

# Generative Modelling of BRDF Textures from Flash Images

PHILIPP HENZLER, University College London, United Kingdom

VALENTIN DESCHAINTE, Adobe Research and Imperial College London, United Kingdom

NILOY J. MITRA, University College London and Adobe Research, United Kingdom

TOBIAS RITSCHEL, University College London, United Kingdom

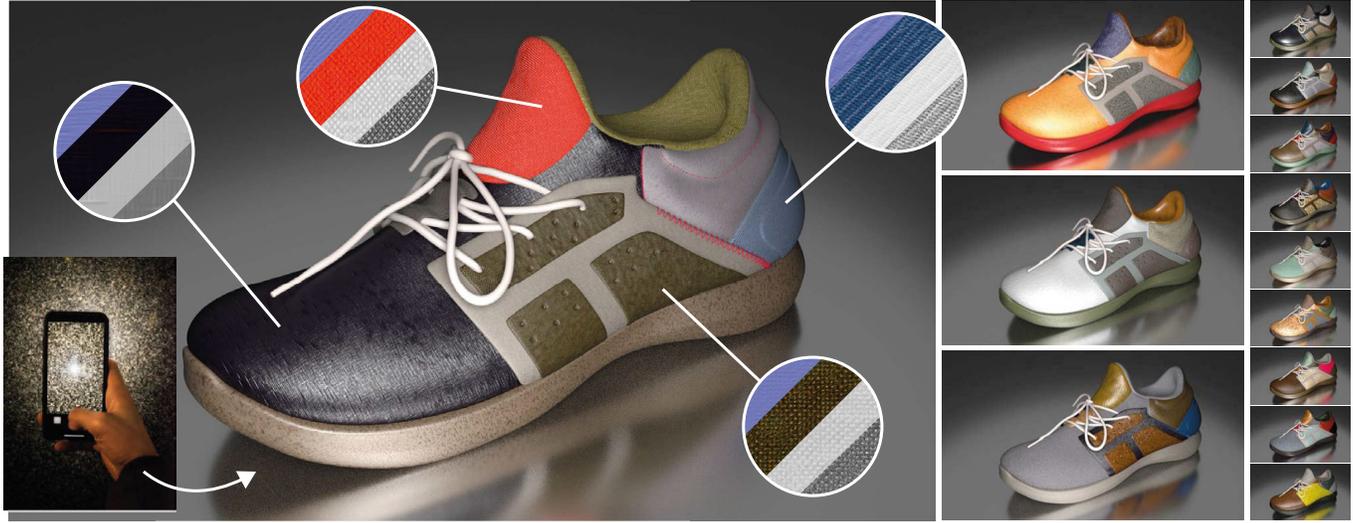


Fig. 1. Results from our generative model of BRDF maps, assigned to a 3D object of a shoe. Circular insets show the diffuse, normal, roughness and specular maps. Our model outputs a space of BRDF materials that can be sampled from, projected to, and interpolated across. The BRDF generative model is trained exclusively from a set of (306) RGB flash images (example shown in the lower left inset) without any BRDF supervision, and shortly fine-tuned in the case of material capture to best match the input picture.

We learn a latent space for easy capture, consistent interpolation, and efficient reproduction of visual material appearance. When users provide a photo of a stationary natural material captured under flashlight illumination, first it is converted into a latent material code. Then, in the second step, conditioned on the material code, our method produces an infinite and diverse spatial field of BRDF model parameters (diffuse albedo, normals, roughness, specular albedo) that subsequently allows rendering in complex scenes and illuminations, matching the appearance of the input photograph. Technically, we jointly embed all flash images into a latent space using a convolutional encoder, and –conditioned on these latent codes– convert random spatial fields into fields of BRDF parameters using a convolutional neural network (CNN). We condition these BRDF parameters to match the visual characteristics (statistics and spectra of visual features) of the input under matching light. A user study compares our approach favorably to previous work, even those with access to BRDF supervision. Project webpage: <https://henzler.github.io/publication/neuralmaterial/>.

CCS Concepts: • **Computing methodologies** → **Reflectance modeling**; **Image processing**.

