Texture Generation Using A Graph Generative Adversarial Network And Differentiable Rendering

Dharma KC^{1[0000-0003-0676-2391]}, Clayton T. Morrison^{1[0000-0002-3606-0078]} and Bradley Walls^{2[0000-0002-5484-7587]}

 ¹ The University of Arizona, Tucson AZ 85721, USA {kcdharma,claytonm}@arizona.edu
² Areté Associates, Tucson AZ 85712, USA bwalls@arete.com

Abstract. Novel photo-realistic texture synthesis is an important task for generating novel scenes, including asset generation for 3D simulations. However, to date, these methods predominantly generate textured objects in 2D space. If we rely on 2D object generation, then we need to make a computationally expensive forward pass each time we change the camera viewpoint or lighting. Recent work that can generate textures in 3D requires 3D component segmentation that is expensive to acquire. In this work, we present a novel conditional generative architecture that we call a graph generative adversarial network (GGAN) that can generate textures in 3D by learning object component information in an unsupervised way. In this framework, we do not need an expensive forward pass whenever the camera viewpoint or lighting changes, and we do not need expensive 3D part information for training, yet the model can generalize to unseen 3D meshes and generate appropriate novel 3D textures. We compare this approach against state-of-the-art texture generation methods and demonstrate that the GGAN obtains significantly better texture generation quality (according to Fréchet inception distance). We release our model source code as open source.³

Keywords: 3D texture synthesis \cdot Graph neural networks \cdot Differentiable rendering.

1 Introduction

Synthesizing novel photorealistic textures for 3D mesh models is an important task for the generation of novel scenes in static images or realistic 3D simulations. Such generated textures can be applied to 3D mesh models and rendered with different lighting conditions and camera angles quite easily. The generative adversarial network (GAN) framework [11] is a promising approach to training models capable of novel image generation. However, extending the GAN framework to support texture generation in 3D that can generalize to novel, previously

³ https://github.com/ml4ai/ggan

unseen 3D models poses interesting challenges. Generating textures in 2D space (u, v coordinate system) and then wrapping to 3D mesh models won't generalize to unseen meshes because the UV mapping function is different for different 3D meshes. However, humans are able to identify the components of a 3D object and could texture them consistently. This raises an interesting research question: can we design an algorithm that can generate realistic textures for unseen 3D mesh models by identifying object components as humans do? We present here a system that addresses this challenge. Our model can learn to distinguish 3D part information shared across instances of an object class (e.g., wheels, doors,hood, windows, tail, headlights, etc. of a car) in an unsupervised way that supports generating specific texture features for these components. Recent work presented a new model called TM-NET [9] that can generate textures in 3D but requires prior supervised segmentation of 3D parts, while our model can identify 3D part information in an unsupervised way. Another closely related work to ours is [34], but they work on 2.5D space rather than on the original 3D space, forcing us to make an expensive forward pass each time the viewpoint changes. We use PyTorch [30] and PyTorch3D [35] for the implementation of our system and use the ShapeNet dataset [3] to train and evaluate our framework. We adopt the commonly used Fréchet Inception Distance (FID) [15] to assess the quality of textures generated by model. FID is typically applied to 2D images. To extend the FID measure to texture map generation for 3D models, we apply the generated texture map that is to be evaluated to the given 3D model and render it from multiple views, producing multiple 2D images, and the FID scores across these images are aggregated to produce a summary FID score. The major contributions of our paper are as follows:

- We present a simple solution to the challenging problem of 3D texture generation for 3D mesh models rather than 2D or 2.5D images.
- Our framework is capable of unsupervised learning of part information shared across a class of objects, therefore avoiding a costly, separate supervised learning task in order to learn textures appropriate to object parts.
- We present a thorough review, analysis, and evaluation of various techniques that can be used for texture generation for 3D meshes and demonstrate the advantages and disadvantages of each.

2 Related work

In this section, we describe five threads of work related to our problem and proposed framework.

