

000  
001  
002  
003  
004  
005  
006  
007054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

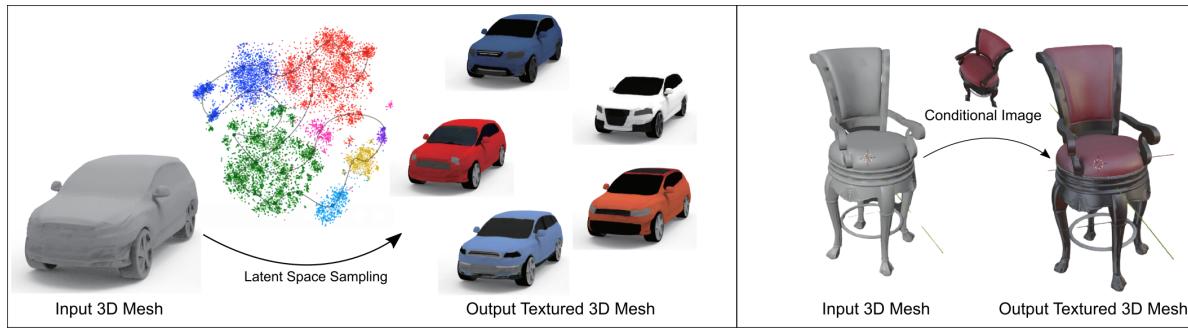
# VJAM Mesh Texture: Learning continuous texture representation for conditional and unconditional texture synthesis of 3D meshes

008 Vikas Thamizharasan

009 Joshua Pierce

010 Angelina Grosso

011 Moyi Tian



## Abstract

We endeavor to design an architecture that can learn detailed textures across a latent space of images and shapes, thus capable of probabilistically generating realistic, novel textures for previously unseen meshes. Taking inspiration from the recent work *TextureFields* by Oechsle et al. [10] and others, we explore varying methods of shape-representations, alternative loss metrics, and high-frequency learning techniques. We demonstrate that applying Fourier feature transformations through a positional encoding is an effective means for learning more detailed textures.

## 1. Introduction

Generative models over the past few years have exploded in popularity and become a hallmark of research problems within the vision community since the seminal work of Goodfellow et al. [4] and Jürgen Schmidhuber [13]. Deep learning models have become a cornerstone of image-based tasks such as image-to-image translation. However adapting these models to 3D data is not trivial as naive implementations can grow in space and time complexity. Unlike images which are regular, 3D data have numerous representations, each with their pros and cons. Fortunately over the past couple of years, there have been numerous works designing learning-based techniques for 3D geometry representations and that take into consideration

the trade-offs across fidelity, efficiency and generalizability. Some of these works include PointNet[12], AtlasNet[5], DeepSDF[11], Occupancy Net[8] and MeshCNN[7]. They have shown promising results in various tasks such as single image 3D reconstruction, surface reconstruction from partial and noisy 3D data, interpolation, generative modeling and learning high quality shape representation.

Learning texture representations has been a less explored problem. Like in the case of geometry, how we represent our input textures is not obvious. This is important, not just for deciding what kind of data should be collected for training, but also informing the design of the architecture and the representations learned by the neural network. If we attempt to leverage the regular nature of voxels data, we run into issues such as: the discretization of voxels presents inherent limits in expressing high frequency, sharp details; furthermore, voxels are very memory inefficient, increasing by a cubic factor. Some of the above problems can be resolved by using techniques such as texture atlas, which provides a mapping between 3D data and UV space (2D), or using Geometry Images [6]. While these representations alleviate some of the concerns of using voxels, they come with substantial distortion associated with orientation as well as discontinuities at boundaries. Furthermore, convolutions on such maps would make learning canonical features difficult due to having irregular and randomly distorted receptive fields.

With this motivation, our Deep Learning Project builds from a recent work called *TextureFields* by Oechsle et al. [10]. *TextureFields* seeks to solve the aforemen-

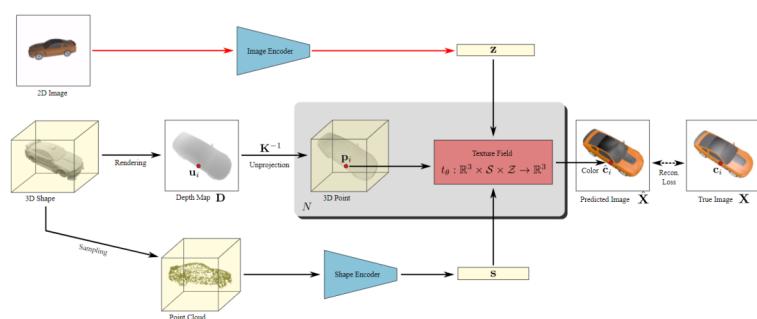


Figure 1. TextureFields's Conditional Model. Source : [10]

tioned problems by using neural networks to approximate a parameterized continuous function for representing texture information in 3D space. In this report, we first review the framework proposed in TextureFields and subsequently explore some ideas to improve certain aspects of this paradigm. Ultimately, we seek to train a network that produces more meaningful latent representations for appearances and shapes, and accordingly generating more realistic textures.

In our project, we attempt to incorporate:

- More state-of-the-art techniques for learning shape representations (like DeepSDF, OccupancyNets, PointNet++ [12]) by either replacing the TextureFields shape encoder entirely, or by initializing the shape encoder weights with a PointNet encoder pre-trained on a geometry-specific task with stronger priors.
- Methods for approximating high-frequency data, such as mapping inputs to Fourier features via a 'positional encoding' [15]. Ideally, this would afford TextureFields a more detailed expression of textures from the latent space.
- Alternative loss functions such as the learned Perceptual Image Patch Similarity [16]. The goal would be to clarify blurry textures generated with L1 loss, which is more sensitive to large perceptual changes
- Additional experiments with ShapeNet [2] and the 3D Future dataset [3], with corresponding pre-processing scripts.

## 2. Background

### 2.1. TextureFields

The authors define a function  $t$  parameterized by a neural network  $t_\theta$ , with learnable parameters  $\theta$ , which maps an input 3D point  $p$  conditioned on shape embedding  $s \in S$  (calculated from the input geometry) and latent variable  $z \in Z$

(derived from either a conditioned image or sampled from a learned distribution), to produce an output colour  $c$ :

$$t_\theta : \mathbb{R}^3 \times S \times Z \longrightarrow \mathbb{R}^3 \quad (1)$$

As the idea of texturing a mesh is ill-posed, the authors choose to condition TextureFields with a shape embedding  $s$ : generated from a shape encoder that ideally captures contextual geometric information about the input shape; and also a viewpoint invariant global-feature representation  $z$ : encoded from an image that constraints this task by providing information about the mesh appearance.

TextureFields, built around a shape encoder and the actual texture fields architecture, implements a total of three different models: a conditional model, a generative adversarial network and a variational auto-encoder to capture ambiguity in a random latent code  $z$ . In this work we focus on the conditional model which is illustrated in figure [1]. This model takes as input both a mesh (with points and surface normals uniformly sampled) and a conditional 2D image, and outputs a textured mesh modelled after the conditioned image. The conditional model is composed of a shape encoder that generates the shape embedding, an image encoder that generates the latent variable and the TextureField model itself.

### 2.2. Occupancy Nets

Extracting shape-information is critical to the task of TextureField. Occupancy Networks implicitly represent the 3D surface of an object as the continuous decision boundary of a deep neural network classifier. In later sections we use a pretrained occupancy net to initialize the shape embedding of texturefields.

### 2.3. PointNet++

Pointnet is another method for extracting shape information: specifically, offering a permutation-invariant approach to analyzing point cloud data and thus allowing one to directly extract geometric information with deep learning



Figure 2. Samples from the chair class of 3D FUTURE.

techniques. Pointnet++ offers notable improvements over its predecessor by introducing a method for applying convolutions to the point cloud data. While Pointnet relied on aggregating individual features of points, Pointnet++ also incorporates the local structures, facilitating markedly better shape representations.

#### 2.4. Losses

TextureFields conditional model uses  $\ell_1$ -loss between the predicted and rendered images to train the model. Such setting is commonly used in network training, but it cannot efficiently capture the perceptual similarity in image features. In recent years, a network-trained loss function, perceptual losses, is proposed as a measure more-aligned with human perceptual judgement by better capturing low-level perceptual similarity across the latent features of images. Studies have shown that perceptual losses excel in further fine-tuning the learned features of a pre-trained network.

