

Deep Deformable 3D Caricatures with Learned Shape Control

Yucheol Jung
POSTECH

Pohang, Republic of Korea
ycjung@postech.ac.kr

Jiaolong Yang
Microsoft Research Asia
Beijing, China
jiaoyan@microsoft.com

Wonjong Jang
POSTECH

Pohang, Republic of Korea
wonjong@postech.ac.kr

Xin Tong
Microsoft Research Asia
Beijing, China
xtong@microsoft.com

Soongjin Kim
POSTECH

Pohang, Republic of Korea
kimsj0302@postech.ac.kr

Seungyong Lee
POSTECH

Pohang, Republic of Korea
leesy@postech.ac.kr

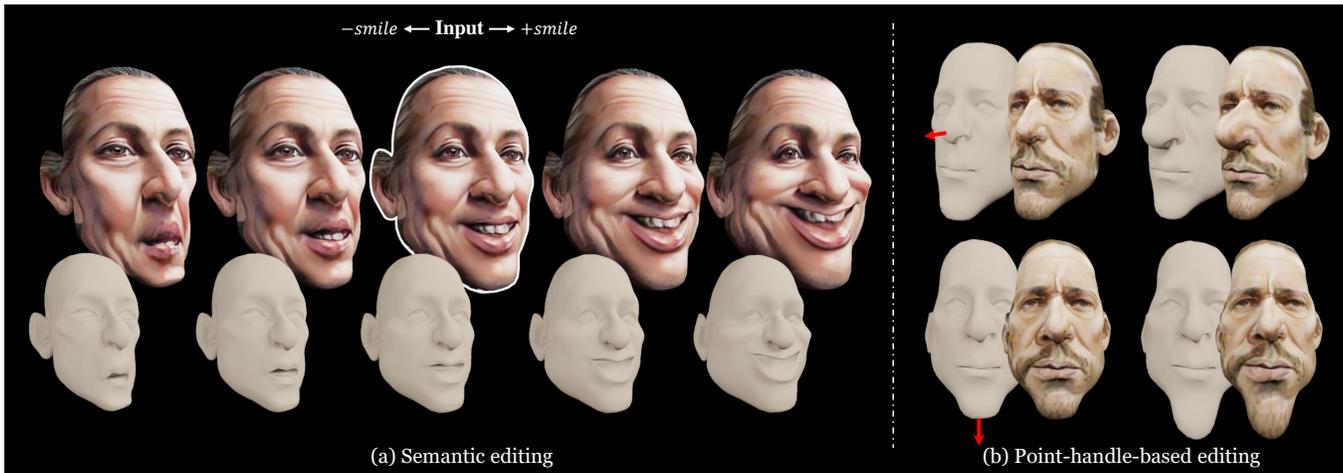


Figure 1: Our deformable model provides a data-driven editing space for 3D caricature shapes. (a) Semantic editing of 3D caricatures with varying degrees. (b) Point-handle-based deformation using learned latent space. The edited 3D caricatures have been reconstructed from 2D caricatures generated by StyleCariGAN [Jang et al. 2021]. 3D caricatures can also be reconstructed from real-world 2D caricatures, e.g., created by artists.

ABSTRACT

A 3D caricature is an exaggerated 3D depiction of a human face. The goal of this paper is to model the variations of 3D caricatures in a compact parameter space so that we can provide a useful data-driven toolkit for handling 3D caricature deformations. To achieve the goal, we propose an MLP-based framework for building a deformable surface model, which takes a latent code and produces a 3D surface. In the framework, a SIREN MLP models a function that takes a 3D position on a fixed template surface and returns a 3D displacement vector for the input position. We create variations of 3D surfaces by learning a hypernetwork that takes a latent code and produces the parameters of the MLP. Once

learned, our deformable model provides a nice editing space for 3D caricatures, supporting label-based semantic editing and point-handle-based deformation, both of which produce highly exaggerated and natural 3D caricature shapes. We also demonstrate other applications of our deformable model, such as automatic 3D caricature creation. Our code and supplementary materials are available at <https://github.com/ycjungSubhuman/DeepDeformable3DCaricatures>.

CCS CONCEPTS

• **Computing methodologies** → **Mesh models; Neural networks.**

KEYWORDS

Deformable model, Parametric model, 3D face model, Semantic 3D face control, 3D face deformation, Auto-decoder

ACM Reference Format:

Yucheol Jung, Wonjong Jang, Soongjin Kim, Jiaolong Yang, Xin Tong, and Seungyong Lee. 2022. Deep Deformable 3D Caricatures with Learned Shape Control. In *Special Interest Group on Computer Graphics and Interactive Techniques Conference Proceedings (SIGGRAPH '22 Conference Proceedings)*, August 7–11, 2022, Vancouver, BC, Canada. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3528233.3530748>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGGRAPH '22 Conference Proceedings, August 7–11, 2022, Vancouver, BC, Canada
© 2022 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-9337-9/22/08...\$15.00
<https://doi.org/10.1145/3528233.3530748>

1 INTRODUCTION

Deformable models have been widely used in graphics for shape creation and control. They provide compact parameter spaces and *a priori* information for the parameters, enabling effective manipulation of shapes. Human face is one of the domains that benefits from deformable models. One of the most widely used formulation for human face is the 3D Morphable Model (3DMM) [Blanz et al. 1999]. 3DMM and its variants have been adopted for challenging tasks, including 3D face reconstruction from 2D images [Deng et al. 2019; Garrido et al. 2016], semantic 3D face editing [Ghafourzadeh et al. 2020], and face reenactment [Thies et al. 2016].

A 3D caricature is an exaggerated 3D depiction of a human face. It has broad applications ranging from professional cartoon film production to low-cost avatar creation for games, social media, and AR/VR. Conventionally, 2D caricatures are created by skilled artists, and crafting 3D caricatures would require even more expertise. On the other hand, similarly to the case of 3D faces, a deformable 3D caricature model can provide a useful toolkit for handling exaggerated 3D faces by modeling the variations of 3D caricatures (Fig. 1).

3DMM has been used for effectively modeling variations of regular faces with a linear parameter space for 3D vertex positions. When applied to 3D caricatures, 3DMM results in large reconstruction errors as shown in a previous work [Wu et al. 2018]. To reduce the reconstruction errors, Wu *et al.* [2018] proposed an explicit non-linear mapping from shape parameters to 3D vertex positions for extrapolating regular faces to 3D caricatures. However, they did not exploit at all the editing capability of the deformable model defined by the mapping. Recently released 3DCaricShop dataset [Qiu et al. 2021] provides a set of high-quality 3D caricature meshes sculpted by 3D artists. Such a dataset enables learning a latent space together with the mapping to 3D caricatures automatically using a neural network. Still, learning from the dataset is challenging as 3D caricatures have highly diverse styles, yet the number of examples is limited to around 2K. We provide an analysis on the complexity of the dataset in the supplementary material.

Our goal is to learn a deformable model from 3DCaricShop dataset [Qiu et al. 2021] to create a toolkit for controllable 3D caricatures. With highly diverse and sparse data, we believe the network architecture plays a central role in the learning. We choose a multi-layer perceptron (MLP) for the network architecture, which has been successfully used for various applications, including shape reconstruction [Park et al. 2019; Saito et al. 2019] and neural rendering [Chan et al. 2021; Lombardi et al. 2018, 2019; Mildenhall et al. 2020].

Another important decision for the learning is the output representation of the MLP (Fig. 2). An array of vertex position is a straightforward approach when all dense correspondence is given in the dataset. However, we observed an MLP for generating vertex array suffers from slow convergence and high reconstruction error. This approach also limits the sampling positions on the surface only to vertices. Continuous SDF in a volume is widely used recently for general 3D shape modeling, but we observed an MLP for SDF tends to miss small details such as the shape of eyes. Consequently, our framework lets the MLP to model continuous deformation function

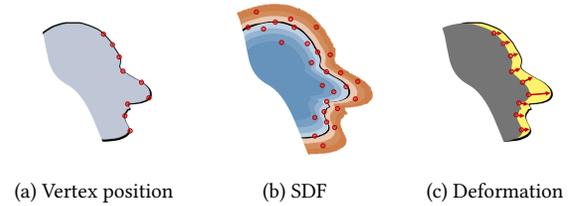


Figure 2: Different ways to represent a 3D shape on a network.

on a template surface. With this approach, we can generate 3D caricatures with important details.

To this end, we propose an MLP-based framework for building a deformable surface model, which takes a 128-dimension latent code and produces a 3D surface. In our framework, a SIREN MLP [Sitzmann et al. 2020] models a function that takes a 3D position on a fixed template surface and returns a 3D displacement vector for the input point. We create variations of 3D surface by learning a hyper-network that takes a latent code and produces the parameters of the MLP. The MLP and the hyper-network are learned jointly to build an auto-decoder [Park et al. 2019], which is used as our deformable 3D surface model.