Additional Key Words and Phrases: material capture, appearance capture, SVBRDF, deep learning, generative model, unsupervised learning

Authors' addresses: Philipp Henzler, Department of Computer Science, University College London, United Kingdom, [p.henzler@cs.ucl.ac.uk](mailto:p.henzler@cs.ucl.ac.uk); Valentin Deschaintre, Adobe Research and Imperial College London, United Kingdom; Niloy J. Mitra, University College London and Adobe Research, United Kingdom; Tobias Ritschel, Department of Computer Science, University College London, United Kingdom.

## 1 INTRODUCTION

Rendering realistic images for feature films or computer games requires adequate simulation of light transport. Besides geometry and illumination, an important factor is material appearance.

Material appearance has three aspects of variation: First, when view or light direction change, reflected light changes. The physics of this process is well-understood and can be simulated provided the input parameters are available. Second, behavior changes across materials. For example, leather reacts differently to light or view changes than paper would, yet, different forms of leather clearly share visual properties, i.e., form a (material) space. Third, appearance details depend on spatial position. Different locations in the same leather exemplar behave differently but share the same visual statistics [Portilla and Simoncelli 2000], i.e., they form a *texture*.

Classic computer graphics captures appearance by *reflection models*, which predict for a given i) light-view configuration, ii) material, and iii) spatial position, how much light is reflected. Typically, the first variation (light and view direction) is covered by *BRDF models*, analytic expressions, such as Phong [1975] which map the light and view direction vector to scalar reflectance. The second variation (material) is covered by choosing *BRDF model parameters*, such as specular or roughness. In practice, it can be difficult, given a desired appearance, to choose those parameters e.g., how to make a

leather look more like the one on a nice jacket. One can measure BRDF model parameters, but it traditionally requires complex capture hardware for accurate results. The third variation (spatial) is addressed by storing multiple BRDF model parameters in images of finite size –often referred to as Spatially-varying Bi-directional Reflectance Distribution Function (svBRDF) maps– or writing functional expressions to reproduce their behaviour. It is even more challenging to choose these parameters to produce something coherent like leather, in particular over a large to very large spatial extent. Additionally, storing all these values requires substantial memory and programming functional expressions to mimic their statistics requires expert skills and time. Capturing the spatial variation of BRDF model parameters over space using sensors requires even more complex hardware [Schwartz et al. 2013].

Addressing those issues, we provide a reflectance model to jointly generalize across all of these three axes. Instead of using analytic parameters, we parametrize appearance by latent codes from a learned space and our decoder weights, allowing for acquisition, interpolation and generation. Without involved capture equipment, these codes are produced by presenting the system a simple 2D flash image, which is then embedded into the latent space. Avoiding to store any finite image texture, we learn a second mapping to produce svBRDF maps from the infinite random field (noise) on-the-fly, conditioned on the latent material code and decoder weights. Instead of using any advanced capture device for learning, flash images will be the only supervision we use. This unsupervised approach allows us to consider our decoder weights as part of the latent representation, which we fine-tune at test time in a few minutes.

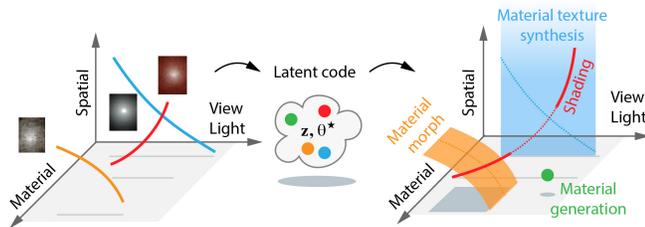


Fig. 2. **BRDF space.** From a flash image, which contains sparse observations across material, space and view-light (**left**) we map to a latent code  $z / \theta^*$  (**middle**) so that changes in these code can be decoded to enable (**right**) material synthesis (holding material fixed and moving spatially), material morphing (holding space and view/light fixed and changing material), or classical shading and material generation (points in the latent space).

A use case of our approach is shown in Fig. 1. First, a user provides a “flash image”, a photo of a flat material sample under flash illumination. This sample is embedded as a code into a latent space using a CNN and used to fine-tune our decoders weights. This code and weights can then be manipulated, e.g., interpolated with a different material. Conditioned on this code, our fine-tuned decoder can generate an infinite field of BRDF to be directly used in rendering.

For training, we solely rely on real flash images. The key insight, inspired by Aittala et al. [2016], is that these flash images reveal the same material at different image locations –they are stationary– but under different view and light angles. Using this constraint, Aittala

et al. [2016] were able to decompose a small patch of a single input image to capture the parameters of a material model that could then be rendered under novel view or light directions. However, this covers only part of the generalization we are targeting: it generalizes across view and light, but not across location or material. Further, they perform an optimization for every exemplar, requiring time in the order of an hour, while ours takes minutes only.