#### 2.5. Positional Encoding

Recent research has made notable strides in enhancing neural networks' ability to approximate high-frequency data, such as Sitzmann's Sinusoidal Representation Networks[14] and Tancik's Fourier Feature Networks[15]. One notable application of the latter method was NeRF's[9] learning of photo-realistic, 3D scene representations. With this as inspiration, we explore whether these techniques could afford TextureFields a more granular, detailed expression of textures from its learned latent space.

### 3. Methodology

#### 3.1. Data

We use data from the ShapeNet dataset [2] and the 3D FUTURE Dataset [3]. ShapeNet is a popular large scale 3D repository covering various object classes. In this work, we use the car category along with the textures. 3D FUTURE was developed by professional designers and contains high quality 3D instances of furniture with high resolution textures. This dataset contains both synthetic scenes and 3D



Figure 3. Samples from the car class of ShapeNet.

instances of furniture, but for our purposes we will only be dealing with 3D instances of furniture, of which there are 9,992 instances across a total of 34 different categories. These 34 categories can be divided into broad categories, such as "bed", "cabinet" or "chair". Looking at figure [2] and [3], it can be seen that 3D FUTURE contains objects with more complex topology and higher resolution textures containing fine details.

TextureFields requires as input: uniformly sampled points of the 3d mesh and the corresponding surface normals at each sampled point,  $N$  depth maps and albedo images rendered from various camera angles along the viewing hemisphere, the camera intrinsic and camera extrinsic for each of the  $N$  camera poses in order to unproject the depth pixel to a 3d point, and optionally a set of conditional images of objects similar to the input mesh. TextureFields provides us with the pre-processed data for ShapeNet, but not the corresponding script that accomplishes this. We thus write our own script to be able to generate data for 3D FUTURE that adheres to the requirements of TextureFields.

#### 3.2. Adjusting the Shape Encoding Architecture

As 2D image encoders are a well-studied field, we attempted to improve the performance of the conditional model by modifying the shape encoder.

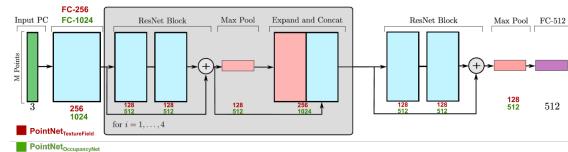


Figure 4. Shape encoder used by TextureFields, Source : [10]

We hypothesized that the shape embedding outputted by the shape encoder should ideally contain pertinent topological information about the mesh, such as local connectivity, symmetry, and part decomposition. We conjecture that having a latent vector which is more topologically aware would benefit the overall output of the model. The shape encoder originally implemented by the authors was inspired by Pointnet and uses ResNet blocks; it takes point clouds as

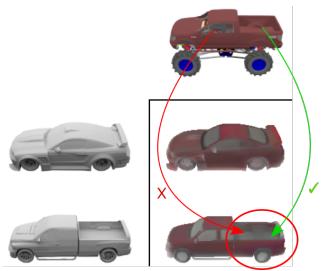
324  
325  
326  
327  
328  
329  
330  
331  
332

Figure 5. An example of incorrect texture prediction. While the example of the input car looks correct, the second example of the pickup truck has an instance of incorrect texture inferred from the conditioned image. Ideally, the shape embedding coupled with the encoded image representation should inform TextureFields of semantic-level information such as windows, chassis, trunk etc. But here we see texture for the trunk being drawn from the windows as opposed to the trunk visible in the conditioned image.

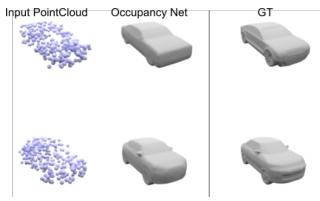
341  
342  
343  
344  
345  
346  
347  
348  
349  
350

Figure 6. Occupancy net trained on surface reconstruction.

351  
352  
353  
354  
355  
356  
357  
358  
359

inputs and outputs a 512 dimensional vector. This PointNet encoder was trained from scratch within the pipeline shown in figure [1], where all the losses are only computed on the output of the TextureFields model. This leads to some uncertainty as to what degree the shape encoder encodes meaningful information as described above into the low dimensional shape embedding. Some of the results shown in figure [5], makes us believe there is room for improvement.

We propose some of the following changes:

360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377

- Initialize the PointNet encoder with pre-trained weights. The idea behind this is to use weights from a model that was learned for a 3D-geometry task such as reconstructing surfaces from pointclouds. With such a task, the network is able to better utilize priors and is penalized for not capturing pertinent topological information.

This goal is never fully-realized in TextureFields and thus motivates us to employ pre-trained weights, and then subsequently fine tune the PointNet encoder within TextureFields, thus giving the network a head start and hopefully helping it converge to an optimal minima.

- We replaced this shape encoder by a Pointnet++ encoder.

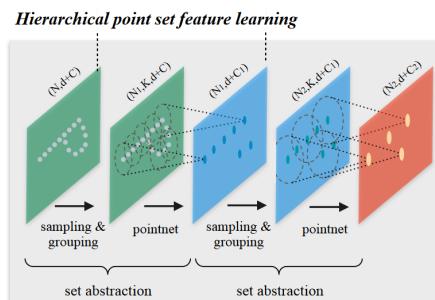


Figure 7.

## Training and Evaluation

Ideally, the changes we made to the TextureFields architecture would be evaluated by training the model multiple times on a large dataset. One training run would use the original TextureFields architecture, while the other training runs would incorporate our various modifications. This experimental setup would allow for results that could clearly evaluate the impact of our modifications on performance. After reaching out to the authors of TextureFields, we learned that training their model on the ShapeNet car dataset took approximately one week. As this was an infeasible task for our project timeline, we modified our original training and evaluation plan.

Instead of training on a large dataset, we decided to overfit on single examples. Using a single example as the input into the model, we trained the model until the loss value appeared to converge and the output did not appear to change between epochs. We then used a mixture of qualitative and quantitative criteria to evaluate the performances of the different architecture. We essentially test the representation power of the network.

Similarly, to compare the performances of the standard  $\ell_1$ -loss and the perceptual losses, we perform the representation power experiments on both  $\ell_1$ -loss and perceptual loss [16]. Specifically, the perceptual loss function we use is based on the VGG network, initialized from a pre-trained classification model.

### 3.3. Positional Encoding

Recent research has proposed effective techniques for enabling neural networks to approximate high-frequency data, such as SIREN[14] and Fourier feature mapping[15]. The latter uses a positional encoding technique: specifically, mapping the inputs of a neural network to higher-dimensional Fourier features. This technique has presented demonstrable improvements in learning high-frequency inputs such as in the case of NeRF's photorealistic, three-dimensional scenes[9]. Given the conceptual similarity between TextureFields and NeRF's learned mapping of colors

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

432 in a 3D space, we chose to implement NeRF’s positional  
 433 encoding method. Furthermore, this approach is limited to  
 434 mapping the input coordinates to higher dimensions, thus  
 435 allowing us to avoid substantially modifying the Texture-  
 436 Fields architecture.  
 437

438 Recall that TextureFields takes as input a shape-vector,  
 439 image-vector, and a 3D coordinate on the provided mesh,  
 440 and outputs a corresponding color. Following the imple-  
 441 mentation in NeRF, our positional encoding technique maps  
 442 each dimension of the 3D coordinate into higher dimen-  
 443 sional space with a Fourier features transformation.  
 444

$$\begin{aligned} \gamma(v) = & (a_0 \sin(\pi b_0^T v), a_0 \cos(\pi b_0^T v), \\ & a_1 \sin(\pi b_1^T v), a_1 \cos(\pi b_1^T v), \dots, \\ & a_L \sin(\pi b_L^T v), a_L \cos(\pi b_L^T v)) \end{aligned} \quad (2)$$

451 Note that each input coordinate is mapped to multiple  
 452 pairs of sin and cos transformations- the number of such  
 453 pairs is represented by  $L$  in the above equation. Therefore,  
 454 the resulting transformation maps the 3D input coordinates  
 455 to  $3 * L * 2$  dimensions.  $L$  is a hyperparameter, and from  
 456 the applications outlined in NeRF[9] and Fourier Features  
 457 Networks[?] we chose  $L = 8$ , thus mapping TextureFields’  
 458 3D input to a 48D vector.  
 459