2.1 Generative adversarial networks

Generative adversarial networks (GANs) [11] are known for their ability to generate photorealistic images with very high resolution [20]. The GAN framework consists of a generator \mathcal{G} and discriminator \mathcal{D} . The generator \mathcal{G} attempts to generate realistic images that can fool the discriminator. At the same time, the discriminator \mathcal{D} tries to predict whether the image is "real" (comes from the real data distribution) or "fake" (generated by the generator). This constitutes a two-player minimax game with the following value function:

$$\mathcal{V}(\mathcal{G}, \mathcal{D}) = \mathbb{E}_{\boldsymbol{x} \sim p_{data}(\boldsymbol{x})} [log D(\boldsymbol{x})] \\ + \mathbb{E}_{z \sim p_{z}(\boldsymbol{z})} [log (1 - \mathcal{D}(\mathcal{G}(\boldsymbol{z})))]$$
(1)

Here, \boldsymbol{x} denotes a sample from a real distribution, p_{data} , and \boldsymbol{z} denotes a "noise" vector from distribution p_z . D(x) denotes the probability that x comes from the real distribution. Multiple applications have adapted GANs for the task of generating realistic images. Specifically, Radford et al. [33] propose deep convolutional GAN (DCGAN), which uses convolutional neural networks (CNNs) to generate low-resolution photorealistic images. Recently, Karras et al. [18, 19, 20] developed a novel architecture for the generator. Arjovsky et al. [1] propose the Wasserstein GAN (WGAN) to improve the stability of training. Xian et al. [39] propose image synthesis with texture, but this only works on 2D images. Mirza et al. [29] propose a conditional GAN framework that can generate samples from a specified class. In this framework, the noise vector is combined with a class label to generate a sample image from that class. This work is closely related to our work as we seek to generate a texture conditioned on an input mesh. This work is a simplified version of our problem as they work on 2D images and the combination of the noise vector with the class label (a one-hot vector) can be easily achieved with a simple concatenation while a simple concatenation of the noise vector with a 3D mesh model is not possible. This makes our problem challenging and requires a new architecture.

2.2 Differentiable rendering



Fig. 1: Differentiable rendering

The second line of work related to ours is *differentiable rendering*. Rendering in computer graphics is a process of generating a 2D image from a 3D mesh,

light source, camera properties, texture properties, and other scene properties. Classical rendering using rasterization, or ray tracing, is not differentiable. This means we cannot propagate the gradients of the loss in image space (2D space) with respect to mesh properties such as vertices and textures. Given that we want to generate a realistic texture for a given 3D mesh model with supervision from 2D images, we need a way to propagate the gradients of the loss from these projected (rendered) 2D images back to the 3D scene properties. Differentiable rendering is a process that enables backpropagating these gradients from the 2D image loss back into the 3D scene properties. Figure 1 illustrates the differentiable rendering part of our architecture. Recent methods propose approximate solutions for making the rendering process differentiable [21, 25, 26, 27, 35]. We use PyTorch3D [35] for differentiable rendering.

2.3 Texturing

Texturing is the process of applying a texture to a given 3D mesh model. There are multiple ways to apply a texture to a 3D mesh. Given that we want to generate textures for 3D polygonal meshes that can be applied directly to 3D shapes, we have the following options for texturing [35]:

UV textures: A *UV texture* is a 2D image that can be mapped to 3D mesh model. For UV textures to work, we need a 2D UV-coordinate image and a mapping function that maps every vertex in the 3D object space to a (u, v) coordinate in a 2D UV image. The advantage of this method is that it can represent high-resolution textures, but the mapping function is different for different meshes, making the texturing process hard to generalize across varieties of 3D models. We refer to this method as TEXTUREUV in the following sections.

Vertex textures: A vertex texture defines a texture per vertex (e.g. r, g, b color). If the mesh has V vertices, and the dimension of texture per vertex is D, the texture can be represented by a tensor of shape (V, D) for a given 3D mesh. In this approach, the texture within faces between vertices has to be interpolated from vertex textures. This makes the vertex texture suitable only for low-resolution textures. We refer to this method as TEXTUREVERTEX in the following sections.