In summary, our main contributions are

- a generative model of a BRDF material texture space;
- generation of maps that are diverse over the infinite plane;
- a flash image dataset of materials enabling our training with no BRDF parameter supervision or synthetic data

Our implementation will be publicly available upon acceptance.

## 2 PREVIOUS WORK

Our work has background in texture analysis, appearance modelling and design spaces as summarized in Tbl. 1 and discussed next.

Table 1. Comparison of features between different previous methods. We distinguish methods producing RGB from those generating **BRDF** or **svBRDF**, whether those can be **Non-Stationary** and **Infininitely** sampled. We also distinguish if their results for one input can be **Diverse**, if they form a **Space** which can be queried and how **Fast** direct sampling is.

Method	Supervision	BRDF	svBRDF	Non-Stat.	Infinite	Diverse	Space	Fast Gen.
Classic texture synth	RGB	×	×	✓	✓	✓	×	✓
Matusik et al. [2005]	RGB	×	×	✓	✓	✓	✓	×
Matusik [2003]	BRDF	✓	×	×	×	✓	✓	✓
Georgoulis et al. [2017]	BRDF	✓	×	×	×	×	✓	✓
Deschaintre et al. [2018]	svBRDF	✓	✓	×	×	×	×	✓
Zhao et al. [2020]	Flash image	✓	✓	✓	×	×	×	×
Aittala et al. [2016]	Flash image	✓	✓	✓	×	×	×	×
Gao et al. [2019]	svBRDF	✓	✓	✓	×	×	×	×
Guo et al. [2020b]	svBRDF	✓	✓	✓	×	×	✓	×
Ours	Flash image	✓	✓	×	✓	✓	✓	✓

### 2.1 Textures in Graphics

A classic definition of texture is defined by Julesz [1965]: *a texture is an image full of features that in some representation have the same statistics.* Portilla and Simoncelli [2000] have provided a practical method to compute representations in which to do statistics on, using linear filters on multiple scales.

Perlin [1985] was first to capture the fractal [Mandelbrot 1983] stochastic variation of appearance in a model applicable to Computer Graphics. His approach is simple –a linear combination of noise at different scales– yet extremely powerful, and has led to extensive use in computer games and production rendering. Wavelet noise [Cook and DeRose 2005] moved this idea further by band-limiting the noise that is combined. Such methods can be used for materials, e.g., gloss maps, bump maps, etc. It however does not provide a solution to acquire a texture from an exemplar, which is left to manual adjustment.

To generate textures from exemplars, non-parametric sampling [Efros and Leung 1999], vector quantization [Wei and Levoy 2000], optimization [Kwatra et al. 2005] or nearest-neighbour field synthesis (PatchMatch [Barnes et al. 2009]) have been proposed. They however have issues in computational scalability and lack intuitive control, limiting their adoption in production rendering or games.

The word “texture” can be ambiguous to mean stochastic variation, as well as images attached on surfaces to localize color features. Here, we focus on stochastic variation in the sense of Julesz [1965] or Portilla and Simoncelli [2000].

Our approach is inspired by deep learning-based texture synthesis [Bergmann et al. 2017; Gatys et al. 2015; Johnson et al. 2016; Karras et al. 2019; Sendik and Cohen-Or 2017; Shaham et al. 2019; Simonyan and Zisserman 2014; Ulyanov et al. 2016, 2017; Zhou et al. 2018] which ideas we extend and apply to BRDFs. We detail their background in Sec. 3.2

## 2.2 Material Modeling

Representing appearance in simulation-based graphics has been an active research field for decades. The survey by Guarnera et al. [2016] presents detailed discussion of the many different material model and BRDF acquisition approach. In our method, we use a state-of-the-art micro-facet BRDF model [Cook and Torrance 1982], and focus on deep-learning based material modelling and acquisition.