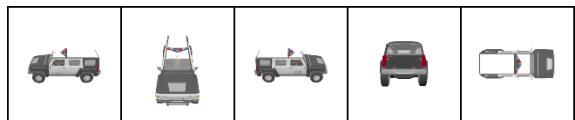
## 460 Training and Evaluating

461 As noted above, fully training an updated TextureFields  
 462 model requires a prohibitive amount of training time. Fur-  
 463 thermore, in this case it’s not sufficient to overfit to a single  
 464 texture because we want to evaluate both the positional  
 465 encoding’s effect on learning detailed textures from the train-  
 466 ing set as well as the impact on synthesizing novel textures  
 467 from new shape- and image-vectors. In balancing these  
 468 goals we chose to train a TextureFields model modified with  
 469 a positional encoding for 30k epochs on five car models.  
 470 For comparison, we also trained the original TextureFields  
 471 model for 29.5k epochs on the same five cars. This re-  
 472 quired 10.5 hrs and 7.5 hrs of training on a single NVIDIA  
 473 Tesla V100 GPU for the normal and positional encoding  
 474 models, respectively. Comparison of the two models con-  
 475 sists of a qualitative assessment of their generated images  
 476 as well as an evaluation of image similarity metrics calcu-  
 477 lated from predicted views and their corresponding ground  
 478 truths. Similar to TextureFields, we use the structure sim-  
 479 ilarity image metric (SSIM) as a measure of local image-  
 480 properties and their Feature- $\ell_1$ -metric as a more global per-  
 481 ceptual measure. Generated images and corresponding im-  
 482 age similarity metrics are calculated at every 500 epochs.  
 483

## 484 4. Results

### 485 4.1. Shape Encoders

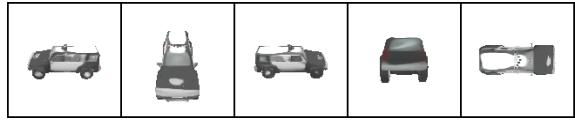
486 In order to evaluate the performance of our modifications  
 487 to the shape encoder architecture, we overfit the model on  
 488 a single example from the shape net dataset. We ran the  
 489 model a total of three time: once as a baseline, once with  
 490 Pointnet++ used as the shape encoder, and once with the  
 491 original shape encoder pre-trained with weights from Oc-  
 492 cupancy Network. Quantitatively, we observe the same two  
 493 model evaluation metrics use by Texture Fields: structure  
 494 similarity image metric (SSIM) and feature  $\ell_1$ -metric.  
 495



496 Figure 8. Ground-truth images for the training



497 Figure 9. Baseline predicted images



498 Figure 10. Pointnet++ predicted images



499 Figure 11. Pretrained Weights predicted images

500 For the baseline, we used the original architecture pro-  
 501 vided by Texture Fields and ran for a total of approximately  
 502 70,000 epochs. Qualitatively, the final result is very close  
 503 to the ground truth which is to be expected when overfitting  
 504 to a single model. Some small details remain unlearned,  
 505 such as the ground truth fact that the seats can be observed  
 506 through the rear window of the car.  
 507

508 We then ran the model again, this time for approximately  
 509 85,000 epochs, with Pointnet++ as the shape encoder. We  
 510 saw that qualitatively, the final result was further from the  
 511 ground truth than the baseline. Many details remained  
 512

540 blurred or smudged, particularly on the roof of the car. As  
 541 we ran this model for more epochs we received worse  
 542 results and concluded that changing the architecture to Point-  
 543 net++ hurt the performance of the model.  
 544

545 The third model used the original architecture provided  
 546 by TextureFields with the shape encoder initialized with  
 547 pre-trained weights from the Occupancy Network. We orig-  
 548 inally tried to initialize all weights with pre-trained weights.  
 549 This yielded a strange output of monochrome texture. At  
 550 first we believed there was a bug within the code we used to  
 551 load the pre-trained weights. Extensive debugging showed  
 552 that when different combinations of weights were initial-  
 553 ized, some would stay frozen at certain values while others  
 554 would change throughout training. We came to the con-  
 555 clusion that there was nothing wrong with our implemen-  
 556 tation. Our theory for this strange behavior is that the pre-  
 557 trained weights from Occupany Network cause the condi-  
 558 tional model to remain stuck in a very non-optimal minima.  
 559

560 To avoid getting stuck in such a minima, we decided  
 561 to only initialize some of the weights within the geo-  
 562 metry encoder. Specifically, all the weights within the first 2  
 563 ResNet blocks. Qualitatively, the results were very similar  
 564 to those from the baseline. The hypothesis here is that, the  
 565 first few layers learn features that are common between the  
 566 two tasks: texture synthesis and surface reconstruction.  
 567

568 Comparing the quantitative metrics, there is no signifi-  
 569 cant difference between the (SSIM) and (FID) of the three  
 570 runs.  
 571

572 Altogether, these results led us to tentatively conclude  
 573 that changing the geometry encoder to Pointnet++ harmed  
 574 the model's performance while initializing some weights  
 575 within the geometry encoder had no effect.  
 576

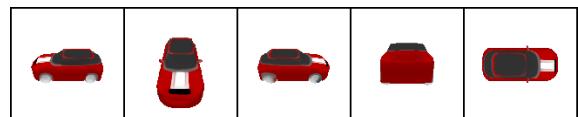
#### 577 4.2. $\ell_1$ -loss vs. Perceptual Loss

578 We trained on one car model from ShapeNet dataset once  
 579 with  $\ell_1$ -loss proposed by TextureFields and once with per-  
 580 ceptual loss function. Each of the trainings ran over ap-  
 581 proximately 29k epochs. Figure [25] shows the real images  
 582 from the actual car model. Figure [13] shows the predicted  
 583 images with respect to  $\ell_1$ -loss and Figure [14] shows the  
 584 predicted images with respect to the perceptual loss.  
 585



586 Figure 12. Ground-truth images for the training  
 587

588 From the generated fake images, we can qualitatively  
 589 conclude that the TextureFields network is unlikely to im-  
 590 prove performance by implementing the perceptual loss.  
 591 Specifically, we observe that some types of patterns in ap-  
 592 pearance, the white stripes on the hood for example, can  
 593



594 Figure 13. Predicted images after training 29k epochs with  $\ell_1$ -loss  
 595



596 Figure 14. Predicted images after training 29k epochs with percep-  
 597 tual loss  
 598

599 no longer be effectively captured after taking the perceptual  
 600 judgement part into the network. We also found that the pre-  
 601 dicted color tends to be brighter in the perceptual network.  
 602 One thing to notice is that the tail lights of the car, which  
 603 have long thin horizontal shapes, are better depicted in the  
 604 perceptual case than the  $\ell_1$ -loss case. We conjecture that  
 605 the perceptual loss is better at learning textures with pat-  
 606 tterns that are thin and long horizontally than the ones that  
 607 are more like vertical strips.  
 608

609 Quantitatively, we can look at the graphs of average loss  
 610 per epoch over training as depicted by Figure [15] and com-  
 611 pare the behaviors of the two networks. We observe that  
 612 perceptual losses result in larger variations and also are  
 613 higher than the  $\ell_1$ -loss on average. Hence, we can also con-  
 614 clude from the data that  $\ell_1$ -loss has better performance than  
 615 the perceptual losses in the task of learning appearance.  
 616

617 In general, the behavior of the network implemented  
 618 with the perceptual loss function does not match our predic-  
 619 tion that the network will be able to learn the texture infor-  
 620 mation better if we simply improve perceptual similarities  
 621 by capturing deep features on the image level. Since the net-  
 622 work with perceptual losses takes approximately double the  
 623 time as the  $\ell_1$ -loss to train for the same number of epochs,  
 624 in practice applying perceptual losses for capturing 3D tex-  
 625 ture is not a practical choice. However, the differences in  
 626 effectively captured features suggests a new potential direc-  
 627 tion to investigate. It is worth studying which features are  
 628 more advantageously learned by these networks, and how  
 629 they may correspond with specific loss functions. This may  
 630 provide a strong basis for strategically combining multiple  
 631 losses to improve performance.  
 632

#### 633 4.3. Positional Encoding

634 When qualitatively comparing the abilities of both mod-  
 635 els' to learn detailed representations of textures in the train-  
 636 ing set, the TextureFields model with the positional encod-  
 637 ing transformation (*PE*) learned more granular textures as  
 638 compared to the original TextureFields model (*ORIG*). Al-  
 639 though the models improved their texture representations  
 640