Once learned, our deformable model provides a nice editing space for 3D caricatures. We demonstrate label-based semantic 3D caricature editing and point-handle-based 3D caricature editing, both of which produce highly exaggerated and natural 3D caricature shapes. We also demonstrate other applications, such as automatic 3D caricature creation from a 3D or 2D regular face and fitting to 2D landmarks to reconstruct 3D geometry out of a 2D caricature.

The contributions of this paper can be summarized as follows;

- We propose an MLP-based framework for building a deformable model that captures the variations of 3D caricatures in a compact parameter space.
- Our deformable model provides a useful data-driven toolkit for handling 3D caricatures, supporting label-based semantic editing and point-handle-based deformation as well as other applications.

2 RELATED WORK

2.1 Deformable face model

3DMM [Blanz et al. 1999] introduces a linear subspace for shape and texture using PCA decomposition from a 3D face database. Wu *et al.* [2018] showed 3DMM is not suitable for defining a subspace for caricatures since the space is not large enough to model large variations in caricatures. Different formulations on the 3DMM have been studied [Lüthi et al. 2017; Ploumpis et al. 2020], but their implication on the 3D caricature domain has not been explored.

Deep learning-based methods for modeling 3D facial mesh has been studied recently. PRNet [Feng et al. 2018] uses a 2D image to represent 3D facial geometry. Lombardi *et al.* [2018] builds a variational auto-encoder of 3D face meshes using fully connected layers. CoMA [Ranjan et al. 2018] proposes an auto-encoder model using graph convolution together with specialized down-sampling and up-sampling operators. MeshGAN [Cheng et al. 2019] trains

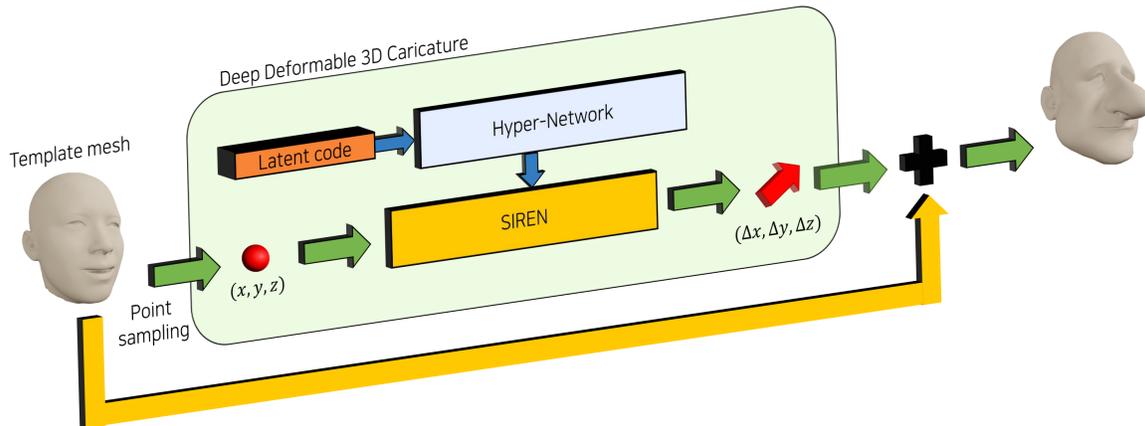


Figure 3: Overall framework for our deep deformable 3D caricature model.

two CoMA-based mesh decoders for facial identity and expression. Jiang *et al.* [2019] disentangles 3D facial meshes into identity and expression, both of which are encoded by graph-convolution-based mesh auto-encoders. In contrast to these methods generating fixed set of discrete vertices, we propose to model shape variations as a continuous deformation function using MLP to build a deformable 3D caricature model.

2.2 Deformable caricature model

One approach for building a deformable 3D caricature model is to expand regular face dataset with computational exaggeration. DeepSketch2Face [Han *et al.* 2017] extends regular face dataset by exaggerating faces synthetically using [Sela *et al.* 2015] to build a bi-linear model of 3D caricatures. CaricatureShop [Han *et al.* 2018] uses synthetic 3D caricatures as training data. DeepSketch2Face and CaricatureShop only focuses on sketch-based 3D caricature editing. Alive Caricature [Wu *et al.* 2018] models a space of 3D caricatures by building a carefully-designed blending operations for regular 3D faces. 3DCaricShop [Qiu *et al.* 2021] proposes a 3D caricature dataset sculpted by artists, and use the 3D caricatures to build a PCA model, which is used as an artifact-filtering tool. Alive Caricature and 3DCaricShop focuses on 3D mesh reconstruction from 2D caricatures through free-form vertex deformation using an optimization guided by parametric models. Although these solution produces 3D caricatures faithful to input images, editing on the result is non-trivial since the result is not fully represented in an editable parameter space. In contrast, our data-driven parametric model for 3D caricatures provides readily available controls on the result when applied to 3D caricature reconstruction.

Cai *et al.* [2021] learn caricature landmark detection and 3D caricature reconstruction using the results of Alive Caricature as training data. 3DMagicMirror [Guo *et al.* 2019] trains a graph-convolutional variational auto-encoder using regular 3D faces and the results of Alive Caricature. The latent space is used for mapping a regular 3D face to a 3D caricature. These methods focus on specific tasks of landmark detection and automatic caricature creation. We

provide a controllable parametric 3D caricature model that can be generally used for various sub-tasks including semantic editing.

2.3 Deep-learning-based 3D shape model

DeepSDF [Park *et al.* 2019] proposed using MLP to learn SDF to represent a 3D shape using an auto-decoder framework. Following DeepSDF, effective learning of SDFs has been researched. Implicit geometric regularization [Gropp *et al.* 2020] is widely used to regularize networks for SDF. DIF-Net [Deng *et al.* 2021] learns a template SDF and volumetric deformation function using MLPs, providing dense correspondence between generated SDFs. AtlasNet [Groueix *et al.* 2018] represents a 3D shape as multiple surface patches in 2D domain to model the surfaces of general 3D objects. In contrast, our method represents a 3D shape as a deformed template mesh.

3DN [Wang *et al.* 2019] reconstructs a 3D mesh from an input image by deforming a given fixed template mesh using MLPs. The function modeled by the MLP is similar to our framework; Both approaches model a continuous surface deformation function of a template surface. 3DN only discusses the effectiveness of the network in the context of 3D reconstruction via deformation of a given template. We use template deformation to effectively learn a controllable latent space of 3D shapes when dense correspondence between the shapes is given.

3 DEEP DEFORMABLE 3D CARICATURES

We build a data-driven deformable model for 3D caricatures using an auto-decoder framework [Park *et al.* 2019]. Once learned, the latent code can be sampled from distribution or optimized for representing an arbitrary input, e.g., unseen 3D caricature or sparse landmarks. Various controls can be achieved using the learned latent space. Fig. 3 shows our overall framework.

3.1 3D caricature representation

The key design choice in the training of the auto-decoder is the representation of a 3D caricature. While other choices may exist, we discuss three possible representations: vertex position array,

signed distance function (SDF), and surface deformation function, where we choose surface deformation function.

Vertex position array. Since the training data comes in the form of 3D mesh with uniform connectivity, one of the simplest approaches to model the data is to represent a shape as an array of per-vertex 3D position. Using this representation, we can generate a 3D caricature from a latent code by mapping the code to a large array of size $V \times 3$ using an MLP, where V is the number of vertices. We observe this architecture suffers from slow convergence and high reconstruction error (See Sec. 4.1).

SDF. DeepSDF [Park et al. 2019] and its variants have been widely used for modeling general shapes with large and diverse variations, even including topology changes. However, MLPs for SDF result in inefficiency when applied to our case of 3D caricature, where dense vertex correspondence between meshes is provided in the dataset. In terms of applications, it is not trivial to retrieve uniform meshing and correspondence from the generated SDFs. UV parameterization may also not be trivially shared between generated 3D caricatures. Also, the MLPs for SDF should produce valid function values outside of the surface, while our region of interest is only on the surface shape. When we trained DeepSDF with the 3D caricature dataset, we observed a loss of crucial details such as the shape of the eyes (See Sec. 4.1).

Surface deformation. We find modeling each shape as a deformation of a fixed template surface effective for learning the latent space of 3D caricatures, compared to using absolute position. We model deformation for each shape as a continuous function defined on the template surface. The continuous function is represented as an MLP and the variation in the function is modeled using the *hypernetwork* [Ha et al. 2017] framework, which has recently been used in modeling variations in 3D shapes using continuous functions [Deng et al. 2021; Liu et al. 2020; Sitzmann et al. 2019]. The hypernetwork architecture enables fast convergence, and the continuous representation enables finer sampling on the surface (See Sec. 4.3). For a more detailed analysis of the convergence of hypernetwork, refer to Sec. 4.2.