Face textures: The *face texture* method defines a separate texture per face. The texture per face can be an RxR dimensional texture, where R is the resolution of a texture image of a single face. This can be modeled using a (F, R, R, D) tensor where F is the number of faces, R is the resolution of a texture and D is the dimension of a texture (*e.g.*, D = 3: r, g, b colors). This allows us to learn very high-resolution textures. We refer to it as TEXTUREFACE in the following sections.

4

2.4 Deformable models

This family of work learns to generate the 3D mesh along with textures from 2D images. The method generally starts with a fixed geometry (e.g. sphere) and a fixed UV mapping. Given input images, the model extracts information from these images, represented as a latent vector. This vector is then used to predict the deformation of the vertices of the sphere template to approximate the 3D mesh and the texture image. Then the estimated 3D shape, estimated texture, and fixed UV map go through a differentiable rendering step to generate a 2D image. The main idea then is to make these generated images similar to the original images, which can be achieved using reconstruction loss and adversarial loss. Figure 2 shows the architecture.



Fig. 2: General architecture of deformable models

Recent work has explored variations of this idea and has achieved good results [5, 10, 14, 17, 31, 32, 32, 41]. This is a really good approach when we don't have ground truth 3D meshes. But, when the 3D meshes are available, as in our case, the major disadvantage of this method is that the predicted mesh tends to be relatively poor quality compared to the ground-truth mesh.

2.5 Graph neural networks

Graph neural networks (GNNs) are powerful models for learning from graphstructured data. They work on the theory of message passing, where a node gets some information from its neighbors and updates its state. Consider a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} is the set of nodes and \mathcal{E} is the set of edges. Let, $\mathcal{X} \in \mathbb{R}^{|v|*d}$ be the set of node features where each node $v \in \mathcal{V}$ has a *d* dimensional feature. The k^{th} message passing iteration of a GNN can be modeled as a variation of the following equation [12]:

$$\begin{aligned} h_v^{(k+1)} &= \texttt{update}^{(k)}(h_v^{(k)}, \texttt{aggregate}^{(k)}(h_u^{(k)}), \\ &= \texttt{update}^{(k)}(h_v^{(k)}, \boldsymbol{m}_{\mathcal{N}(v)}^{(k)}) \quad \forall u \in \mathcal{N}(v)) \end{aligned}$$
(2)

Here, $\mathcal{N}(v)$ denotes the neighbors of node v. At any iteration of the GNN, the **aggregate** function takes the embedding of the neighbors of node v and combines them into one embedding vector. The **update** function takes the embedding of the node v at the previous time step and the output embedding vector of the **aggregate** function to give us the new embedding for the node v. Here, **update** and **aggregate** can be any differentiable functions. In our work here, we convert an input 3D mesh model to a graph and use the power of GNNs to learn latent part information of a class of objects.

3 Models

In this section, we describe multiple different approaches to addressing the problem of how to generate novel but realistic textures for variant 3D meshes of an object class. We discuss the advantages and disadvantages of these methods. In summary:

- MODEL-BASELINE: This model utilizes simple UV mapping. We found that this architecture doesn't generalize to unseen 3D meshes.
- MODEL-UV: This model takes UV layout as extra input information but suffers similar problems to MODEL-BASELINE.
- MODEL-DEFORMABLE: This method is based on the idea of deformable models. The disadvantage of this method is that the approximated mesh is of low quality compared to the original 3D mesh.
- MODEL-GRAPH: This model consists of two variants, MODEL-GCN, and MODEL-GGAN. These models transform an input 3D mesh into a graph and generate a texture conditioned on the graph. The final variant of this model, called MODEL-GGAN, produced higher quality and more diverse results than the other previous methods.

In the following section, we describe each model in detail.

3.1 Model-Baseline

In this first model, the texturing is performed using the TEXTUREUV mechanism. The framework is summarized in Figure 3. For this model, we adapted a deep convolutional generative adversarial network (DCGAN) architecture [33] for the generator and the discriminator. We modified the DCGAN architecture to support the generation of higher-resolution textures. For applying the texture map to these diverse 3D models, we need a way to map the 3D vertices in these models into a 2D texture map (a procedure called UV mapping). We use the *smart UV project* feature of the Blender Python API to automatically generate these mappings for a given 3D model. Recall that our challenge is to adapt the GAN framework so that we can take advantage of training on multiple real-world examples and have the learned generation capability transfer to new 3D objects. However, the UV image for each 3D model has a different coordinate system. This means that, for example, the features of the given 3D model, such as the



tires or windshields of different cars, project to different regions in the UV space for each 3D model.