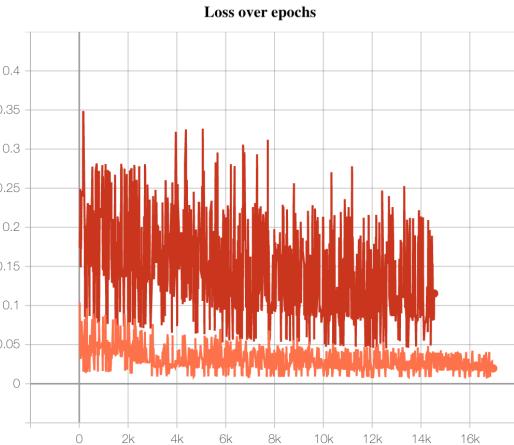
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663

Figure 15. Average loss per epoch (collected from the initial part of the training). Perceptual loss is in red and  $\ell_1$ -loss is in orange

over the course of training, both models occasionally generate abruptly distorted textures for one or more car models before adjusting in the subsequent epochs. (For examples see figure 24 in the appendix: epochs 24.5k and 29.5k for *PE* and *ORIG*, respectively.) The last epoch for *ORIG* generated substantially distorted textures for three of the car models. To avoid exaggerating the degree of *PE*'s improved accuracy relative to *ORIG*, in figure 16 we compare the textures generated from *PE*'s last epoch with the qualitatively 'best' images selected from across *ORIG*'s training history. Even with these *ORIG*-favorable terms, *PE* generates comparatively more detailed textures. For further context, figure 22 in the appendix compares the generated textures from the last epoch of both models for all five cars in the training set.

Additionally figures 23-25 in the appendix show the progression of generated images at an interval of 2.5k epochs. (These generated images are the ones from which *ORIG*'s 'best' results were selected in figure 16.) Figures 23-25 also highlight a distinction between the learning processes of *ORIG* and *PE*. In broad terms, details appeared to be learned more rapidly for the *PE* model. Specifically, note that even the initial texture representations of *PE* include some variability in allocating color, while initial *ORIG* representations are largely monochrome, reflecting the most prevalent color of each car. This comports with the functional purpose of positional encoding- ie, *PE* more readily learns to allocate high-frequency changes in color-values across relatively short distances.

The comparative speed with which *PE* approximates textures in the training set is evident when evaluating the training loss. See figure 17. The *PE* model also shows better SSIM and Feature- $\ell_1$ -metric values over the course of the 500-epoch reporting intervals, indicating higher-fidelity

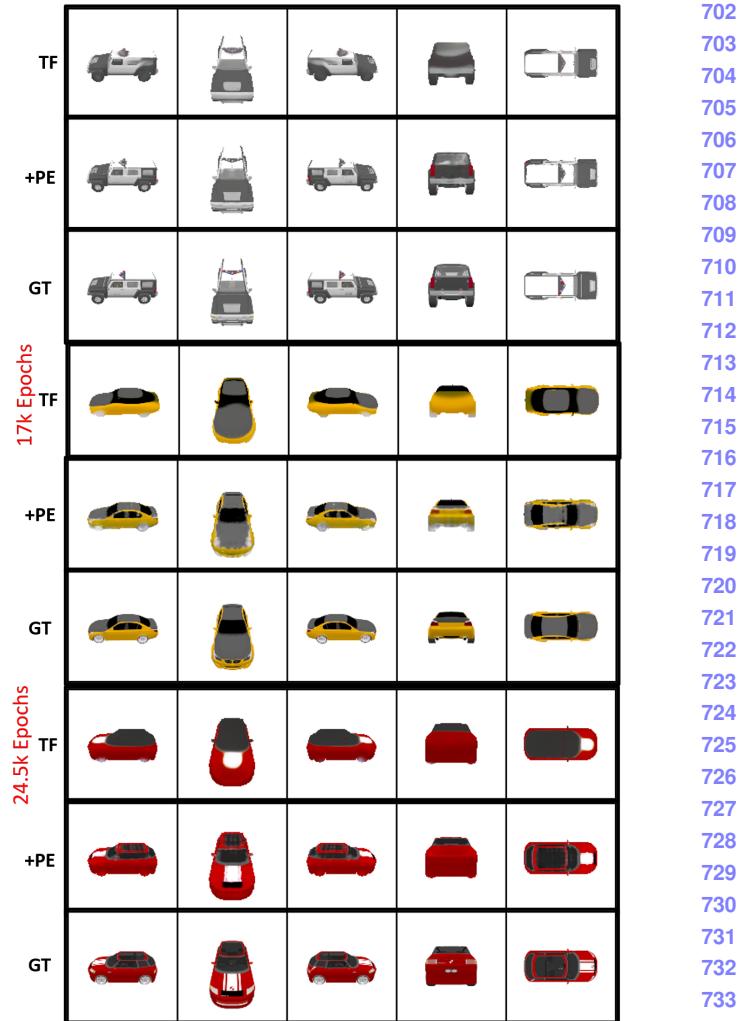


Figure 16. **Comparison of textures generated from three of the five cars in the training set.** The model with positional encoding (+PE) learned more detailed representations of the training textures. For example, note the red border around the roof of the bottom car; and on extremely close inspection, the white stripe on the hood can be seen to be segmented into three sections similar to the ground truth (GT). This model was trained for 29.5k epochs. The normal TextureFields model (TF) happened to generate substantially distorted textures for the red and yellow cars on the final 29.5k epoch. (This is a normal periodic occurrence for both models over the course of training- for reference the generated textures across training is visible in figures 23-25 in the appendix.) So as to not over-emphasize the difference between the two models, the above figure selects the best texture over the course of training for the TextureFields model, with the corresponding epoch indicated in red text. (Unless otherwise indicated, the 29.5k epoch is used in the above figure.) A comparison of generated textures from the 29.5k epoch for both models is visible in figure 22 of the appendix

representations of the training textures. See figures 18 and 19.

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

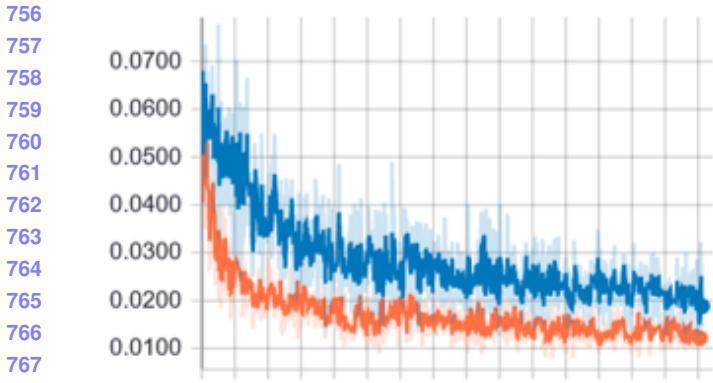


Figure 17. **Training loss over the course of 29.5k epochs.** The blue line represents the normal TextureFields model (*ORIG*) while orange indicates the model with positional encoding (*PE*). Note the steeper initial drop in loss for the *PE* model, as well as the diminished variation across epochs. The steeper slope in the loss curve comports with the comparatively earlier appearance of texture details during *PE*'s training (see figures 23-25 in the appendix)

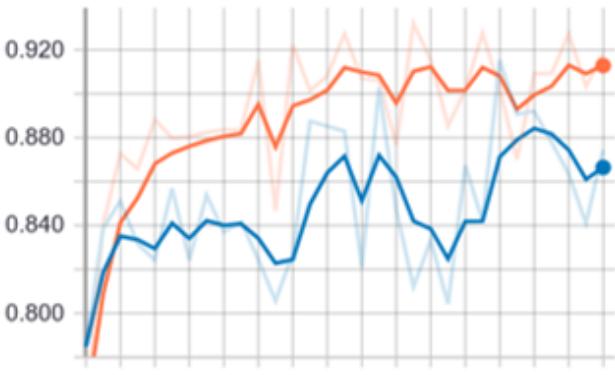


Figure 18. **Structure similarity image metric (SSIM) calculated at 500-epoch intervals over the course of 29.5k epochs.** The blue line represents the normal TextureFields model (*ORIG*) while orange indicates the model with positional encoding (*PE*). Higher values indicate more image similarity between between the model's predicted images and the ground truth for cars in the training set. As compared to the metric reported in figure 19, SSIM is a measure of local image-properties.