Compared to using SDFs, our representation fully utilize the correspondence information in the dataset. Our representation has three advantages over an MLP for SDF: 1) Initial meshing and the UV parameterization is shared trivially among all generated shapes. 2) Unlike SDF defined on volume, the template deformation function is only required to have valid values on the template surface, so the network capacity can focus on a smaller domain. 3) The template contains essential components of human face, so the network can focus only on the variations of shapes and skip the learning to produce overall shapes of facial components.

In conclusion, we represent a 3D caricature as a continuous displacement function defined on a fixed template surface. The function is modeled by an MLP that maps a 3D position sampled on the template surface to a 3D displacement vector to be applied to the input point. To generate a 3D caricature from a latent code, we map a latent code to the parameters of the MLP using a hypernetwork. The trained network can model diverse 3D caricatures with various expressions (Fig. 4).



Figure 4: Random 3D caricatures sampled from our deformable model.

3.2 Network structure

To learn the latent space of 3D caricature deformations, we adapt the *Deform-net* network architecture of DIF-Net [Deng et al. 2021], which was used to deform and edit signed distance functions. Originally, the network takes a 3D position in volume to produce deformation for the coordinate to be fed into SDF. We reformulate the network by changing the goal to learning surface deformation defined on a fixed surface. Our network consists of two modules: SIREN MLP [Sitzmann et al. 2020] for mapping a point on the template to a displacement vector and a hypernetwork that introduces variations in shapes by generating parameters for the SIREN MLP. We use the SIREN network since it converges fast and handles large shape variations effectively, as demonstrated in [Deng et al. 2021].

A possible alternative to the hypernetwork is simply concatenating the latent code to the input and intermediate features as in [Park et al. 2019; Schwarz et al. 2020]. We observed hypernetwork provides faster convergence. The conditioning via latent code concatenation took much longer time to converge. More details on our network structure and comparison between the latent code concatenation and the hypernetwork are in Sec. 4.2.

3.3 Network training

Dataset. We use 3DCaricShop dataset [Qiu et al. 2021]. The dataset contains 2K meshes sculpted by 3D artists. We obtained 1,409 registered 3D caricature meshes from the authors. The registered meshes have similar vertex connectivity as FaceWarehouse [Cao et al. 2013], except the neck area is removed, and the holes in the eyes and the mouth are closed. We used 1,268 meshes as the training set, 14 meshes as the validation set, and 127 meshes as the test set.

Along with the dataset, we obtained the mean face of FaceWarehouse that has the vertex connectivity of 3DCaricShop. We use the mean face as our template mesh. The mean face has an open mouth to avoid self-intersection around the lips. A template with an open mouth is important. If the template shape has closed mouth, since two different points on the upper lip and the lower lip have almost identical coordinates, mouth opening is impossible. With our template mesh, the mouth opens well with no problem.

Point sampling. To train the network, we need point samples on the template mesh and its corresponding displacements. A straightforward approach is to use the vertices of the template mesh. The vertices contain important samples on locations where details are required, such as eyes. On the other hand, triangles on the template mesh have different sizes, and the shapes around large triangles, e.g., the shape of the cheek, can be more accurately captured by uniform random sampling on the surface. We use a mixture of

these two approaches. Details on the point sampling algorithm are described in the supplementary material.

Our hybrid approach enables the trained network to produce more accurate shapes around large triangles compared to using vertices only. When the test set reconstruction error is calculated using uniform samples on the surface to take into account different triangle areas, the hybrid approach showed lower error of 0.0171, compared to 0.0188 of the vertices-only case.

Loss. We train the ReLU hypernetwork H and the M -dimension latent code z_k simultaneously, similarly to [Park et al. 2019]. z_k is initialized randomly using a zero-centered normal distribution of standard deviation 0.01 for each k -th training sample. A SIREN MLP S takes a template 3D position and hypernetwork-generated parameters $H(z_k)$. The training is done using the loss

$$L(p, \hat{p}, z_k) = \frac{\lambda_{mse}}{N} \sum_i \|p_i - S(\hat{p}_i, H(z_k)) - \hat{p}_i\|_2^2 + \frac{\lambda_{reg}}{M} \|z_k\|_2^2, \quad (1)$$

where p denotes the points on an example surface, \hat{p} denotes the corresponding points on the template surface, λ_{mse} controls the importance of data term, λ_{reg} controls the regularization on the latent code, N is the number of point samples. We use $\lambda_{mse} = 3.0 \times 10^3$, $\lambda_{reg} = 1.0 \times 10^6$, $N = 23132$, and $M = 128$ for training the network.

4 EXPERIMENTS

4.1 Comparison between different shape representations

Vertex position array. Fig. 5 shows the reconstruction of an unseen 3D caricature shape using two different representations for the 3D caricature: vertex position array and our proposed continuous deformation function on the surface. Using vertex position array, we find it hard to train the model to generalize well to unseen caricatures. The design of the model using the vertex position array is inspired by [Guo et al. 2019] and [Jiang et al. 2019]. Refer to the supplementary material for details on the network design and training of the model.

SDF. Fig. 6 compares visual results of training when SDF-based method is used to model 3D caricatures. We compare with DeepSDF [Park et al. 2019] and DIF-Net [Deng et al. 2021] using the authors' implementations. When we directly applied the auto-decoder framework for SDFs by converting all 3D caricatures into SDF, important details of the input were lost. Our method suffers less from the loss of details. One of the reasons is that we formulate our 3D caricature as a deformable model. Since we use a fixed template mesh that contains many crucial details, even when MLPs generate a low-frequency function, we readily obtain a reasonable 3D caricature shape due to existing details in the template.

4.2 Ablation study

In this work, we designed the network architecture based on *Deform-Net* of DIF-Net [Deng et al. 2021]. We keep the network structure of *Deform-Net* such as hypernetwork, while tweaking the training goal for the MLP to model template deformation. To inspect the

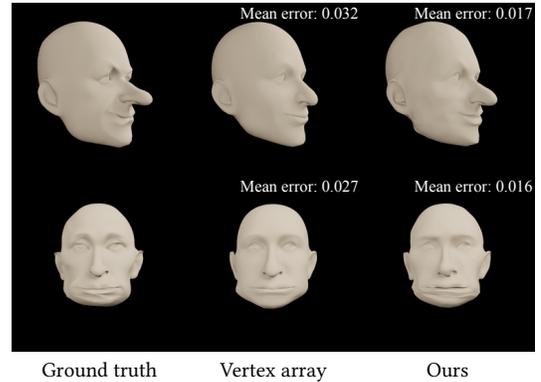


Figure 5: Comparison to using vertex position array. Using our proposed representation and network architecture, the model provides faster convergence and better generalization. The reported error is the mean l^2 distance between corresponding vertices of the ground truth and the reconstruction. On average, the error for the vertex array case was 0.031 while ours achieved 0.017 on the test set.

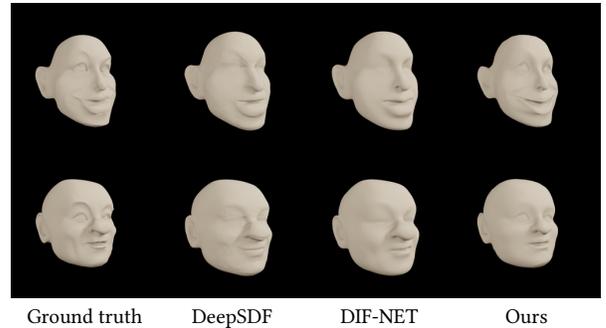


Figure 6: Visual comparison to SDF-based auto-decoder methods. We show reconstructions of 3D shapes in the training set. While DeepSDF and DIF-Net reconstruct overall shapes reasonably, the reconstructions miss important details such as the shape of eyes or the expression around the mouth.

effects of the components in our final model, we provide an ablation study of four different models.

Vertex position array MLP is the baseline model discussed in Sec. 4.1. The details for this model are in the supplementary document.

Vertex displacement array MLP has the same network architecture as *Vertex position array MLP*, but the model differs in the training goal so that the network's output is the per-vertex displacement from the template mesh.

SIREN MLP with conditioning is a tweak of our final model. Instead of using hypernetwork, similarly to [Park et al. 2019], we introduce the shape variation by concatenating the latent code to the input template position and the features of the middle layer.