Fig. 3: Initial architecture for texture synthesis

To establish a baseline, we directly applied the UV mapping and found that this led the generator to converge to a mean texture across examples, rather than learning how to produce varied textures constrained by the input 3D mesh. This happens because the generator has no information about the UV mapping function and 3D mesh model to which the texture will be applied.

3.2 Model-UV

To address the above issue, we next explored the idea of injecting the UV map layout into the generator with the hypothesis that the generator might adapt during training in order to learn where the different parts of the given 3D model project in the 2D texture image. We used the same texturing mechanism, TEX-TUREUV, as described in the above MODEL-BASELINE. The architecture produced results similar to MODEL-BASELINE. The reason was the model couldn't learn complicated UV mapping information just from the UV layout image.

3.3 Model-Deformable

In this model, we use the idea of deformable models similar to Figure 2, but with some modifications. We use the same TEXTUREUV method described above. The

main problem with the above two models is that the UV mapping function is different for different 3D mesh models, making it harder for the generator to learn features that can work across different 3D mesh models. To mitigate this problem, we explored the idea of starting from a common mesh model with a fixed UV mapping. We used a sphere template 3D model as the starting point and used azimuth and elevation as the UV map function. We then used 3D chamfer loss [7] to predict the deformations of the sphere template to approximate the 3D mesh. This model more directly addresses our overall challenge by learning a generalized mapping from the space of textures to different 3D meshes enabling us to swap the texture learned from one model to another. However, the generated textures must still be applied to the approximate model, which reduces the quality of the 3D mesh model and the texture. The distortions in the model shape mesh are significant as demonstrated in the example in Figure 6. Another disadvantage of this method is that we need to approximate the deformation for every new 3D model, creating extra computational overhead for training and inference.

3.4 Model-Graph

8



Fig. 4: Graph-based methods for texture synthesis

In this section, we describe our architecture that incorporates information from the 3D mesh model to guide the generator. We first convert the input 3D mesh model into a graph by taking each face as a node in the graph and connecting neighboring faces using graph edges. The graph neural network is then used to learn a latent representation of the structural components of the given 3D model. In the latent representation, the topological features of the 3D mesh graph can be clustered, so as to learn features that could correspond to structural components that tend to share texture properties, such as wheels, windows, lights, and hood of a car. In turn, this latent component representation can then provide an inductive bias for the generator to produce a texture for the given 3D mesh model. The architecture is shown in Figure 4. An interesting aspect of this design is that the generator can take the unsupervised latent part representation as node features and combine it with the input noise vector to generate a texture for the particular 3D mesh. Node features is a 2D tensor of shape v, f where v is the number of nodes and f is the dimension of the node feature. The noise is a 1D vector of shape d. We sample a noise vector $z \in \mathbb{R}^d$ from a multivariate normal distribution. This d dimensional noise vector is then replicated to have a shape of v*d. This allows our model to process 3D mesh models with a different number of nodes. This noise tensor is then concatenated with the node feature tensor v * f. The concatenated tensor is then input to the generator (e.g. MLP) that, in turn, generates the textures for the given 3D mesh model. We use TEXTUREFACE for texturing as it enables us to generate higher quality textures than TEXTURE-VERTEX. We represent the faces of the 3D mesh as the nodes in the graph. The x, y, z face position and its normal (n_x, n_y, n_z) form the initial node features. We use a graph convolutional neural network (GCN) [23] to learn the latent part representation as shown in Figure 4. The generator generates a tensor of shape F * 3, where F is the number of faces (nodes in the graph), and 3 represents the three r, g, b colors (texture) per face. We use TEXTUREFACE for texturing the 3D mesh. We explored the following two variants of this MODEL-GRAPH that differ only in the design of the generator. MODEL-GCN uses GCN [23] as a generator to generate the texture from a combination of latent part representation and noise vector. And MODEL-GGAN uses a multi-layer perceptron with residual connections [13] as a generator. MODEL-GGAN is our best-performing model.