While *PE* performs better at learning textures from the training set, that does not necessarily equate to improved performance when synthesizing novel textures. Figures 20 and 21 show two examples where *PE* and *ORIG* generate novel textures after being given a new mesh and corresponding shape and image-vectors as inputs. Given the limited training data and training time this analysis is only useful for the sake of comparison between the two models; in other words, we wouldn't expect either to perform well with this task. In the first example, *ORIG* does a notably better job



Figure 19. **Feature- $\ell_1$ -metric calculated at 500-epoch intervals over the course of 29.5k epochs.** The blue line represents the normal TextureFields model (*ORIG*) while orange indicates the model with positional encoding (*PE*). This is a metric that was introduced by TextureFields, where higher values indicate less image similarity between between the model's predicted images and ground truth for cars in the training set. This metric was designed to capture more global features of image similarity.

of incorporating the image-vector into the texture, however *ORIG* completely misses the mark with the second example. While *PE* displays more texture-variance with its mottled allocation of colors, this does not clearly correspond to generating more realistic textures.

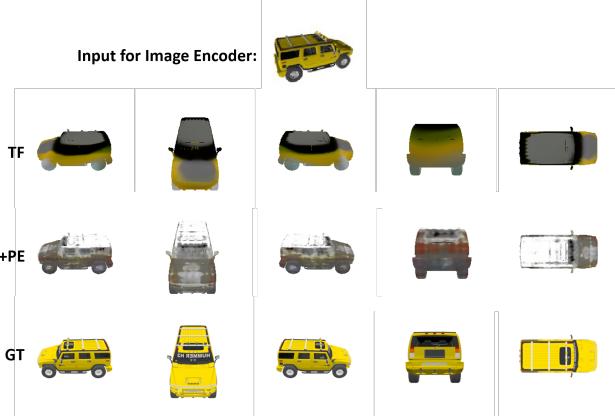


Figure 20. **Comparison of novel textures generated after training for 29.5k epochs on five car models.** The models are given as input the shape- and image-encodings from a previously unseen car, and then predict a corresponding texture. The image-encoding is calculated from a single view of the car, which is displayed at the top of the figure. As expected, neither model did particularly well at this task due to limited training; however, *ORIG* better approximated the color of the target texture, seeming to appropriately lift the coloring-scheme of the yellow-car from the training set. On extremely close inspection, it's apparent both models made slight attempts at the lettering on the windshield of the truck, with the *PE* model's version slightly more pronounced.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

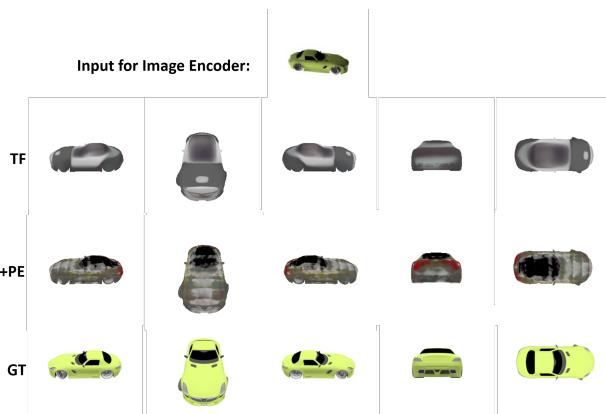
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877

Figure 21. **Comparison of novel textures generated after training for 29.5k epochs on five car models.** In this instance *ORIG* entirely missed the mark of the target coloring scheme. However, *PE* is not much better, though arguably closer to the slightly darker shade of green in the input-image. Similar to the previous figure, *PE*'s textures display higher-frequency variation in color, while *ORIG*'s textures are consistently more smooth. This fits with the design of the positional encoding, however at least in these examples this does not result in discernibly more realistic textures. Altogether, these two results fail to suggest that the improved granularity of the positional encoding in learning the training set will naturally translate to improved synthesis of novel textures.

## 5. Discussion

### 5.1. Shape Encoding

With regards to the shape encoder modifications, we had hoped to see loss converging faster and to a smaller value than the overfitting performance yielded by the baseline. Quantitatively we see this is not the case; the differences in loss convergence and metrics is not significant. Qualitatively we see that pretrained weights make no difference and that Pointnet++ seems to yield inferior results. We will need to put more thought into what sort of geometric details are captured by the shape embedding.

### 5.2. $\ell_1$ -loss vs. Perceptual Loss

With regards to changing the loss, quantitative results in different metrics yielded no significant differences when compared to the baseline, but the variances and loss values are in general higher than the original network. Through qualitative observations, we note that perceptual losses can capture some minor horizontal texture features slightly better. In contrast to our expectation, the experiment indicates that capturing deep features or perceptual similarities at image level does not contribute to learning essential texture information in 3D. One future direction is to investigate which features are better learned by which networks/loss functions, and how we can combine these advantages to achieve better performance.

### 5.3. Positional Encoding

In our experiment, the positional encoding transformation corresponded with a clear improvement in TextureField's ability to learn detailed textures from the training set. However, this did not result in a corresponding improvement in performance when synthesizing novel textures. Altogether, there are a couple of caveats to keep in mind with this experiment:

- The limitations in data and training prohibit us from conclusively evaluating the impact of the positional encoding on synthesizing novel textures. For example, it's conceivable that the improved performance of *PE* in learning textures in the training set may come at the cost of generalizing to new textures. In other words, it's possible the positional encoding as it's currently structured may present the inherent trade-off of overfitting: learning high-frequency details in the training set versus generalizing to novel textures. On the other hand, the failure in both these examples could very reasonably be the result of the minuscule training data.

- Given the limitations in training time, we cannot speak conclusively about the relative potential of these models in less-restrictive scenarios. For example, it's quite possible that *ORIG* would always eventually learn a similar level of granular textures as the positional encoding model given enough training time.

- The larger input vector of *PE* (48 dimensions instead of three) inherently provides considerably more trainable parameters in the first fully connected layer of TextureFields. It's conceivable that some amount of improved performance with regards to learning training textures is the result of simply having more trainable parameters. In subsequent studies this could be controlled for by increasing the size of *ORIG*'s initial, fully-connected layer.

None-the-less, the clear improvement in the positional encoding's ability to learn detailed textures suggests that this is a worthy avenue of further research. Future studies could include:

- Expanding training time. This will help investigate whether or not the original model is in fact capable of learning the same level of detail, albeit at a slower pace.
- Expanding the training data. This can help inform whether the positional encoding's approximation of high-frequency details could be applied to the task of novel texture synthesis after learning a more comprehensive latent space.

864		918
865		919
866		920
867		921
868		922
869		923
870		924
871		925
872		926
873		927
874		928
875		929
876		930
877		931
878		932
879		933
880		934
881		935
882		936
883		937
884		938
885		939
886		940
887		941
888		942
889		943
890		944
891		945
892		946
893		947
894		948
895		949
896		950
897		951
898		952
899		953
900		954
901		955
902		956
903		957
904		958
905		959
906		960
907		961
908		962
909		963
910		964
911		965
912		966
913		967
914		968
915		969
916		970
917		971

972 • Exploring other techniques for learning high-  
973 frequency data. NeRF’s positional encoding was  
974 partly selected for its simplicity in assimilating with  
975 TextureFields. Subsequent research with SIREN[14]  
976 and Fourier transformations[15] have demonstrated  
977 improved efficacy in learning high-frequency data. For  
978 example, recent research by Chan et al. [1] has applied  
979 neural radiance fields in the context of a generative  
980 adversarial network while incorporating the SIREN  
981 method for approximating high-frequency data. This  
982 research may provide relevant insight for how to better  
983 incorporate high-frequency techniques into Texture-  
984 Fields as it bridges the conceptual similarity of neural  
985 radiance fields’ implicit 3D scene-representation with  
986 TextureFields’ operation over a latent space. An  
987 interesting avenue of research would be studying the  
988 varying benefits and drawbacks of these methods for  
989 learning and synthesizing detailed textures.  
990

- 991 • Expanding the size and impact of the latent space for  
992 the shape and image encodings. One can imagine that  
993 the improved expressiveness of these high-frequency  
994 learning techniques should correspond with a richer  
995 latent space to draw from. Incorporating these tech-  
996 niques into the generation of the shape and image en-  
997 codings may thus be a worthwhile endeavor.  
998

## 1000 6. Challenges

1001 We encountered the biggest challenges of this project  
1002 when we learned from the authors of TextureFields that  
1003 training on the shapenet cars data took approximately a  
1004 week. Our original plan was to train on the Shapenet cars  
1005 for comparison, and then train on the 3DFuture data. Re-  
1006 alizing that this plan was infeasible due to time constraints,  
1007 we had to develop an entirely new training and evaluation  
1008 plan. Unfortunately this led to us not using the 3D Future  
1009 data, meaning the time we spent preprocessing the data be-  
1010 came a sunk cost.  
1011

1012 Another challenge came when we attempted to imple-  
1013 ment DeepSDF as the shape encoder within the Texture-  
1014 Fields architecture. DeepSDF, where a neural networks ap-  
1015 proximates the signed distance field, belongs to the same  
1016 class of implicit representation networks as Occupany Net.  
1017 Although DeepSDF takes a point cloud as input, it requires  
1018 an essential preprocessing before passing the input into the  
1019 model. After spending a substantial amount of time at-  
1020 tempting to get the DeepSDF preprocessing code operational,  
1021 we decided that our efforts would be better focused  
1022 elsewhere. Another bottleneck was that DeepSDF uses an  
1023 auto-decoder model as opposed to auto-encoder. This made  
1024 it less clear how best to integrate it with TextureFields.  
1025

## 7. Reflection

A limitation of our experimental process is that overfitting on a single example is not the ideal way to evaluate model performance. A better task would be to train on the entire dataset and evaluate generalization. Additionally, we also seek to incorporate an architecture like MeshCNN that can operate directly on meshes as opposed to requiring point cloud sampling with the rendering and unprojection pipeline.