Many methods have been proposed to acquire materials using data-driven approaches. Matusik [2003] proposed a data driven BRDF linear model. More recently, Rematas et al. [2016] extract reflectance maps from 2D images using a CNN trained in a supervised manner. Materials and illuminations acquisition were further explored by Georgoulis et al. [2017]. Deschaintre et al. [2018] proposed a rendering loss to capture svBRDFs from flash images. Nam et al. [2018] jointly reconstructed svBRDF, normals, and 3D geometry in an iterative inverse-rendering setup towards a practical acquisition setup, while different methods relied on deep learning to estimate object shape and svBRDF from one or multiple images [Boss et al. 2020; Deschaintre et al. 2021; Li et al. 2018b]. Li et al. [2019] propose a weakly supervised learning-based method for generating novel category-specific 3D shapes and demonstrate that it can help in learning material-class specific svBRDFs from images distributions. Ye et al. [2018] used a mixture of images and procedural material maps to train a network for modeling svBRDFs. Hu et al. [2019] developed a reduced svBRDF model, using only diffuse and normal channels, towards solving inverse procedural textures matching from reference, while Guo et al. [2020a] used Bayesian inference for material synthesis. Recently, Shi et al. [2020] developed a differentiable material graph nodes library to optimize material parameters to match an input material, given material graphs.

U-net [Ronneberger et al. 2015] inspired many approaches for image to image translation to translate RGB pixels to material attributes [Deschaintre et al. 2018; Li et al. 2017, 2018a,b]. Most work now includes a differentiable shading step [Deschaintre et al. 2018, 2019; Guo et al. 2020b; Li et al. 2018b; Liu et al. 2017] such as we do here. Gao et al. [2019] and Guo et al. [2020b] propose to use a post-optimization in an encoded latent space, improving an initial

material estimation, and comparing renderings of their results directly to their input pictures. Deschaintre et al. [2020] propose to fine tune their material acquisition network on svBRDF parameter examples to transfer them to a larger scale. Zhou and Kalantari [2021] propose a partially unsupervised training approach, allowing to use additional real data without ground truth maps. With our approach, we completely remove the need for ground truth maps as it solely relies on real flash photographs. Guo et al. [2021] address the issue of strong highlights baked into svBRDF maps through highlight-aware convolutions and an attention-based feature selection module. Our design, focused on stationarity of textures, inherently prevents any flash residual to be left in the results.

All these approaches focus on capturing a single instance of a svBRDF map, but with little or no editing options across materials (space) or generalization across the spatial domain (diversity). For rapid materials generation, Zsolnai-Fehér et al. [2018] propose to use Gaussian process regression.

However, most of these methods require synthetic svBRDF supervision for training, while we focus on directly learning from flash images without access to channel-level supervision. In particular this removes the risk of domain gap between synthetic and real materials and enables our fine-tuning approach.

We take inspiration from Aittala et al. [2016] who extended the approach of Gatys et al. [2015] to generate svBRDF parameter maps from a single picture of a stationary material exemplar and propose an approach for improved diversity, generation and quality.

## 2.3 Spaces-of

Spaces of color [Nguyen et al. 2015], materials [Gao et al. 2019; Guo et al. 2020b; Matusik 2003], textures [Matusik et al. 2005], faces [Blanz et al. 1999], human bodies [Allen et al. 2003], and more have been useful in graphics for content creation and edition. Matusik et al. [2005] has devised a space of textures. Here, users can interpolate combinations of visually similar textures. They warp all pairs of exemplars to each other and constructs graph edges for interpolation when there is evidence that the warping is admissible. To blend between them, histogram adjustments are made. Consequently, interpolation between exemplars does not take a straight path in pixel space from one to the other, but traverses only valid regions. Photoshape [Park et al. 2019] learns the relation of given material textures over a database of 3D objects. Serrano et al. [2016] allow users to semantically control captured BRDF data. They represent BRDFs using the derived principal component basis [Matusik 2003] and map the first five PCA components to semantic attributes through learned radial basis functions. Similar to our method, Guo et al. [2020b] and Gao et al. [2019] produce spaces of materials that can be interpolated. We take inspiration from this body of work and build a space allowing svBRDFs generation and interpolation.

# 3 BACKGROUND

## 3.1 Flash Images

Aittala et al. [2016] leveraged the fact that a single flash image of a stationary material reveals multiple realizations of the same reflectance statistics under different light and view angles. We will now recall a simplified definition of their approach.

A flash image is an RGB image of a material, taken in conditions where a mobile phone’s flashlight is the dominant light source. We write  $L(\mathbf{x})$  to denote the RGB radiance value at every image location  $\mathbf{x}$ . The illumination is expected to be an isotropic point light collocated with the camera. Further, the geometry is assumed to be flat and captured in a fronto-parallel setting, so that the direction from light to every image location in 3D is known. Self-occlusion and parallax are assumed to be negligible.

