source meshes. Then we compare the Chamfer distance of the network generated mesh, with appropriate latent code exchanged, to the SMPL generated ground truth. These results are summarized in Table 1. It is clear that CFAN-VAE outperforms GDVAE, particularly on the SMAL dataset. This nice property of the CFAN-VAE latent space enables flexible control of the geometry of generated surfaces. We stress that these two datasets lack isometric pairs, so other methods (Cosmo et al., 2020; Zhou et al., 2020) are not of use for comparison. Further analysis of the learned latent space is available in Section A.3.

**Surface Generation** Figure 4 plots random generations where latent codes are sampled from a multivariate Gaussian. This Gaussian,  $\mathcal{N}(\mu, \Sigma)$ , is a product of marginals fit to the embedding of the test set for each latent variable. We sample from  $\mathcal{N}(\mu, 0.8\Sigma)$ . Here the normal latent code and conformal code is fixed along the vertical axis and horizontal axis respectively. Clearly, pose and identity are successfully fixed. The generated meshes are of comparable visual quality to the reconstructions from the test set. Then the latent space of CFAE-VAE allows for high-quality generation while providing interpretability of geometric features.

## 4 CONCLUSION

We propose a novel architecture, CFAN-VAE, for unsupervised geometric disentanglement in mesh convolutional autoencoders. This is accomplished by utilizing the conformal factor and normal (CFAN) feature to separate intrinsic and extrinsic geometry. There is clearly strong geometric disentanglement in CFAN-VAEs. We achieve state-of-the-art performance on transfer tasks and disentanglement compared to prior work. The continued integration of geometric theory into neural network architectures can only lead to more interesting results in the future.

## ACKNOWLEDGEMENTS

J. Tatro and R. Lai’s work is supported in part by NSF CAREER Award (DMS—1752934). J. Tatro’s work is also supported in part by IBM-RPI AIRC program. Part of this research was performed while the authors were visiting the Institute for Pure and Applied Mathematics (IPAM), which is supported by the NSF DMS-1440415, during the *Geometry and Learning from Data in 3D and Beyond* long program.

## REFERENCES

- Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pp. 308–318, 2016.
- Tristan Aumentado-Armstrong, Stavros Tsogkas, Allan Jepson, and Sven Dickinson. Geometric disentanglement for generative latent shape models, 2019.
- Federica Bogo, Javier Romero, Gerard Pons-Moll, and Michael J. Black. Dynamic FAUST: Registering human bodies in motion. In *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, July 2017.

- Michael M. Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst. Geometric deep learning: Going beyond euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, Jul 2017. ISSN 1053-5888. doi: 10.1109/msp.2017.2693418. URL <http://dx.doi.org/10.1109/MSP.2017.2693418>.
- R. Qi Charles, Hao Su, Mo Kaichun, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jul 2017. doi: 10.1109/cvpr.2017.16. URL <http://dx.doi.org/10.1109/cvpr.2017.16>.
- Shiyang Cheng, Michael Bronstein, Yuxiang Zhou, Irene Kotsia, Maja Pantic, and Stefanos Zafeiriou. Meshgan: Non-linear 3d morphable models of faces, 2019.
- Luca Cosmo, Antonio Norelli, Oshri Halimi, Ron Kimmel, and Emanuele Rodolà. Limp: Learning latent shape representations with metric preservation priors, 2020.
- Manfredo P Do Carmo. *Differential Geometry of Curves and Surfaces: Revised and Updated Second Edition*. Courier Dover Publications, 2016.
- Boyi Jiang, Juyong Zhang, Jianfei Cai, and Jianmin Zheng. Learning 3d human body embedding, 2019a.
- Zi-Hang Jiang, Qianyi Wu, Keyu Chen, and Juyong Zhang. Disentangled representation learning for 3d face shape, 2019b.
- Jürgen Jost and Jèurgen Jost. *Riemannian geometry and geometric analysis*, volume 42005. Springer, 2008.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Jake Levinson, Avneesh Sud, and Ameesh Makadia. Latent feature disentanglement for 3d meshes, 2019.
- Or Litany, Alex Bronstein, Michael Bronstein, and Ameesh Makadia. Deformable shape completion with graph convolutional autoencoders. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Jun 2018. doi: 10.1109/cvpr.2018.00202. URL <http://dx.doi.org/10.1109/CVPR.2018.00202>.
- Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. Smpl: a skinned multi-person linear model. *ACM Trans. Graph.*, 34:248:1–248:16, 2015.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Qianli Ma, Siyu Tang, Sergi Pujades, Gerard Pons-Moll, Anurag Ranjan, and Michael J. Black. Dressing 3d humans using a conditional mesh-vae-gan. *ArXiv*, abs/1907.13615, 2019.
- Carsten Moenning and Neil A Dodgson. Fast marching farthest point sampling. Technical report, University of Cambridge, Computer Laboratory, 2003.
- Anurag Ranjan, Timo Bolkart, Soubhik Sanyal, and Michael J. Black. Generating 3d faces using convolutional mesh autoencoders. *Lecture Notes in Computer Science*, pp. 725–741, 2018. ISSN 1611-3349.
- Stefan C. Schonsheck, Bin Dong, and Rongjie Lai. Parallel transport convolution: A new tool for convolutional neural networks on manifolds, 2018.
- Qingyang Tan, Lin Gao, Yu-Kun Lai, and Shihong Xia. Variational autoencoders for deforming 3d mesh models. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Jun 2018. doi: 10.1109/cvpr.2018.00612. URL <http://dx.doi.org/10.1109/cvpr.2018.00612>.
- Gül Varol, Javier Romero, Xavier Martin, Naureen Mahmood, Michael J. Black, Ivan Laptev, and Cordelia Schmid. Learning from synthetic humans. In *CVPR*, 2017.