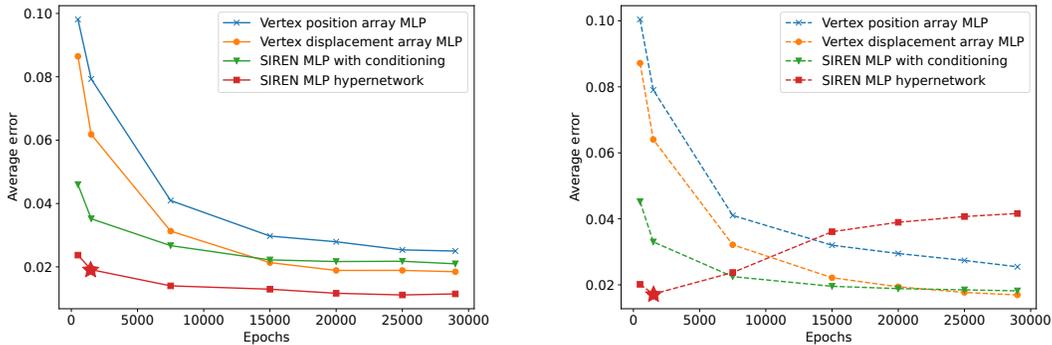


Figure 7: Ablation study on the network architecture. We inspect error on the training set (left) and test set (right) of each model at different epochs. Star (★) indicates early stopping at epoch 1500 determined by checking the validation error.

SIREN MLP hypernetwork is our proposed model in the main paper. We performed early-stopping at epoch 1500 using validation error to prevent over-fitting.

Fig. 7 provides a comparison of the average mean l^2 reconstruction error of different models at different epochs. Refer to the supplementary material for a visual comparison between the four models.

Comparing *Vertex position array MLP* and *Vertex displacement array MLP*, we observed lower training and test set error on the *Vertex displacements array MLP*. The change in the representation from absolute position to deformation plays an important role in modeling 3D caricatures accurately.

Comparing *SIREN MLP with conditioning* and *SIREN MLP hypernetwork*, we observed lower training error on the *SIREN MLP hypernetwork*. The hypernetwork provides higher network capacity. Furthermore, by comparing the error curve on the test set, we observed *SIREN MLP hypernetwork* produces much faster convergence.

In conclusion, *SIREN MLP hypernetwork* represents the complex 3D caricature data accurately by using the deformation representation. Moreover, the model converges fast by adopting the hyper-network architecture.

4.3 Continuous template deformation function

Although we trained the network on a template mesh consisting of triangles, we can generate 3D caricatures of any resolution as the MLP models a continuous function defined on the template. We visualize a continuous deformation function generated from our network by showing a 3D caricature generated from an upsampled template mesh (Fig. 8). To upsample the template mesh, we subdivide each triangle into four triangles by splitting each edge, as in Loop subdivision [Loop 1987], but without moving initial vertex positions. We apply the subdivision operation three times to upsample the template mesh of 11,551 vertices to 739,183 vertices. The result shows that the generated deformation function favors spatially smooth deformation even though the training samples have been generated from the vertices on the mesh and the samples on each triangle.

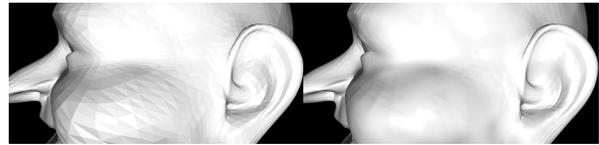


Figure 8: Visualization of the continuous deformation function learned by the network. Left: A generated 3D caricature shape sampled with the original template mesh resolution. Right: The same caricature with a denser sampling on the template mesh. The continuous function modeled by the MLP is generally smooth, so the interiors of large triangles are filled with smooth shapes.

5 APPLICATIONS

In this section, we demonstrate various applications using the learned latent space of our deep deformable 3D caricature model. In addition, we provide a supplementary video that demonstrates our applications¹. In the following, we refer to the deformable model D as a function

$$D(\hat{p}, z) = S(\hat{p}, H(z)) + \hat{p}, \quad (2)$$

where S is SIREN MLP with parameters $H(z)$, z is a latent code, and \hat{p} is a point on the template surface.

5.1 3D caricature reconstruction from 2D landmarks

In Fig. 9, we demonstrate 3D caricature reconstruction from 2D landmarks. Unlike previous work on 3D caricature reconstruction, our reconstruction is readily editable using the latent space. Alive Caricature [Wu et al. 2018] and 3DCaricShop [Qiu et al. 2021] almost perfectly fit the position constraints in the image by allowing free-form deformation of mesh vertices. Our approach reconstructs a 3D caricature by optimizing an editable latent code to reasonably

¹https://www.youtube.com/channel/UC3N03KqwKo_5gD4fISLMTg

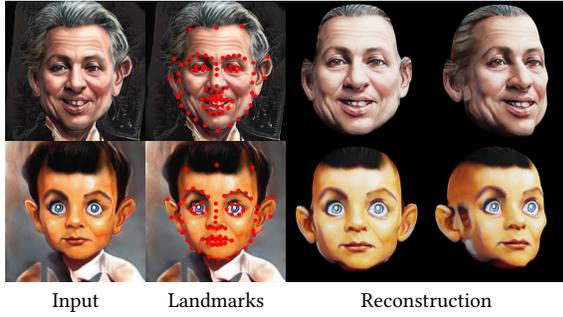


Figure 9: 3D caricature reconstruction from 2D caricature landmarks. Using manually marked sparse landmarks on the input, we optimize the latent code for our model to fit the landmark constraints. The inputs are 2D caricatures generated by StyleCariGAN [Jang et al. 2021].

fit the constraints in the input. For visual comparison with Alive Caricature and 3DCaricShop, refer to our project page².

Given 2D landmark positions b and the corresponding 3D landmark vertex indices l , the fitting is done via an alternating optimization. In a similar fashion to [Wu et al. 2018], we alternate between optimizing pose parameters Π, R, t (Π is an orthographic projection matrix, R is a rotation matrix, t is a translation vector) and the latent code for the shape z to minimize the loss:

$$L_{rec}(b, \hat{p}, z) = \frac{1}{B} \sum_i^B \|\Pi R D(\hat{p}_{l(i)}, z) + t - b_i\|_2^2 + \frac{\lambda_{reg}}{M} \|z\|_2^2, \quad (3)$$

where B is the number of 2D landmarks, and l is landmark vertex indices. We iterate the pose and the shape optimizations four times. We provide more details on the optimization algorithm in the supplementary material.

5.2 Semantic editing

We can apply semantic editing by manipulating the latent code of an input 3D face in a similar fashion to 2D caricature editing demonstrated in StyleCariGAN [Jang et al. 2021]. The WebCariA dataset [Ji et al. 2020] provides a set of semantic labels, e.g., whether a caricature is smiling or not, for each caricature in the training set. Using InterFaceGAN technique [Shen et al. 2020] to our latent space, we obtain the direction for each attribute in the WebCariA dataset. Fig. 10 demonstrates the semantic editing results on the reconstructions of the test split of the 3DCaricShop dataset.

5.3 Point-handle-based editing

Using the learned latent space, we can apply point-handle-based editing to a 3D caricature shape. Fig. 11 shows editing on the 3D caricature reconstructions from 2D caricatures. Even though the input for editing is extremely sparse, e.g., one or two points, it produces plausible deformation to complete local editing. Given a latent code z_0 for the initial 3D shape to be edited, a list of vertex indices h for the handles, and the 3D displacement vectors δ for

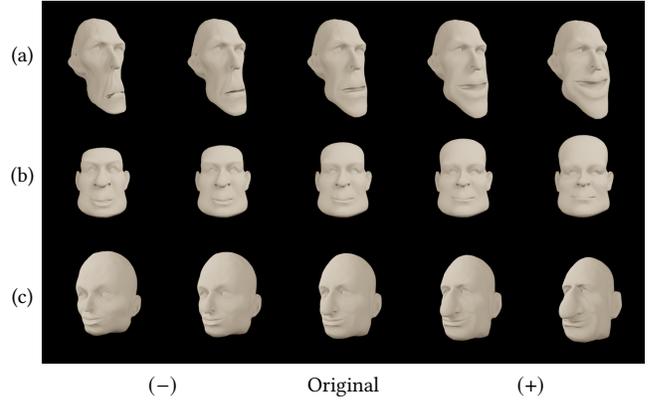


Figure 10: Semantic editing using InterFaceGAN technique. (a) Editing on the label *Smile*. (b) *Large forehead*. (c) *Big nose*. (+) denotes adding the editing vector that corresponds to the direction of the label on each row. (-) denotes subtracting the vector.