Reflectance is parameterized by a *material*, represented as a function  $f(\mathbf{x})$  mapping image location  $\mathbf{x}$  to shading model parameters, including the shading normal. Under these conditions, the reflected radiance is  $L = \mathbf{R}f$ , where  $\mathbf{R}$  is the *differentiable* rendering operator, mapping shading model parameters to radiance.

A material  $f$  explains a flash image  $L$  if it is *visually similar* to  $L$  when rendered. Unfortunately, without further constraints, there are many materials to explain the flash image. This ambiguity can be resolved when assuming that the material  $f$  is *stationary*. We say a material is stationary, if local statistics of the shading model parameters  $f$  do not change across the image.

Putting both –visual similarity and stationarity– together, the best material from a family  $f_\theta$  of material mapping functions parameterized by a vector  $\theta$ , can be found by minimizing a loss:

$$\mathcal{L}'(\theta) := \mathcal{T}(L, \mathbf{R}f_\theta) + \lambda \mathcal{S}(f_\theta), \quad (1)$$

where  $\mathcal{T}(L, \mathbf{R}f_\theta)$  is a metric of visual similarity between a flash image  $L$  and a differentiable rendering  $\mathbf{R}f_\theta$ , and  $\mathcal{S}(f)$  is a measure of stationarity of a material map  $f$ .

Comparison,  $\mathcal{T}$ , of two textures is not trivial. Pixel-by-pixel comparison is typically not suitable to evaluate visual statistical similarity. Instead, images are mapped to a feature space in which images that are perceived as similar textures, map to similar points [Portilla and Simoncelli 2000]. Different mappings are possible here. Classic texture synthesis [Heeger and Bergen 1995] uses moments of linear multi-scale filters responses. Gatys et al. [2015] proposed to use Gram matrices of non-linear multi-scale filters responses such as those of the VGG [Simonyan and Zisserman 2014] detection network. Such a characterization of textures was also used by Aittala et al. [2016] and, without loss of generality, will be used and extended in this work as well.

While  $f$  is stationary,  $L$  is not –due to the lighting– and has features at different random positions  $\mathbf{x}$  which are compared as

$$\mathcal{T}'(L_1, L_2) := \mathbb{E}_{\mathbf{x} \sim (0,1)^2, s \sim (0,1)} [|\mathcal{P}(L_1, \mathbf{x}, s) - \mathcal{P}(L_2, \mathbf{x}, s)|_1], \quad (2)$$

where  $\mathcal{P}'(L, \mathbf{x})$  crops a patch of randomly chosen scale  $s$  at the location  $\mathbf{x}$  and resamples it to the input resolution of VGG [Simonyan and Zisserman 2014], computes the filter responses and their Gram matrices:

$$\mathcal{P}(L, \mathbf{x}, s) := \text{gram}(\text{vgg}(\text{resample}(\text{crop}(L, \mathbf{x}, s)))). \quad (3)$$

Minimizing  $\theta$  with respect to Eq. 1 for a given  $L$  results in a material.  $f_\theta$  can represent different approaches. Aittala et al. [2016] directly use the pixel basis and optimize discrete material maps for  $\theta$  using a single input flash image  $L$ . With their approach, optimizing for both visual similarity and stationarity is challenging. In particular, the reflectance stationarity term  $\mathcal{S}$ , requires a “spectral

preconditioning” step as explained in their paper. Instead, we propose an approach in the form of a neural model  $f$  that is (i) defined on the infinite domain and (ii) stationary by construction. Thus, our loss does not need to include a stationarity term.

Next, we describe how to generate RGB textures using deep learning (Sec. 3.2), before combining the two components (flash images and NN texture (spaces)) into our approach (Sec. 4).

### 3.2 Deep Texture Synthesis

Julesz [1965] define textures by their feature statistics across space. The choice of which features to use remains an important open problem. With the advent of deep learning, Gatys et al. [2015] suggested to use Gram matrices of activations of filters learned in deep convolutional neural networks (e.g., VGG [Simonyan and Zisserman 2014]), for neural style transfer. Aittala et al. [2016] rely on the same statistics to recover material parameters of stationary materials. By optimizing directly over pixel values, their method can produce images with the desired texture properties. These methods however require a different long optimization to be ran for each material.