Keyang Zhou, Bharat Lal Bhatnagar, and Gerard Pons-Moll. Unsupervised shape and pose disentanglement for 3d meshes. In *16th European Conference on Computer Vision*. Springer, 2020.

Silvia Zuffi, Angjoo Kanazawa, David Jacobs, and Michael J. Black. 3D menagerie: Modeling the 3D shape and pose of animals. In *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, July 2017.

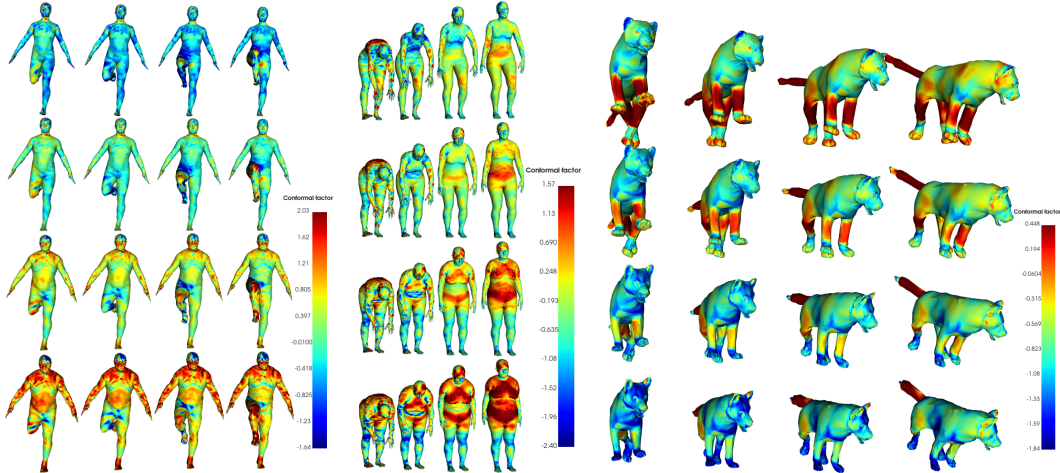


Figure 5: Additional geometrically disentangled interpolations using CFAN-VAE as in Figure 1.

Table 2: Disentanglement metrics on the learned latent spaces. For each mesh embedding, we find its nearest neighbor in either the intrinsic or extrinsic space,  $z_c / z_n$ . We report the norm of the difference in the identity and pose SMPL parameters,  $\beta$  and  $\theta$ , of the embedded mesh and its neighbor. We report mean distance with standard error on the test set.

Networks	Param.	Desired		SURREAL		SMAL	
		Magnitude		$z_c$	$z_n$	$z_c$	$z_n$
GDVAE <sup>2</sup>	$\beta$	↓	↑	$4.42 \pm 0.03$	$4.44 \pm 0.03$	$2.40 \pm 0.02$	$3.48 \pm 0.02$
	$\theta$	↑	↓	$4.21 \pm 0.04$	$4.11 \pm 0.04$	$2.72 \pm 0.01$	$2.71 \pm 0.01$
CFAN-VAE	$\beta$			<b><math>4.05 \pm 0.02</math></b>	$4.38 \pm 0.02$	<b><math>1.73 \pm 0.01</math></b>	$3.73 \pm 0.02$
	$\theta$			$3.69 \pm 0.03$	<b><math>2.81 \pm 0.04</math></b>	$2.74 \pm 0.01$	<b><math>2.63 \pm 0.01</math></b>