Although we were not able to train on the entire dataset as initially hoped, we did manage to implement different shape encoders, perceptual losses and a positional encoding. Even within these limitations, the positional encoding implementation offered promising improvements over the base model. Although are other various experiments did not yield results better than the existing works, some experiments produced results with similar quality and most importantly, we saw pros and cons for different methods and processed the reasoning about why certain issues would occur. We accomplished the base and target goals that we modified the architecture, tried applying more state-of-the-art techniques and observed the qualitative and quantitative results.

This is a novel and challenging problem and thus despite not achieving the ultimate goal, we still learned some valuable lessons both on the engineering and research aspect and got an opportunity to ask important questions in the field of 3D deep learning.

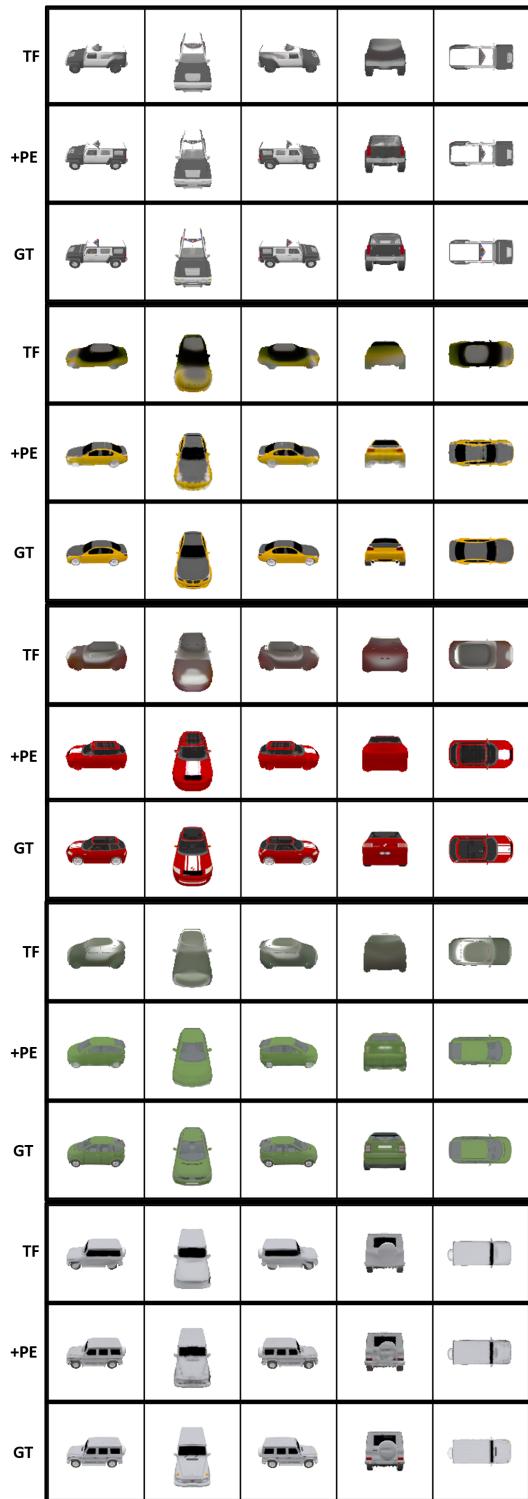
## References

- [1] Eric Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis, 2020. 10
- [2] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. Shapenet: An information-rich 3d model repository, 2015. 2, 3
- [3] Huan Fu, Rongfei Jia, Lin Gao, Mingming Gong, Bin Jiang Zhao, Steve Maybank, and Dacheng Tao. 3d-future: 3d furniture shape with texture. *arXiv preprint arXiv:2009.09633*, 2020. 2, 3
- [4] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014. 1
- [5] Thibault Groueix, Matthew Fisher, Vladimir G. Kim, Bryan C. Russell, and Mathieu Aubry. Atlasnet: A papier-mâché approach to learning 3d surface generation, 2018. 1
- [6] Xianfeng Gu, Steven J. Gortler, and Hugues Hoppe. Geometry images. 2002. 1
- [7] Rana Hanocka, Amir Hertz, Noa Fish, Raja Giryes, Shachar Fleishman, and Daniel Cohen-Or. Meshcnn. *ACM Transactions on Graphics*, 38(4):1–12, Jul 2019. 1

- 1080 [8] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Se- 1134  
1081 bastian Nowozin, and Andreas Geiger. Occupancy networks: 1135  
1082 Learning 3d reconstruction in function space, 2019. 1 1136  
1083 [9] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. 1137  
1084 Ramamoorthi, and R. Ng. Nerf: Representing scenes as neu- 1138  
1085 ral radiance fields for view synthesis, 2020. 3, 4, 5 1139  
1086 [10] Michael Oechsle, Lars Mescheder, Michael Niemeyer, Thilo 1140  
1087 Strauss, and Andreas Geiger. Texture fields: Learning tex- 1141  
1088 ture representations in function space, 2019. 1, 2, 3 1142  
1089 [11] Jeong Joon Park, Peter Florence, Julian Straub, Richard 1143  
1090 Newcombe, and Steven Lovegrove. Deepsdf: Learning con- 1144  
1091 tinuous signed distance functions for shape representation, 1145  
1092 2019. 1 1146  
1093 [12] Charles R. Qi, Li Yi, Hao Su, and Leonidas J. Guibas. Point- 1147  
1094 net++: Deep hierarchical feature learning on point sets in a 1148  
1095 metric space, 2017. 1, 2 1149  
1096 [13] J. Schmidhuber. Learning factorial codes by predictability 1150  
1097 minimization. *Neural Computation*, 4(6):863–879, 1992. 1 1151  
1098 [14] Vincent Sitzmann, Julien N. P. Martel, Alexander W. 1152  
1099 Bergman, David B. Lindell, and Gordon Wetzstein. Im- 1153  
1100 plicit neural representations with periodic activation func- 1154  
1101 tions, 2020. 3, 4, 10 1155  
1102 [15] M. Tancik, P. P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, 1156  
1103 N. Raghavan, U. Singhal, R. Ramamoorthi, J. T. Barron, and 1157  
1104 R. Ng. Fourier features let networks learn high frequency 1158  
1105 functions in low dimensional domains, 2020. 2, 3, 4, 10 1159  
1106 [16] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, 1160  
1107 and Oliver Wang. The unreasonable effectiveness of deep 1161  
1108 features as a perceptual metric. In *CVPR*, 2018. 2, 4 1162  
1109 1163  
1110 1164  
1111 1165  
1112 1166  
1113 1167  
1114 1168  
1115 1169  
1116 1170  
1117 1171  
1118 1172  
1119 1173  
1120 1174  
1121 1175  
1122 1176  
1123 1177  
1124 1178  
1125 1179  
1126 1180  
1127 1181  
1128 1182  
1129 1183  
1130 1184  
1131 1185  
1132 1186  
1133 1187