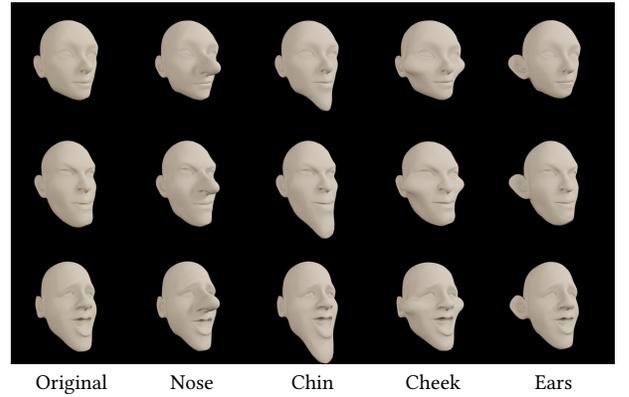


Figure 11: Point-handle-based editing. **Nose:** Pick a point on the nose tip. Move the point to the front. **Chin:** Pick a point on the bottom of the chin. Move the point to the bottom. **Cheek:** Pick a point at each side of the cheek. Stretch the two points sideways. **Ears:** Pick a point at each side of the ears. Stretch the two points.

each handle, we solve an optimization problem that minimizes

$$L_{edit}(\delta, \hat{p}, z) = \frac{\lambda_{con}}{H} \sum_i^H \|p_{h(i)} + \delta_i - D(\hat{p}_{h(i)}, z)\|_2^2 + \frac{\lambda_{pre}}{M} \|z_0 - z\|_1, \quad (4)$$

where H is the number of handles, p is the initial shape from z_0 , the weight for the constraint term $\lambda_{con} = 3.0 \times 10^3$, and the weight for term that enforces the preservation of initial shape $\lambda_{pre} = 1.0 \times 10^5$.

Our handle-based deformation can be compared to other traditional handle-based surface deformation technique such as as-rigid-as-possible (ARAP) mesh deformation [Sorkine and Alexa 2007]. Since ARAP energy is insensitive to global translation by design, the optimal deformation given a single point handle displacement is naturally a global translation without any local deformation. To

²<https://ycjungsubhuman.github.io/DeepDeformable3DCaricatures>

create local deformation using a single point in ARAP, manual selection of the region of interest (ROI) is required. Our method produces natural local deformation even with a single point constraint. Also, expansion or exaggeration, which we believe is crucial for caricatures, is not easily allowed in ARAP deformation. Our method naturally favors local expansion using the learned latent space. We provide visual comparison with ARAP in the supplementary material.

5.4 Automatic 3D caricature creation

Fig. 12 demonstrates automatic caricature creation for a regular 3D face through semantic editing. By training our model with both 3D caricatures and regular 3D faces, we obtain a model that spans both the 3D caricatures and regular faces. By calculating the editing direction from the class of regular faces to caricature faces using InterFaceGAN [Shen et al. 2020], we enable the editing operation to be applied to the latent code of a regular face.

To enable this editing, we train our model using both 3DCaricShop [Qiu et al. 2021] and FaceWarehouse [Cao et al. 2013] that provide 3D caricature examples and regular 3D face examples, respectively. The model was trained using the training set of 3DCaricShop and all neutral faces in FaceWarehouse. Automatic caricature generation is done starting from a neutral face in the FaceWarehouse dataset. We describe details on FaceWarehouse dataset processing in the supplementary material.

Another approach for automatic caricature creation is to utilize a 2D caricature generator. Given an image of a subject, we generate a 2D caricature using StyleCariGAN [Jang et al. 2021]. Using manually marked landmarks on the generated caricature, we can create an editable automatic 3D caricature (Fig. 13).

5.5 Running time

We implemented and tested our model using PyTorch [Paszke et al. 2019], and run it on a PC with an AMD Ryzen 9 3950X and an NVIDIA Titan Xp GPU. Generating a 3D shape with 11,551 vertices takes 17 ms, and a 3D shape with 739,183 vertices takes 987 ms (Sec. 4.3). The alternating optimization for our 3D caricature reconstruction from 2D landmarks takes approximately 500 ms (Sec. 5.1). The optimization for our point-handle-based editing took approximately 20 ms (Sec. 5.3).

6 CONCLUSION

To build a useful editing space for 3D caricatures, we proposed a deformable 3D caricature framework. By adopting MLPs to model a continuous function of template surface deformation, we model the complex variations of 3D caricatures effectively. Once the model has been learned, it can be readily used for other tasks by editing latent codes or optimizing the latent codes to meet various constraints.

An interesting implication of our method is that each caricature is modeled as a continuous and differentiable function. This setting provides easy access to differential properties of deformation without using finite element methods. For example, by inspecting MLP parameters, we may easily inspect the smoothness of deformation function as demonstrated for SDFs in [Liu et al. 2022].

Our network may not completely restore high-frequency details originally present in ground-truth 3D caricatures (See Figs. 5 and 6).

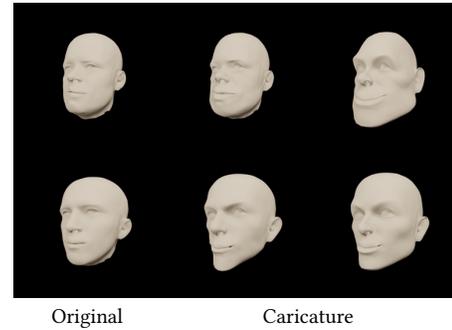


Figure 12: Automatically created caricatures of varying degrees. Given a regular face, we apply semantic editing towards caricatures. The third column is edited three times further towards caricatures compared to the second column.

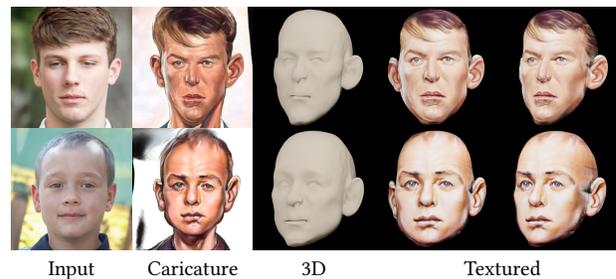


Figure 13: Automatic 3D caricature creation via 3D caricature reconstruction on the results of StyleCariGAN. The inputs are photos generated by StyleGAN [Karras et al. 2019]. This application also works for real-world photo inputs.

We also observed self-intersection around the chin from some reconstruction results, where training data often contain artifacts. Our point-handle-based editing may result in self-intersection given some input since we do not put any explicit protection from self-intersection. Improving the details and preventing self-intersections are left as future work. An explicit guarantee for identity preservation in editing would also be an interesting future work.

We believe our method can be applied to other domains where the data to be modeled are surfaces in dense correspondence. Our method assumes the template mesh does not have self-intersections. When the assumption does not hold, a different position encoding scheme would be required to differentiate two different points on the template surface sharing the same coordinates.

ACKNOWLEDGMENTS

We would like to thank the anonymous reviewers for their constructive comments. This work was supported by IITP grants (SW Star Lab, 2015-0-00174; Collaborative Research Project with MSRA, IITP-2020-0-01649; AI Innovation Hub, 2021-0-02068; AI Graduate School Program (POSTECH), 2019-0-01906) and KOCCA grant (R2021040136) from Korea government (MSIT and MCST).