## A APPENDIX

### A.1 LOSS FUNCTION

The training loss for our model is given by:

$$\mathcal{L} := \|\mathbf{P} - \hat{\mathbf{P}}\|_1 + \sum_{\mu^{(0)}, \mu^{(1)} \in \mu} \left( \lambda_D \mathcal{L}_D(\mu^{(0)}, \mu^{(1)}) + \lambda_M \mathcal{L}_M(\mu^{(0)}, \mu^{(1)}) \right) + \lambda_{KL} \|\sigma^2 + \mu^2 - 1 - \log(\sigma^2)\|_1. \quad (2)$$

The first term is simply the  $L_1$  error of the reconstruction which promotes robustness to vertex outliers.  $\mathcal{L}_D$  and  $\mathcal{L}_M$  are disentanglement and metric penalties that we define shortly. The fourth term of the loss function is the KL-divergence of the latent representation from the unit normal distribution, given by the Bayesian prior assumption (Kingma & Welling, 2013).

In CFAN-VAE, geometric disentanglement is a result of separately encoding intrinsic and extrinsic geometric information. We stress that structural conditions prevent guaranteeing complete independence of the separate latent vectors. This is intuitive, considering motion can be restricted by intrinsic properties such as body mass distribution.

To promote independence, we consider a pair of encoded meshes. We create a latent variable with the conformal latent code of the first mesh, while the normal code is a random convex combination of the two normal codes. We decode and re-encode this new latent variable. We then have an  $L_2$  penalty on the resulting change in the conformal code upon re-encoding as it should not be affected by this change in the normal code. We penalize a change in the normal code analogously. This



disentanglement penalty,  $\mathcal{L}_D$ , is formalized below,

$$\mathcal{L}_D := \mathbb{E} \sum_{\substack{i,j \text{ in} \\ \{c,n\}}} \|\mu_i^{(0)} - E_i \left( D(\mu_c^{(0)} + \delta_{j,c}\epsilon_c, \mu_n^{(0)} + \delta_{j,n}\epsilon_n) \right)\|_2^2, \quad (3)$$

$$\epsilon_k := \alpha(\mu_k^{(1)} - \mu_k^{(0)}), \quad \alpha \in U[0, 1]. \quad (4)$$

Inspired by the metric regularization introduced in (Cosmo et al., 2020), we add a metric penalty,  $\mathcal{L}_M$ , to the loss function,

$$\mathcal{L}_M := \mathbb{E}_\alpha \|l_\alpha(\phi_c \circ D(\mu^{(0)}), \phi_c \circ D(\mu^{(0)})) - \phi_c \circ D(l_\alpha(\mu^{(0)}, \mu^{(1)}))\|_1, \quad \alpha \in U[0, 1], \quad (5)$$

where  $l_\alpha$  is the  $\alpha$  convex combination of its entries. This promotes smoothness in the metric of reconstructions.

## A.2 NETWORK HYPERPARAMETERS

Each VAE contains 5 layer encoders and decoder. For the PTC kernels, we use a 13 point stencil with the support of the kernel about a point being its 9 nearest geodesic neighbors. The size of the conformal/normal latent vectors are (32, 32) for DFAUST and (8, 32) for SMAL. For visualization, we use models with a latent dimension of (32, 32). For corresponding xyz-VAEs, the latent dimension is the sum of the two latent dimensions.

The latent dimension of the GDVAE models are (5, 15) and (5, 11) for SURREAL and SMAL respectively. When comparing to CFAN-VAE models for SURREAL, we choose a corresponding latent dimension of (5, 15).

We use batch normalization (BN) layers in the encoder. When training with the disentanglement penalty  $\mathcal{L}_D$ , we freeze the BN layers during reencoding. The activation functions are ReLU and ELU for the encoders and decoder. We use the AdamW optimizer (Loshchilov & Hutter, 2017) with a learning rate of 1E-3 and a weight decay of 5E-5. The value of  $\lambda_{KL}$  is set at 1E-4, while  $\lambda_D$  and  $\lambda_M$  are 5E-2. Each model is trained for 300 epochs with a batch size of 32. We train two instances of each network using different random seeds.

## A.3 LATENT SPACE ANALYSIS

To better assess the geometric disentanglement in our networks, we analyze the latent spaces of the CFAN-VAEs. We quantitatively measure the disentanglement in latent space in the manner of (Zhou et al., 2020). First, we embed the test sets for SURREAL and SMAL in the latent space. Given a mesh embedding, we find its nearest neighbor in either the intrinsic or extrinsic latent spaces,  $Z_c$  or  $Z_n$ . As these meshes were created via SMPL shape and pose parameters,  $\beta$  and  $\theta$ , we measure the distance between these parameters of the embedded mesh and its neighbor. These distances are summarized in Table 2. If the network latent spaces truly possess disentanglement, we expect the  $z_c$  neighbor to be closer in  $\beta$  distance and the  $z_n$  neighbor to be closer in  $\theta$  distance.