4 Experiments

We use the ShapeNet [3] car data set for all of our experiments. This data set consists of a total of 3,514 3D mesh models of cars with textures. We use 3314 mesh models for training, 100 mesh models for validation, and 100 mesh models for testing. The features extracted from intermediate layers of the pre-trained deep neural networks are known to correspond to the perceptual metrics of human vision [16, 40]. We found that incorporating this perceptual loss into the generative adversarial loss improved the qualitative appearance of the generated textures. Thus our overall loss function is as follows:

$$Loss(L) = gan_loss + \lambda * perceptual_loss$$
(3)

10 Dharma KC, Clayton T. Morrison, and Bradley Walls

We use the validation dataset to select the best value of λ . We use the library of Zhang et al. [40] to extract features from the intermediate layer of pre-trained AlexNet architecture [24] that are in turn used to calculate the perceptual loss. All of the models are trained with the above loss function. At each minibatch iteration, we render a 3D mesh with real and synthetic textures from eight different viewpoints. The loss is calculated from these real and synthetic ("fake") images. We use a learning rate of 0.0001 for both the generator and the discriminator. We use a hidden size of 64 for both GNN and the MLP generator. Here, hidden size is the dimension that's being used to project the node features of the graph. We render images of size 512×512 from the differentiable renderer. We use the Adam optimizer [22] for training all of our models. We used d = 16 for the random noise vector. For MODEL-GRAPH variants, we use 3 graph convolution layers [23] for learning the latent part representation. Each convolutional layer has a hidden size of 64. The noise vector is concatenated to the output of the last graph convolutional layer. We use TEXTUREFACE to texture the 3D mesh model. For the MODEL-GCN, we use a 7-layered GCN [23] as a generator with a hidden size of 64. The generator does not use residual connections [13]. For the MODEL-GGAN architecture, we use a 7-layered MLP as a generator with a hidden size of 64. The generator uses residual connections [13].

5 Results

The MODEL-BASELINE and MODEL-UV were only able to learn to generate textures for a single mesh model and were not able to generalize across unseen 3D mesh models. The MODEL-DEFORMABLE was able to generate a texture for unseen 3D mesh models, but the quality of approximated mesh and texture was not good, as demonstrated in contrast between Figures 5 and 6. Moreover, the approximation of the 3D mesh model created extra computational overhead. Thus we didn't move forward with this approach for the full ShapeNet car experiment. The graph-based models based on the MODEL-GRAPH architecture were



Fig. 5: Original mesh with original texture

Fig. 6: approximate mesh with learned texture using MODEL-DEFORMABLE

able to learn textures across different 3D mesh models. This general approach has multiple advantages compared to existing solutions that generate textures in 2D. First, the model is able to learn about the parts of the given 3D mesh model in an unsupervised way. This removes the effort and cost required for manual labeling of the 3D part segmentation. Second, the approach generates textures in 3D, so that a texture can be applied once and the 3D model can be viewed from multiple directions and under multiple light conditions without a need to generate texture each time we change these parameters. We evaluated these models by applying the synthetic texture generated from respective models and rendering them from multiple viewpoints. We then calculated the FID score based on these projected images and actual original images rendered from the same views with the original texture. The MODEL-BASELINE, MODEL-UV, and MODEL-DEFORMABLE architectures were not suitable for learning textures across different 3D mesh models, so we did not compute FID scores for these. Table 1 shows the average FID values (lower is better) for different models per 3D mesh model. For further comparison, we also experimented with a simple variant of the NeRF [28] model as a generator (MODEL-NERF), but it did not produce results as good as the graph neural network approaches. The low quality of results is reasonable because it doesn't have a way for learning part information like our MODEL-GRAPH.