Another group of recent methods [Johnson et al. 2016; Ulyanov et al. 2016] introduce neural networks capable of producing RGB textures directly, in milliseconds. While these approaches use a network to generate the textures, they are still limited to the input texture exemplar, and do not show further variations in their results. Ulyanov et al. [2017] introduced an explicit diversity term enforcing results in a batch to be different. This diversity is however limited and restrict the results quality. Indeed, they add a diversity term to the loss, but the architecture is not modified to enable it. Alternatively, adversarial training has been used to capture the essence of textures [Bergmann et al. 2017; Shaham et al. 2019], including the non-stationary case [Zhou et al. 2018] or even within a single image [Shaham et al. 2019]. In particular, StyleGAN [Karras et al. 2019] generates images with details by transforming noise using adversarial training. As opposed to these approach we do not rely on challenging adversarial trainings, by directly learning a Neural Network to produce VGG statistics.

Instead of incentivizing stationarity in the loss, Henzler et al. [2020] suggest a learnable texture representation that is built on mapping an infinite noise field to a field that has the statistics of the exemplar texture. Their method is a point operation, implemented by an MLP that is fed exclusively with noise sampled at different scales as done by Perlin [1985]. By explicitly preventing the network to access any absolute position, this approach is stationary by-design. Inspired by this approach our architecture enforces a convolutional stationary by-design constraint.

## 4 NOISE TO BRDF TEXTURE SPACES

An overview of our approach is shown in Fig. 3. We train a neural network which acts as a decoder  $f_\theta(\mathbf{x}|\mathbf{z})$  that generalizes across spatial positions  $\mathbf{x}$  as well as across materials, expressed as latent material codes  $\mathbf{z}$ . The material codes  $\mathbf{z}$  are produced by an encoder  $g$  with  $\mathbf{z} = g(L)$ . Both encoder and decoder are trained jointly over a set of flash images using the loss:

$$\mathcal{L}(\theta) := \mathbb{E}_L [\mathcal{T}(L, \mathbf{R}f_\theta(\cdot|g_\theta(L)))]. \quad (4)$$

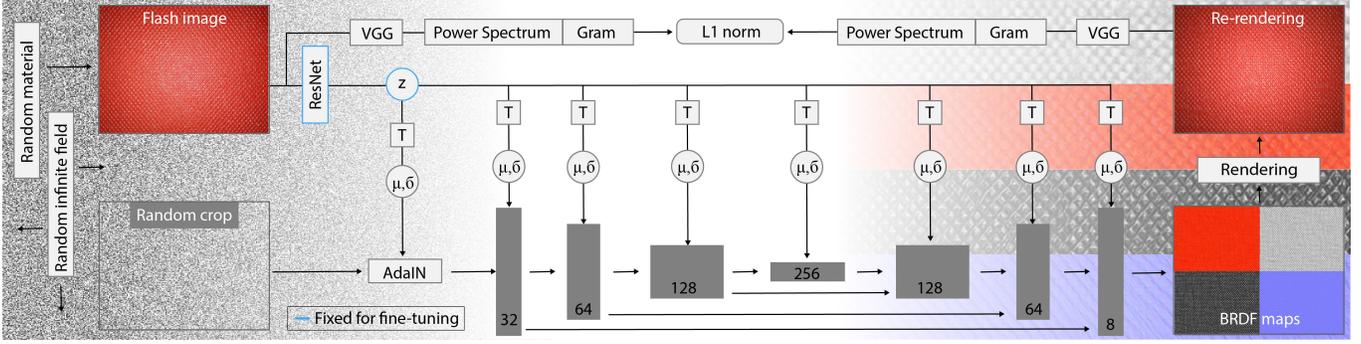


Fig. 3. **Our architecture.** Starting from an exemplar (top-left) our trained encoder encodes the image to a compact latent space variable  $z$ . Additionally, a random infinite field is cropped with the same spatial dimensions as the flash input image. The noise crop is then reshaped based on a convolutional U-Net architecture. Each convolution in the network is followed by an Adaptive Instance Normalization (AdaIN) layer [Huang and Belongie 2017] reshaping the statistics (mean  $\mu$  and standard deviation  $\sigma$ ) of features. A learned affine transformation  $T$ -s per layer maps  $z$  to the desired  $\mu$ -s and  $\sigma$ -s. The output of the network are the diffuse, specular, roughness, normal parameters of an svBRDF that, when rendered using a camera colocated flash light, look the same as the input. Our unsupervised setting allows us to fine-tune our trained network on materials to acquire.