1188	<b>8. APPENDIX: Supplementary Material</b>	1242
1189		1243
1190		1244
1191		1245
1192		1246
1193		1247
1194		1248
1195		1249
1196		1250
1197		1251
1198		1252
1199		1253
1200		1254
1201		1255
1202		1256
1203		1257
1204		1258
1205		1259
1206		1260
1207		1261
1208		1262
1209		1263
1210		1264
1211		1265
1212		1266
1213		1267
1214		1268
1215		1269
1216		1270
1217		1271
1218		1272
1219		1273
1220		1274
1221		1275
1222		1276
1223		1277
1224		1278
1225		1279
1226		1280
1227		1281
1228		1282
1229		1283
1230		1284
1231		1285
1232		1286
1233		1287
1234		1288
1235		1289
1236		1290
1237		1291
1238		1292
1239		1293
1240		1294
1241		1295

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344



1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

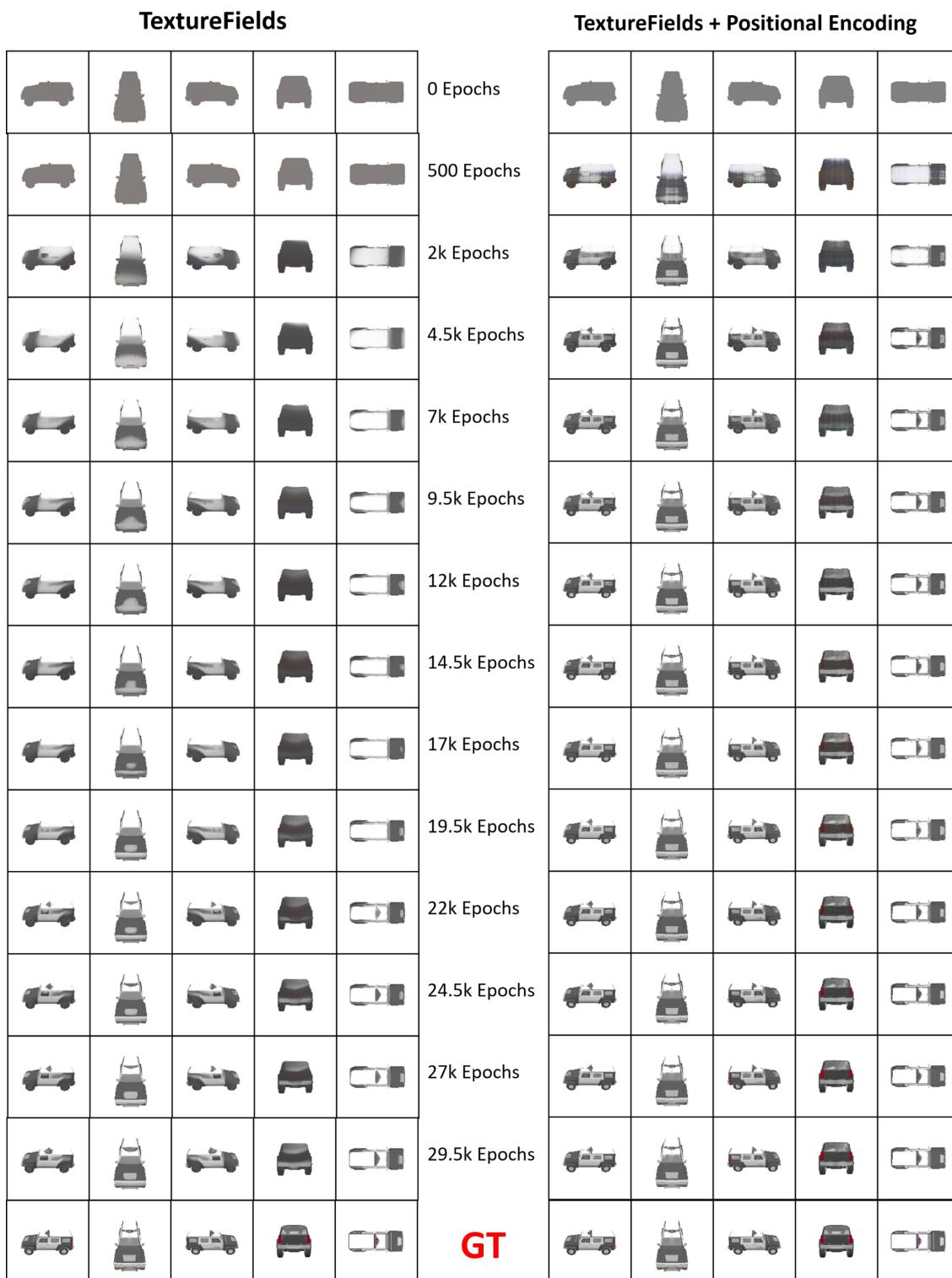
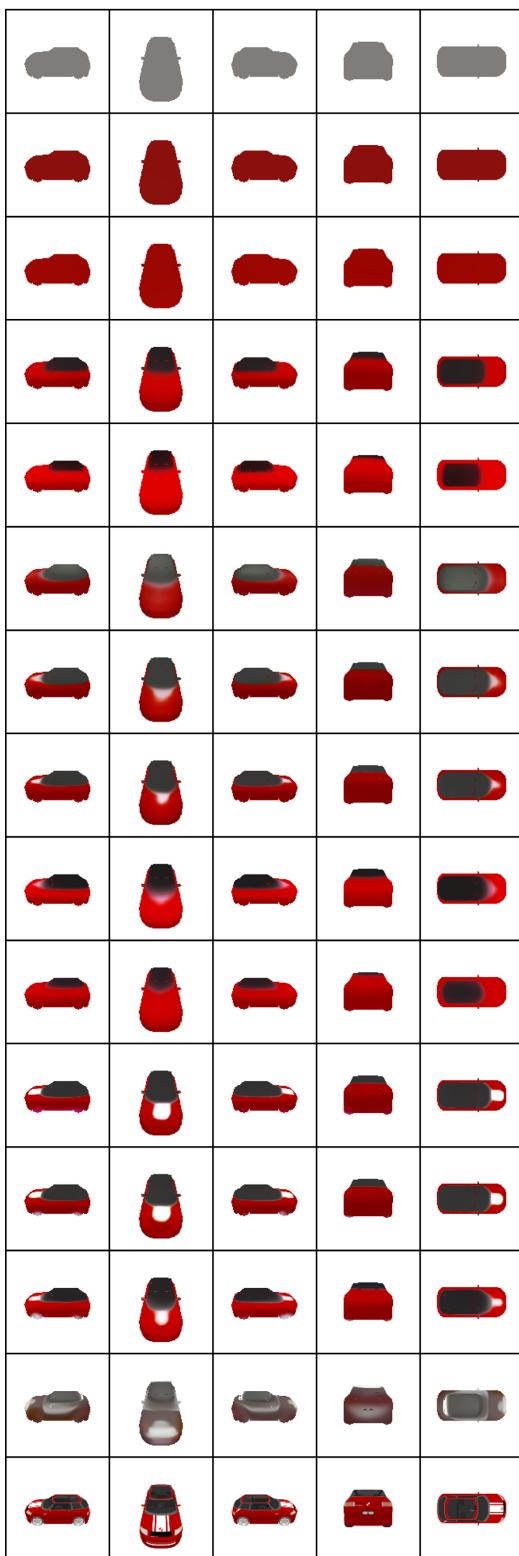
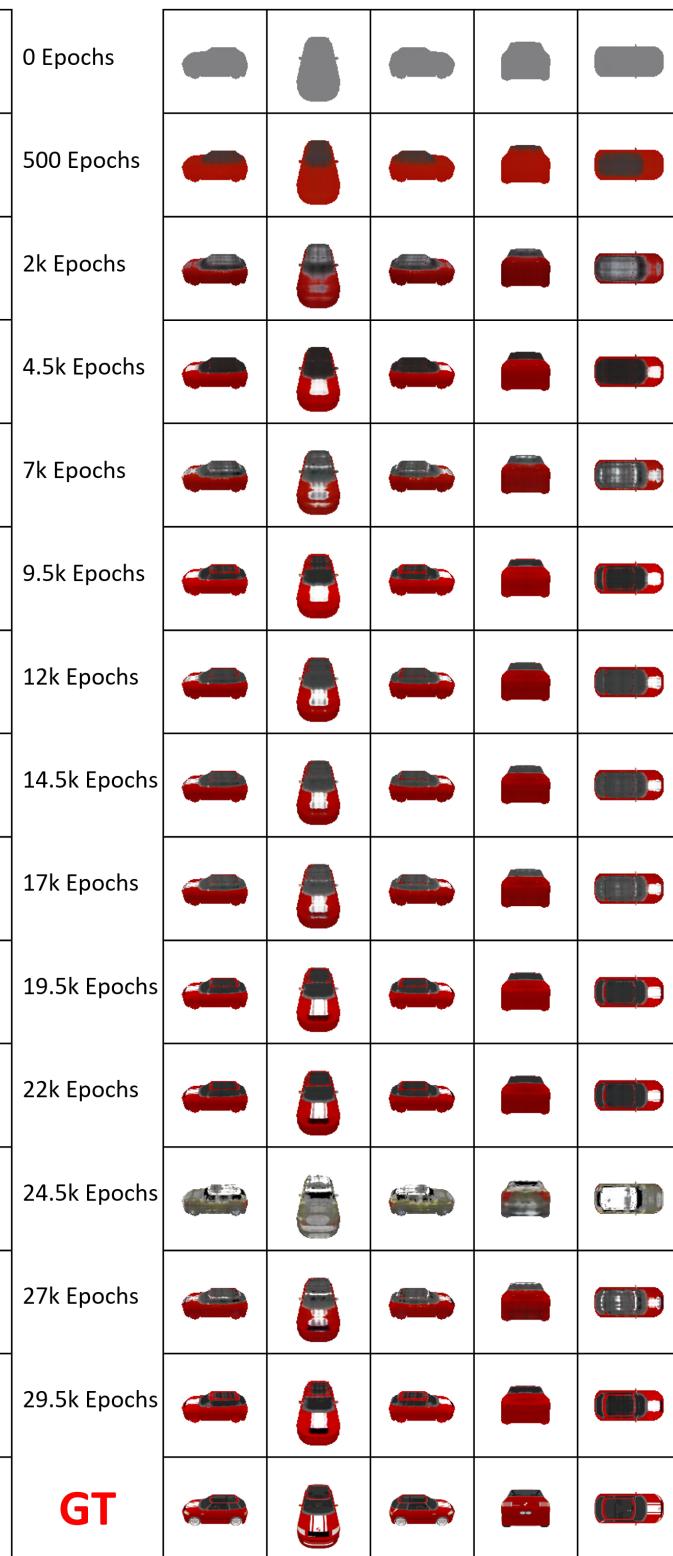