Table 1: FID scores on the test dataset

l	Model	FID
[Model-GCN	0.75
Ì	Model-nerf	0.93
Ĩ	Model-GGAN	0.70

Figure 7 shows a set of selected examples of rendered images generated from different models with a fixed viewpoint. Textures are applied to the 3D mesh models and rendered as projected 2D images for visualization. The first column shows the images rendered with original textures, the second column shows the images rendered with textures generated from the MODEL-GGAN model, and the third column shows the images rendered with textures generated from the MODEL-GCN. In Figure 7, we observe that the MODEL-GCN lacks diversity in the generated images: it produces images with the same texture for every random noise input. We hypothesize that this is due to the over-smoothing problem observed in GNNs [2, 4] as the model uses GNN-only layers for the generator. Finally, the model MODEL-GGAN (GGAN) produces images (second column, Figure 7) that respect the object boundaries, are visually better than other models and are diverse (the model produces new textures on each run with different random noise input). Some of the images generated from our final model MODEL-GGAN (third row, second column) look even better than the original image itself (third row, first column).

6 Conclusion

In this work, we have presented and evaluated the graph generative adversarial network (GGAN), a new architecture that can learn to generate a texture for a given 3D mesh with high fidelity and that can learn 3D part information in an unsupervised way. We think GGAN will be useful in various domains to generate graph-structured representation. However, there are multiple directions



Fig. 7: First column: original images, second column: MODEL-GGAN, third column: MODEL-GCN

for improvement. The first important research direction for future work is to introduce symmetry constraints on the system such that all components with symmetrical structures will generate textures that respect symmetries. Second, we want to increase the diversity of the generated textures. Third, we want to improve the controlled synthesis of part-specific textures. Another important research direction would be to incorporate encoder-decoder graph architectures [8] within our framework. Another important direction would be to couple with a semi-supervised labeling approach. Finally, we would like to explore the use of flow-based models [36, 37] and diffusion models [6, 38] for the generation of texture within our current framework.

Bibliography

- Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein gan. arxiv 2017. arXiv preprint arXiv:1701.07875 30, 4 (2017)
- [2] Cai, C., Wang, Y.: A note on over-smoothing for graph neural networks. arXiv preprint arXiv:2006.13318 (2020)
- [3] Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., et al.: Shapenet: An informationrich 3d model repository. arXiv preprint arXiv:1512.03012 (2015)
- [4] Chen, D., Lin, Y., Li, W., Li, P., Zhou, J., Sun, X.: Measuring and relieving the over-smoothing problem for graph neural networks from the topological view. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 3438–3445 (2020)
- [5] Chen, W., Ling, H., Gao, J., Smith, E., Lehtinen, J., Jacobson, A., Fidler, S.: Learning to predict 3d objects with an interpolation-based differentiable renderer. Advances in Neural Information Processing Systems 32 (2019)
- [6] Dhariwal, P., Nichol, A.: Diffusion models beat gans on image synthesis. Advances in Neural Information Processing Systems 34, 8780–8794 (2021)
- [7] Fan, H., Su, H., Guibas, L.J.: A point set generation network for 3d object reconstruction from a single image. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 605–613 (2017)
- [8] Gao, H., Ji, S.: Graph u-nets. In: international conference on machine learning. pp. 2083–2092. PMLR (2019)
- [9] Gao, L., Wu, T., Yuan, Y.J., Lin, M.X., Lai, Y.K., Zhang, H.: Tm-net: Deep generative networks for textured meshes. ACM Transactions on Graphics (TOG) 40(6), 1–15 (2021)
- [10] Goel, S., Kanazawa, A., Malik, J.: Shape and viewpoint without keypoints. In: European Conference on Computer Vision. pp. 88–104. Springer (2020)
- [11] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. Advances in neural information processing systems 27 (2014)
- [12] Hamilton, W.L.: Graph representation learning. Synthesis Lectures on Artificial Intelligence and Machine Learning 14(3), 1–159
- [13] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
- [14] Henderson, P., Tsiminaki, V., Lampert, C.H.: Leveraging 2d data to learn textured 3d mesh generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7498–7507 (2020)
- [15] Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in neural information processing systems **30** (2017)