This equation is an adapted version of Eq. 1 to fit our objectives. In particular we propose a neural network-based  $f_\theta$ , leveraging the expectation  $\mathbb{E}_L$  over all flash images in our training set and removing the stationarity term as it is enforced by construction in our network architecture. We describe the flash image encoder  $g$  (Sec. 4.1), the material texture decoder  $f$  (Sec. 4.2), the texture comparison model  $\mathcal{T}$  (Sec. 4.3) and our fine tuning approach (Sec. 4.5), next.

#### 4.1 Encoder

The encoder  $g$  maps a flash image  $L$  to a latent code  $z$ . The flash images used by our method are similar to those of recent svBRDF acquisition papers [Aittala et al. 2016; Deschaintre et al. 2018]: we use a phone with a flash colocated with the camera and capture surface in a fronto-parallel way. Our encoder is implemented using ResNet-50 [He et al. 2016]. The ResNet starts at a resolution of  $512 \times 384$  and maps to a compact latent code. Empirically, we find a  $n_z = 64$ -dimensional latent space to work best for our data and present all results using this number.

#### 4.2 Decoder

The decoder  $f$  maps location  $\mathbf{x}$ , conditioned on a material code  $z$  to a set of material parameter maps. The key idea is to provide the architecture with access to noise, as previously done for style transfer [Huang and Belongie 2017], generative modelling [Karras et al. 2019] or 3D texturing [Henzler et al. 2020]. In particular, we sample rectangular patches with edge length of  $n \times m$  pixels from an infinite random field and convert them to material maps using a U-net architecture [Ronneberger et al. 2015]. The U-net starts at the desired output resolution  $n \times m$  and reduces resolution four times using max-pooling before upsampling back to  $n \times m$  through a series of bi-linear upsampling and convolutions. Let  $F$  be the array of input features. For  $i = 0$ , the first level, in full resolution, these features are sampled from the random field at  $\mathbf{x}$ . Then, output features are

$$F' := \text{adaIN}(\text{conv}_\theta(F), T_\theta z), \quad (5)$$

where adaIN is Adaptive Instance Normalization (AdaIN) [Huang and Belongie 2017], conv a convolution (including up- or down-sampling and ReLU non-linearity),  $z$  is a latent material code and  $T$  is an affine transformation. Components with learned parameters are denoted with subscript  $\theta$ .

We use AdaIN as defined by Huang and Belongie [2017] as

$$\text{adaIN}(\xi, \{\mu, \sigma^2\}) = \frac{\sigma}{\sigma_F} (\xi - \mu_F) + \mu \quad (6)$$

and remaps the input features with mean  $\mu_F$  and variance  $\sigma_F^2$  to a distribution with mean  $\mu$  and variance  $\sigma^2$ .

The affine mapping  $T$  is implemented as  $(n_z + 1) \times (2 \times c_i)$  matrices multiplied with the latent code  $z$ . Here  $2 \times c_i$  represent a different mean and variance for each channel dimension  $c_i$  of a layer. It provides the link between the material code and the noise statistics. Each material code  $z$  is mapped to a mean and variance to control how the statistics of features are shaped at every channel on every layer of the decoder.

Our control of noise statistics from latent codes is similar to StyleGAN [Karras et al. 2019], with the key difference that we do not sample noise at different scales, but learn how to produce noise with different, complex, characteristics at different scales by repeatedly filtering it from high resolutions.

#### 4.3 Images Comparison

As mentioned in Sec. 3.1 (Eq. 2 and Eq. 3), we want to evaluate visual similarity and stationarity. To this end, we propose to compare images based on a loss that accounts both for the statistics of activations [Gatys et al. 2015] and their spectrum [Liu et al. 2016] on multiple scales across the infinite spatial field,

$$\mathcal{T}(L_1, L_2) := \mathbb{E}_{\mathbf{x} \sim \mathbb{R}^2, s \sim (s_{\min}, s_{\max})} [|\mathcal{P}(L_1, \mathbf{x}, s) - \mathcal{P}(L_2, \mathbf{x}, s)|_1]. \quad (7)$$

$$\mathcal{P}(L, \mathbf{x}, s) := \text{gram}(V(L, \mathbf{x}, s)) + \lambda \cdot \text{powerSpectrum}(V(L, \mathbf{x}, s)) \quad (8)$$

$$V(L, \mathbf{x}, s) := \text{vgg}(\text{resample}(\text{crop}(L, \mathbf{x}, s))) \quad (9)$$