Figure 23. Comparison of textures generated by *ORIG* and *PE* for a car in the training set at 2.5k epoch intervals across the course 29.5k epochs of training. The bottom row are the ground truth images. Note the early approximation of details for the *PE* model.

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

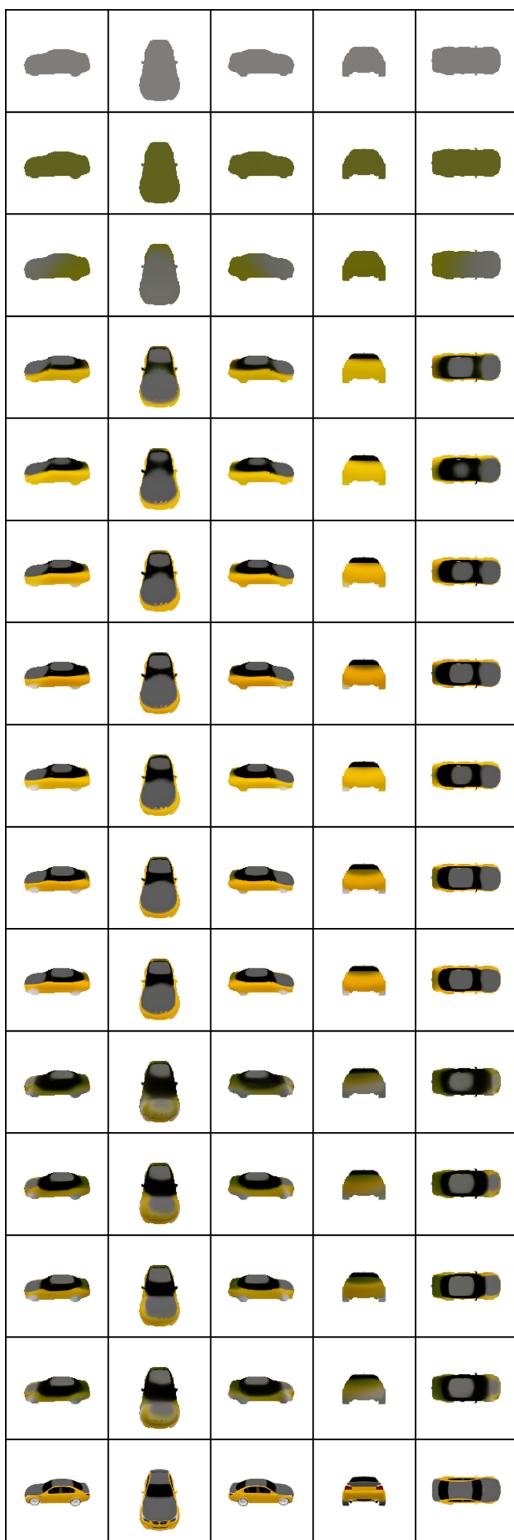
**TextureFields****TextureFields + Positional Encoding**

GT

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

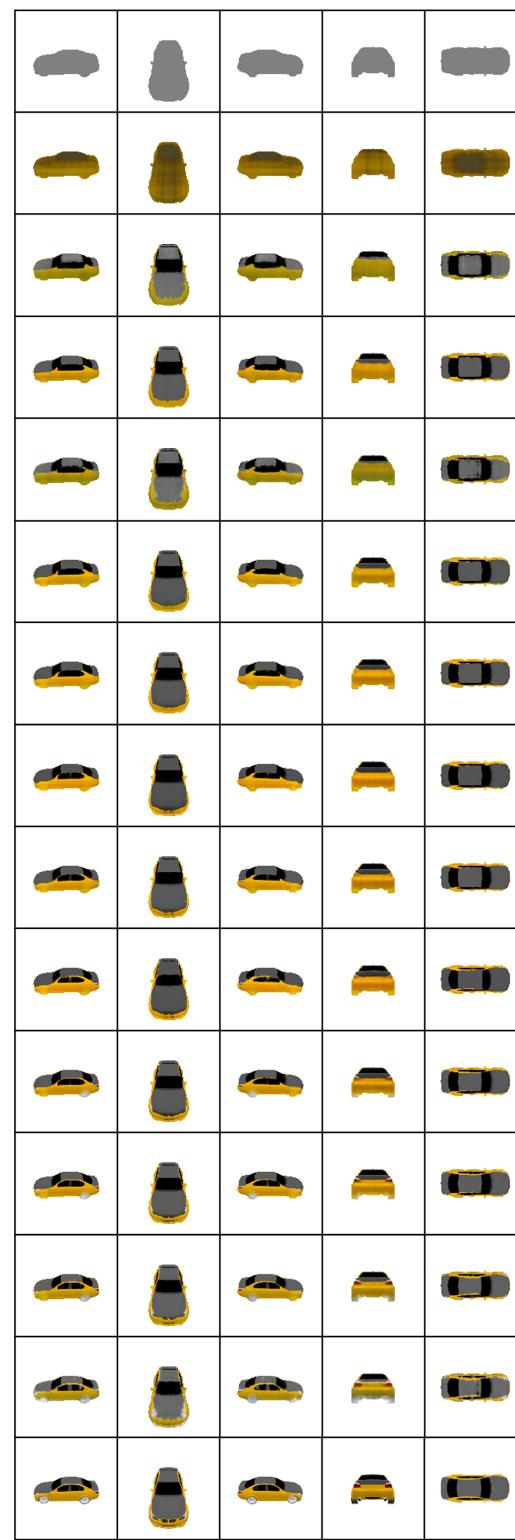
Figure 24. Comparison of textures generated by *ORIG* and *PE* for a car in the training set at 2.5k epoch intervals across the course 29.5k epochs of training. The bottom row are the ground truth images. Note the early approximation of details for the *PE* model. There are abrupt distortions in the generated textures for both *ORIG* and *PE* at epochs 29.5k and 24.5k, respectively

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

**TextureFields****TextureFields + Positional Encoding**

0 Epochs  
500 Epochs  
2k Epochs  
4.5k Epochs  
7k Epochs  
9.5k Epochs  
12k Epochs  
14.5k Epochs  
17k Epochs  
19.5k Epochs  
22k Epochs  
24.5k Epochs  
27k Epochs  
29.5k Epochs

GT



1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727

Figure 25. Comparison of textures generated by *ORIG* and *PE* for a car in the training set at 2.5k epoch intervals across the course 29.5k epochs of training. The bottom row are the ground truth images. Note the early approximation of details for the *PE* model.