REFERENCES

- Volker Blanz, Thomas Vetter, et al. 1999. A morphable model for the synthesis of 3D faces. In *Proc. SIGGRAPH*.
- Hongrui Cai, Yudong Guo, Zhuang Peng, and Juyong Zhang. 2021. Landmark detection and 3D face reconstruction for caricature using a nonlinear parametric model. *Graphical Models* 115 (2021), 101103.
- Chen Cao, Yanlin Weng, Shun Zhou, Yiyong Tong, and Kun Zhou. 2013. Facewarehouse: A 3d facial expression database for visual computing. *IEEE Trans. Vis. Comput. Graph.* 20, 3 (2013), 413–425.
- Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. 2021. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proc. CVPR*.
- Shiyang Cheng, Michael Bronstein, Yuxiang Zhou, Irene Kotsia, Maja Pantic, and Stefanos Zafeiriou. 2019. Meshgan: Non-linear 3d morphable models of faces. *arXiv preprint arXiv:1903.10384* (2019).
- Yu Deng, Jiaolong Yang, and Xin Tong. 2021. Deformed implicit field: Modeling 3d shapes with learned dense correspondence. In *Proc. CVPR*.
- Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. 2019. Accurate 3D Face Reconstruction with Weakly-Supervised Learning: From Single Image to Image Set. In *Proc. CVPRW*.
- Yao Feng, Fan Wu, Xiaohu Shao, Yanfeng Wang, and Xi Zhou. 2018. Joint 3d face reconstruction and dense alignment with position map regression network. In *Proc. ECCV*.
- Pablo Garrido, Michael Zollhöfer, Dan Casas, Levi Valgaerts, Kiran Varanasi, Patrick Pérez, and Christian Theobalt. 2016. Reconstruction of personalized 3D face rigs from monocular video. *ACM Transactions on Graphics* 35, 3 (2016), 28.
- Donya Ghafourzadeh, Cyrus Rahgoshay, Sahel Fallahdoust, Andre Beauchamp, Adeline Aubame, Tiberiu Popa, and Eric Paquette. 2020. Part-Based 3D Face Morphable Model with Anthropometric Local Control. In *Proc. GI*.
- Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. 2020. Implicit Geometric Regularization for Learning Shapes. In *Proc. MLR*.
- Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. 2018. A papier-mâché approach to learning 3d surface generation. In *Proc. CVPR*.
- Yudong Guo, Luo Jiang, Lin Cai, and Juyong Zhang. 2019. 3D Magic Mirror: Automatic Video to 3D Caricature Translation. *arXiv preprint arXiv:1906.00544* (2019).
- David Ha, Andrew Dai, and Quoc V. Le. 2017. HyperNetworks. In *Proc. ICLR*.
- Xiaoguang Han, Chang Gao, and Yizhou Yu. 2017. DeepSketch2Face: a deep learning based sketching system for 3D face and caricature modeling. *ACM Trans. Graph.* 36, 4 (2017), 126.
- Xiaoguang Han, Kangcheng Hou, Dong Du, Yuda Qiu, Shuguang Cui, Kun Zhou, and Yizhou Yu. 2018. Caricatureshop: Personalized and photorealistic caricature sketching. *IEEE Trans. Vis. Comput. Graph.* 26, 7 (2018), 2349–2361.
- Wonjong Jang, Gwangjin Ju, Yuchoel Jung, Jiaolong Yang, Xin Tong, and Seungyong Lee. 2021. StyleCariGAN: caricature generation via StyleGAN feature map modulation. *ACM Trans. Graph.* 40, 4 (2021), 1–16.
- Wen Ji, Kelei He, Jing Huo, Zheng Gu, and Yang Gao. 2020. Unsupervised domain attention adaptation network for caricature attribute recognition. In *Proc. ECCV*.
- Zi-Hang Jiang, Qianyi Wu, Keyu Chen, and Juyong Zhang. 2019. Disentangled representation learning for 3D face shape. In *Proc. CVPR*.
- Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *Proc. CVPR*.
- Hsueh-Ti Derek Liu, Francis Williams, Alec Jacobson, Sanja Fidler, and Or Litany. 2022. Learning Smooth Neural Functions via Lipschitz Regularization. *arXiv preprint arXiv:2202.08345* (2022).
- Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. 2020. Neural sparse voxel fields. *Advances in Neural Information Processing Systems* 33 (2020), 15651–15663.
- Stephen Lombardi, Jason Saragih, Tomas Simon, and Yaser Sheikh. 2018. Deep appearance models for face rendering. *ACM Transactions on Graphics (ToG)* 37, 4 (2018), 1–13.
- Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. 2019. Neural volumes: Learning dynamic renderable volumes from images. *arXiv preprint arXiv:1906.07751* (2019).
- Charles Loop. 1987. *Smooth Subdivision Surfaces Based on Triangles*. Ph. D. Dissertation.
- Marcel Lüthi, Thomas Gerig, Christoph Jud, and Thomas Vetter. 2017. Gaussian process morphable models. *IEEE Trans. PAMI* 40, 8 (2017), 1860–1873.
- Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2020. Nerf: Representing scenes as neural radiance fields for view synthesis. In *Proc. ECCV*.
- Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. 2019. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proc. CVPR*.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan
- Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems* 32, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.). Curran Associates, Inc., 8024–8035. <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
- Stylios Ploumpis, Evangelos Ververas, Eimear O'Sullivan, Stylios Moschoglou, Haoyang Wang, Nick Pears, William AP Smith, Baris Gecer, and Stefanos Zafeiriou. 2020. Towards a complete 3D morphable model of the human head. *IEEE Trans. PAMI* 43, 11 (2020), 4142–4160.
- Yuda Qiu, Xiaojie Xu, Lingteng Qiu, Yan Pan, Yushuang Wu, Weikai Chen, and Xiaoguang Han. 2021. 3DCaricShop: A Dataset and A Baseline Method for Single-view 3D Caricature Face Reconstruction. In *Proc. CVPR*.
- Anurag Ranjan, Timo Bolkart, Soubhik Sanyal, and Michael J Black. 2018. Generating 3D faces using convolutional mesh autoencoders. In *Proc. ECCV*.
- Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. 2019. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proc. ICCV*.
- Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. 2020. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems* 33 (2020), 20154–20166.
- Matan Sela, Yonathan Aflalo, and Ron Kimmel. 2015. Computational caricaturization of surfaces. *Computer Vision and Image Understanding* 141 (2015), 1–17.
- Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. 2020. Interfacegan: Interpreting the disentangled face representation learned by gans. *IEEE Trans. Pattern Anal. Mach. Intell.* (2020).
- Vincent Sitzmann, Julien NP Martel, Alexander W Bergman, David B Lindell, and Gordon Wetzstein. 2020. Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems* 33 (2020), 7462–7473.
- Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. 2019. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *Advances in Neural Information Processing Systems* 32 (2019).
- Olga Sorkine and Marc Alexa. 2007. As-rigid-as-possible surface modeling. In *Proc. SGP*.
- Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. 2016. Face2face: Real-time face capture and reenactment of rgb videos. In *Proc. CVPR*.
- Weiyue Wang, Duygu Ceylan, Radomir Mech, and Ulrich Neumann. 2019. 3dn: 3d deformation network. In *Proc. CVPR*.
- Qianyi Wu, Juyong Zhang, Yu-Kun Lai, Jianmin Zheng, and Jianfei Cai. 2018. Alive Caricature from 2D to 3D. In *Proc. CVPR*.

Supplementary material: Deep Deformable 3D Caricatures with Learned Shape Control

Yucheol Jung
POSTECH
Pohang, Republic of Korea
ycjung@postech.ac.kr

Wonjong Jang
POSTECH
Pohang, Republic of Korea
wonjong@postech.ac.kr

Soongjin Kim
POSTECH
Pohang, Republic of Korea
kimsj0302@postech.ac.kr

Jiaolong Yang
Microsoft Research Asia
Beijing, China
jiaoyan@microsoft.com

Xin Tong
Microsoft Research Asia
Beijing, China
xtong@microsoft.com

Seungyong Lee
POSTECH
Pohang, Republic of Korea
leesy@postech.ac.kr

We provide details of the proposed method and more experimental results in this supplementary document. Sec. 1 provides quantitative analysis on the complexity of the 3DCaricShop dataset. Sec. 2 provides details on our training and the network structure. Sec. 3 provides more details on Sec 4.1. in the main paper. Sec. 4 provides visual results for the ablation study in the main paper. Sec. 5 provides details on the 3D caricature reconstruction from 2D landmarks. Sec. 6 provides details on our semantic face editing. Sec. 7 visually compares our point-handle-based editing to as-rigid-as-possible (ARAP) deformation [Sorkine and Alexa 2007]. Lastly, Sec. 8 explains data processing details in our automatic 3D caricature creation. For a large gallery of results using real-world caricatures and visual comparison with Alive Caricature [Wu et al. 2018] and 3DCaricShop [Qiu et al. 2021], refer to our project page¹.

1 DATASET ANALYSIS

Our key challenge in learning an editing space is the high variation and complexity of 3D caricatures. We observed the 3DCaricShop [Qiu et al. 2021] dataset contains much more complex information compared to other common 3D mesh datasets. We compare 3DCaricShop with a regular face dataset (FaceWarehouse [Cao et al. 2013]), a dynamic human body dataset (DFAUST [Bogo et al. 2017]), and a facial animation dataset (COMA [Ranjan et al. 2018]).

We evaluate the complexity of each dataset by observing the effectiveness of dimension reduction on the dataset. We apply principal component analysis (PCA) [Pearson 1901] to each dataset and calculate the number of principal components required to explain 99% of variations in the dataset. 3DCaricShop requires 178 principal components, which is a huge number compared to FaceWarehouse (52), DFAUST (22), and COMA (41). 3DCaricShop provides 2K examples, which are sparse samples compared to other datasets: FaceWarehouse (7K), DFAUST (41K), and COMA (20K). In addition, 3DCaricShop contains diverse high-frequency variations, while other datasets contain mostly low-frequency variations from shape, pose, or expression.

¹<https://ycjungsubhuman.github.io/DeepDeformable3DCaricatures>

Authors' addresses: Yucheol Jung, POSTECH, Pohang, Republic of Korea, ycjung@postech.ac.kr; Wonjong Jang, POSTECH, Pohang, Republic of Korea, wonjong@postech.ac.kr; Soongjin Kim, POSTECH, Pohang, Republic of Korea, kimsj0302@postech.ac.kr; Jiaolong Yang, Microsoft Research Asia, Beijing, China, jiaoyan@microsoft.com; Xin Tong, Microsoft Research Asia, Beijing, China, xtong@microsoft.com; Seungyong Lee, POSTECH, Pohang, Republic of Korea, leesy@postech.ac.kr.