- 14 Dharma KC, Clayton T. Morrison, and Bradley Walls
- [16] Johnson, J., Alahi, A., Fei-Fei, L.: Perceptual losses for real-time style transfer and super-resolution. In: European conference on computer vision. pp. 694–711. Springer (2016)
- [17] Kanazawa, A., Tulsiani, S., Efros, A.A., Malik, J.: Learning category-specific mesh reconstruction from image collections. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 371–386 (2018)
- [18] Karras, T., Aittala, M., Laine, S., Härkönen, E., Hellsten, J., Lehtinen, J., Aila, T.: Alias-free generative adversarial networks. Advances in Neural Information Processing Systems 34 (2021)
- [19] Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 4401–4410 (2019)
- [20] Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and improving the image quality of stylegan. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 8110–8119 (2020)
- [21] Kato, H., Ushiku, Y., Harada, T.: Neural 3d mesh renderer. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3907–3916 (2018)
- [22] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
- [23] Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. arXiv preprint arXiv:1609.02907 (2016)
- [24] Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. Advances in neural information processing systems 25 (2012)
- [25] Li, T.M., Aittala, M., Durand, F., Lehtinen, J.: Differentiable monte carlo ray tracing through edge sampling. ACM Transactions on Graphics (TOG) 37(6), 1–11 (2018)
- [26] Liu, S., Li, T., Chen, W., Li, H.: Soft rasterizer: A differentiable renderer for image-based 3d reasoning. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 7708–7717 (2019)
- [27] Loper, M.M., Black, M.J.: Opendr: An approximate differentiable renderer. In: European Conference on Computer Vision. pp. 154–169. Springer (2014)
- [28] Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: European conference on computer vision. pp. 405–421. Springer (2020)
- [29] Mirza, M., Osindero, S.: Conditional generative adversarial nets. arXiv preprint arXiv:1411.1784 (2014)
- [30] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. Advances in neural information processing systems **32** (2019)
- [31] Pavllo, D., Kohler, J., Hofmann, T., Lucchi, A.: Learning generative models of textured 3d meshes from real-world images. In: Proceedings of the

IEEE/CVF International Conference on Computer Vision. pp. 13879–13889 (2021)

- [32] Pavllo, D., Spinks, G., Hofmann, T., Moens, M.F., Lucchi, A.: Convolutional generation of textured 3d meshes. Advances in Neural Information Processing Systems 33, 870–882 (2020)
- [33] Radford, A., Metz, L., Chintala, S.: Unsupervised representation learning with deep convolutional generative adversarial networks. arXiv preprint arXiv:1511.06434 (2015)
- [34] Raj, A., Ham, C., Barnes, C., Kim, V., Lu, J., Hays, J.: Learning to generate textures on 3d meshes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops. pp. 32–38 (2019)
- [35] Ravi, N., Reizenstein, J., Novotny, D., Gordon, T., Lo, W.Y., Johnson, J., Gkioxari, G.: Accelerating 3d deep learning with pytorch3d. arXiv preprint arXiv:2007.08501 (2020)
- [36] Rezende, D., Mohamed, S.: Variational inference with normalizing flows. In: International conference on machine learning. pp. 1530–1538. PMLR (2015)
- [37] Weng, L.: Flow-based deep generative models. lilianweng.github.io (2018), https://lilianweng.github.io/posts/2018-10-13-flow-models/
- [38] Weng, L.: What are diffusion models? lilianweng.github.io (2021), https://lilianweng.github.io/posts/2021-07-11-diffusion-models/
- [39] Xian, W., Sangkloy, P., Agrawal, V., Raj, A., Lu, J., Fang, C., Yu, F., Hays, J.: Texturegan: Controlling deep image synthesis with texture patches. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 8456–8465 (2018)
- [40] Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 586–595 (2018)
- [41] Zhang, Y., Chen, W., Ling, H., Gao, J., Zhang, Y., Torralba, A., Fidler, S.: Image gans meet differentiable rendering for inverse graphics and interpretable 3d neural rendering. arXiv preprint arXiv:2010.09125 (2020)