2 TRAINING DETAILS

Our model is optimized with the Adam optimizer [Kingma and Ba 2015] using batch size 128, with learning rate 1.0×10^{-4} , $\beta_1 = 0.9$, and $\beta_2 = 0.999$.

SIREN MLP. We construct the SIREN MLP with five fully connected layers: the first SIREN layer mapping a 3D position to a 128-dimension feature, three hidden SIREN layers with hidden feature dimension 128, and the last layer with feature dimension 3.

hypernetwork. We have two ReLU MLPs for each fully-connected layer in the SIREN MLP: One for weights and the other for biases. Each ReLU MLP produces network parameters given a 128-dimension latent code. The ReLU MLPs have three fully connected layers each: the first two ReLU layers with output feature dimension 256, and the last fully connected layer without activation. For more details on hypernetworks, refer to [Sitzmann et al. 2020].

Point sampling. We sample two groups of points from a mesh; vertex position in the mesh and uniformly sampled points on the surface. Collating the two groups, we obtain 23,132 points for each example. The sampling is done statically before training. Each point p_i on an example is mapped to the corresponding point \hat{p}_i on the fixed template mesh by converting a vertex position on the example mesh into barycentric coordinates, then converting the barycentric coordinates to a template mesh position.

3 DETAILS ON COMPARISON BETWEEN DIFFERENT SHAPE REPRESENTATION (SEC. 4.1)

The design of the compared model using the vertex position array is inspired by [Guo et al. 2019] and [Jiang et al. 2019]. The model is a single MLP, and the training is done using the same auto-decoder architecture as ours. The MLP takes a 128-dimension latent code to produce an array of 3D vertex positions. The MLP is constructed using two fully connected layers with the following channel dimensions: $(128) \rightarrow (400) \rightarrow (3 \times 11551)$. The activation for the first layer is leaky ReLU with a negative slope of 0.1. The number of parameters is 13.9M, which is similar to that of ours with 14M parameters. The training is done using the same loss, optimizer, and batch size as our network. We show the results of the compared model after approximately 17,500 epochs of training (8 hours). The

result for our proposed network is from 1,500 epochs of training (2 hours). Both trainings were done on an Intel(R) Xeon(R) Silver 4114 server using two NVIDIA Titan Xp GPUs. The reconstruction for the comparison using the test set is done by optimizing latent codes using the same loss used in training. In this optimization, we only use vertex positions on the mesh without random sampling on the surface.

4 VISUAL RESULTS OF ABLATION STUDY

In addition to error curves in the main paper, we provide visual results of each model at different epochs. Fig. 1 visualizes the accuracy of the reconstruction on the training set. The fitting on the training set shows high network capacity of our final model. Fig. 2 visualizes the accuracy of the test set reconstruction. We early stop at epoch 1500. The result of our final model shows a good visual reconstruction compared to other models.

5 DETAILS ON 3D CARICATURE RECONSTRUCTION FROM 2D LANDMARKS

The 2D caricature landmarks are annotated on the silhouette. If the face pose contains large rotation, 3D landmark vertices on the 3D mesh may become hidden and not correspond to 2D landmarks. In our setting, when the input 2D caricature is not frontal, 3D landmarks on the chin may not correspond to 2D landmarks. This problem has been pointed out in [Jiang et al. 2018; Wu et al. 2018]. Their solution is to update 3D landmark vertex indices on the chin after every pose update. Similarly, we apply the silhouette landmark update algorithm after each pose estimation.

For pose estimation, we use the 3D-to-2D point alignment algorithm presented in [Kemelmacher-Shlizerman and Seitz 2011]. For shape optimization, we update latent code using the same Adam optimizer with the hyperparameters we used for the training. The shape optimization is done until $(L_t - L_{t-1})/L_t < 1.0 \times 10^{-10}$, where L is the loss value at t -th update.

6 DETAILS ON INTERFACEGAN IN OUR SEMANTIC EDITING

Given semantic labels on the training set, InterFaceGAN [Shen et al. 2020] enables semantic editing on the latent space given a GAN model. Although we do not use GAN architecture, we find their semantic editing is applicable to our network since our model also maps a latent vector to a face.

We use their *single attribute manipulation* technique to perform our semantic manipulation. We have semantic WebCariA [Ji et al. 2020] labels for each latent vector corresponding to the 3D caricature in the training set. We find the latent direction vector for editing each semantic in the WebCariA by training a linear support vector machine (SVM) [Cortes and Vapnik 1995] for classifying each label given a latent vector. The training of SVMs is done using the WebCariA authors' published code. The normal of the hyper-plane for each SVM is the direction for editing each semantic label. We edit each face by scaling the editing vector and adding it to the source shape. The scaling factors are decided empirically.

7 VISUAL COMPARISON TO ARAP DEFORMATION

Handle-based deformation algorithms based on elastic energy on the surface have been extensively studied. One of the most popular of such algorithms is based on *as-rigid-as-possible* (ARAP) deformation [Sorkine and Alexa 2007]. We compare our data-driven point-handle-based deformation with ARAP deformation using the same set of constraints on the shape (Fig. 3). We used the ARAP implementation of libigl [Jacobson et al. 2018] in the comparison.

ARAP deformation cannot support local deformations with high stretch and area expansion that are required for face exaggeration. Without other constraints, the face generated by ARAP deformation is not guaranteed to be a valid caricature face. In contrast, our handle-based editing is constrained by the learned MLP decoder and guarantees the result to be within the latent space of caricature faces.

8 FACEWAREHOUSE DATASET PROCESSING IN AUTOMATIC 3D CARICATURE CREATION

To train the model using both 3D caricatures and regular 3D faces, we use 3DCaricShop [Qiu et al. 2021] and FaceWarehouse [Cao et al. 2013] datasets. As the two sets of meshes differ in mesh connectivity, we use R3DS Wrap [Russian3DScanner 2022] to register our template mesh to the mean FaceWarehouse face. Using the correspondence, we obtain the training set for regular faces in FaceWarehouse using the following procedure: 1) Randomly sample 40K points on each regular face. 2) Find the corresponding point on the FaceWarehouse mean face. 3) Find the closest point on our template mesh. Since some points on the FaceWarehouse mesh do not have a valid corresponding point to our template mesh, e.g., points on the neck, we filter out points with errors higher than 0.3. Then, we select 23,132 random points to obtain point samples for training a regular 3D face.

REFERENCES

- Federica Bogo, Javier Romero, Gerard Pons-Moll, and Michael J. Black. 2017. Dynamic FAUST: Registering Human Bodies in Motion. In *Proc. CVPR*.
- Chen Cao, Yanlin Weng, Shun Zhou, Yiyong Tong, and Kun Zhou. 2013. Facewarehouse: A 3d facial expression database for visual computing. *IEEE Trans. Vis. Comput. Graph.* 20, 3 (2013), 413–425.
- Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning* 20, 3 (1995), 273–297.
- Yudong Guo, Luo Jiang, Lin Cai, and Juyong Zhang. 2019. 3D Magic Mirror: Automatic Video to 3D Caricature Translation. *arXiv preprint arXiv:1906.00544* (2019).
- Alec Jacobson, Daniele Panozzo, et al. 2018. libigl: A simple C++ geometry processing library. <http://libigl.github.io/libigl/>.
- Wen Ji, Kelei He, Jing Huo, Zheng Gu, and Yang Gao. 2020. Unsupervised domain attention adaptation network for caricature attribute recognition. In *Proc. ECCV*.
- Luo Jiang, Juyong Zhang, Bailin Deng, Hao Li, and Ligang Liu. 2018. 3D face reconstruction with geometry details from a single image. *IEEE Trans. Image Process.* 27, 10 (2018), 4756–4770.
- Zi-Hang Jiang, Qianyi Wu, Keyu Chen, and Juyong Zhang. 2019. Disentangled representation learning for 3D face shape. In *Proc. CVPR*.
- Ira Kemelmacher-Shlizerman and Steven M Seitz. 2011. Face reconstruction in the wild. In *2011 international conference on computer vision*. IEEE, 1746–1753.
- Diederik P Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *Proc. ICLR*.
- Karl Pearson. 1901. LIII. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin philosophical magazine and journal of science* 2, 11 (1901), 559–572.
- Yuda Qiu, Xiaojie Xu, Lingteng Qiu, Yan Pan, Yushuang Wu, Weikai Chen, and Xiaoguang Han. 2021. 3DCaricShop: A Dataset and A Baseline Method for Single-view

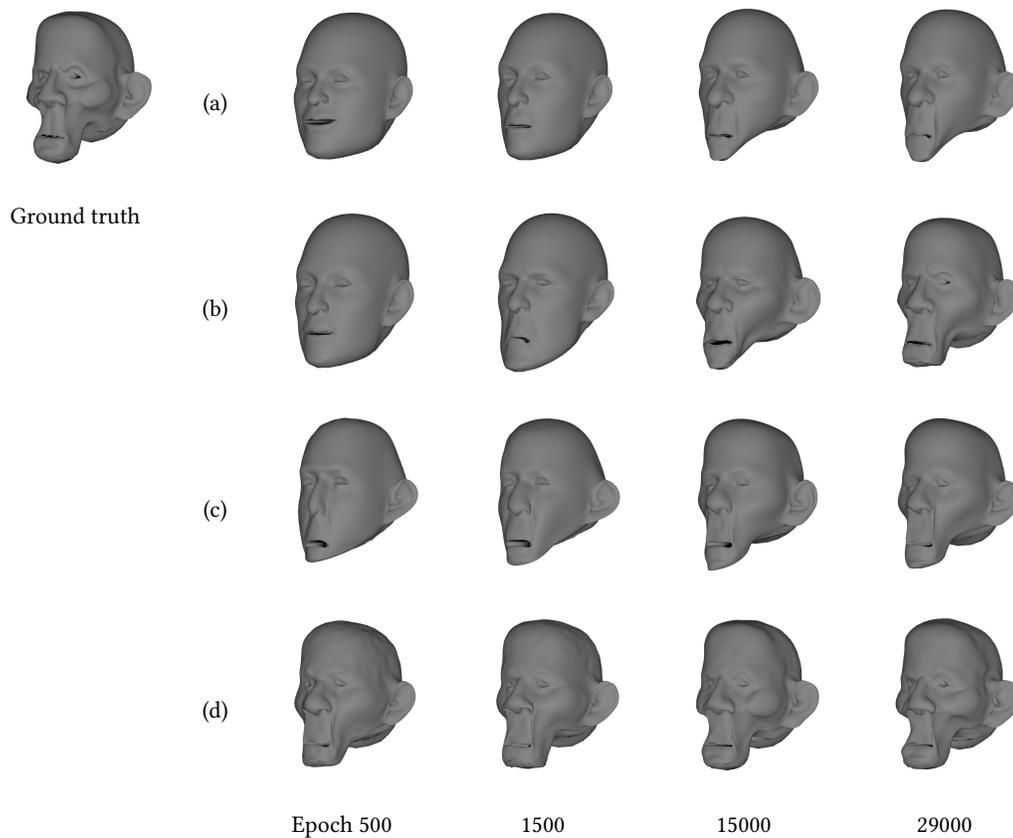


Figure 1: Visual result of training set reconstruction for different model configurations. (a) Vertex position array MLP. (b) Vertex displacement array MLP. (c) SIREN MLP with conditioning. (d) SIREN MLP hypernetwork.

3D Caricature Face Reconstruction. In *Proc. CVPR*.
 Anurag Ranjan, Timo Bolkart, Soubhik Sanyal, and Michael J Black. 2018. Generating 3D faces using convolutional mesh autoencoders. In *Proc. ECCV*.
 Russian3DScanner. 2022. *R3DS Wrap*. <https://www.russian3dscanner.com/>
 Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. 2020. Interfacegan: Interpreting the disentangled face representation learned by gans. *IEEE Trans. Pattern Anal. Mach. Intell.* (2020).

Vincent Sitzmann, Julien NP Martel, Alexander W Bergman, David B Lindell, and Gordon Wetzstein. 2020. Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems* 33 (2020), 7462–7473.
 Olga Sorkine and Marc Alexa. 2007. As-rigid-as-possible surface modeling. In *Proc. SGP*.
 Qianyi Wu, Juyong Zhang, Yu-Kun Lai, Jianmin Zheng, and Jianfei Cai. 2018. Alive Caricature from 2D to 3D. In *Proc. CVPR*.

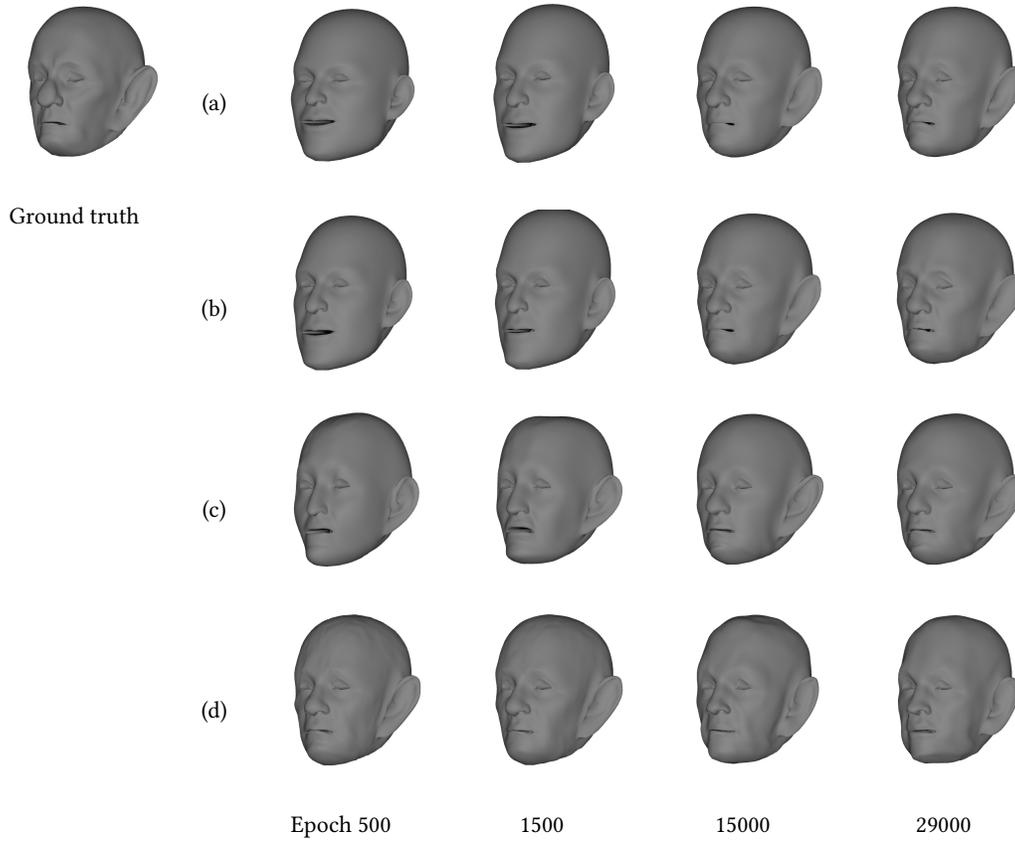


Figure 2: Visual result of test set reconstruction for different model configurations. (a) Vertex position array MLP. (b) Vertex displacement array MLP. (c) SIREN MLP with conditioning. (d) SIREN MLP hypernetwork.

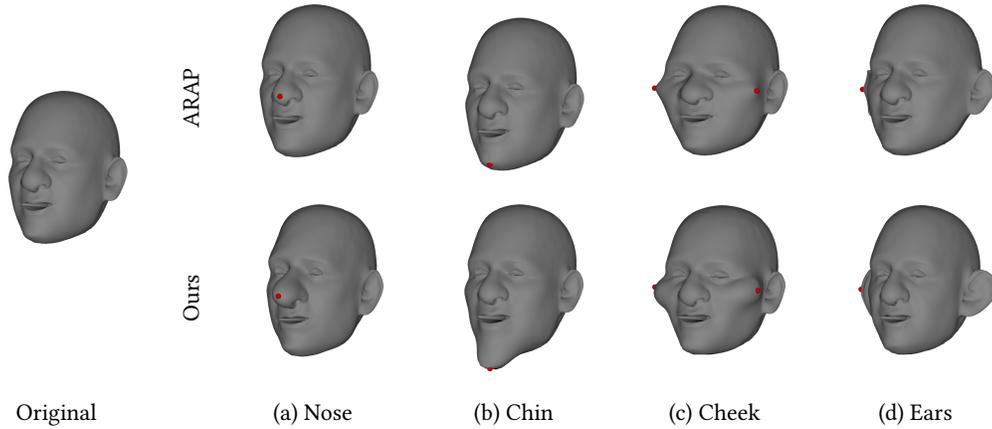


Figure 3: Visual comparison to ARAP deformation. Red points indicate the control points. (a,b) Since ARAP energy is insensitive to global translation by design, the optimal deformation given a single point handle displacement is naturally a global translation without any local deformation. To create local deformation using the single point in ARAP, manual selection of the region of interest (ROI) is required. Our method produces natural local deformation even with a single point constraint. (c,d) The ARAP energy increases as the area of local cells increase, so expansion or exaggeration, which we believe is crucial for caricatures, is not easily allowed in ARAP deformation. Our method naturally favors local expansion using the learned latent space.