*Spectrum.* VGG Gram matrices capture the frequency of a feature appearance, unless it forms a regular pattern [2016]. Liu et al. [2016] proposed to include the L1 norm of the power spectra of RGB images into the texture metric for texture synthesis. We combine both ideas and use VGG, but do not limit ourselves to its Gram matrix statistics, and also leverage its spectrum. We set  $\lambda = 1e - 3$ .

*Scale.* As VGG works at a specific scale of features it was trained for, it behaves differently at different scales. As the material should be visually plausible regardless of its scale we include multiple scales  $s$ , ranging from  $s_{\min} = 0.1$  to  $s_{\max} = 8$  in the loss computation.

*Infinity.* Expectation over the infinite plane is implemented by simply training with different random seeds for the noise field. This results in the generation of statistically similar, but locally different variations of materials. As, given a seed, every generated patch is a coherent material, combinations of multiple patches remains coherent as well. This allows to query an endless, seamless and diverse stream of patches without repetition. It also prevents overfitting and is crucial to guarantee stationarity by-design.

#### 4.4 Training

To enforce a generalizable material prior, we first train the system as a Variational Auto-encoder (VAE) [Kingma and Welling 2013]. Instead of mapping to a single 64-D latent material code, the encoder  $g$  maps to a 64-D mean and variance vector, from which we sample in training. At test time we use the mean for each 64-D. We have omitted the additional VAE terms enforcing  $z$  to be normally distributed from Eq. 4 and Fig. 3 for clarity. We trained our model for 4 days and a batch size of 4 on a NVIDIA Tesla V100 using ADAM optimizer with a learning rate of  $1e-4$  and weight decay  $1e-5$ .

#### 4.5 Fine-tuning

Using the trained encoder-decoder pair we can instantaneously compress a 2D RGB flash image to a latent code and decompress into an infinite svBRDF field. The quality of the decoding can further be improved by adapting the decoder weights to a specific exemplar  $L^*$  with a short one-shot training. To this end, all weights  $\theta$  are held fixed, except for the decoder weights  $\theta^* \subset \theta$ , which are further trained to reproduce a single flash image  $L^*$  at material code  $z^* = g(L^*)$ . This is made possible by our completely unsupervised approach, allowing to fine-tune for any flash image, without requiring ground truth maps. Note that unlike [Guo et al. 2020b] we use a style loss rather than a pixel-wise loss for fine-tuning, preserving the *diversity* properties of our results. In practice we fine-tune for 1000 steps with an increased learning rate by a factor of 10, for about 5 minutes.

Fine-tuning of two materials will result in two different decoders  $f_1$  and  $f_2$  as well as two latent codes  $z_1$  and  $z_2$  produced by the same encoder. We show that despite being a more complex space, interpolating both the latent code and decoder parameters, as in  $\text{lerp}(f_1, f_2)(\text{lerp}(z_1, z_2))$  works well in practice. Unless otherwise specified, we show fine-tuned results in the remainder of this paper and ablate several variants in Sec. 5.4.

#### 4.6 Material model

We use the Cook-Torrance [1982] micro-facet BRDF Model, with Smith’s geometric term [Heitz 2014], Schlick’s [1994] Fresnel and GGX [Walter et al. 2007]. Hence, parameters are diffuse RGB albedo, monochromatic specular albedo, roughness and height, i.e., six dimensions. Instead of learning a normal map, a height field is generated from which normals are computed using finite differences. During the our differentiable rendering step, we assume a FOV of  $45^\circ$  to simulate smartphone cameras.

#### 4.7 Alignment

Many flash images entail a slight rotation as it can be difficult to take a completely fronto-parallel image. This was handled by Aittala et al. [2016] by locating the brightest pixel and cropping, but we found our, more abstract, training to struggle with such a solution.

Instead, we add a horizontal and a vertical rotation angle to the parameter vector generated from the latent code (not shown in Fig. 3 for clarity). During training, these are used to rotate the plane, including the normals. During testing, these angles are not applied meaning that the output is in the local space of the exemplar.

We use a branch of the encoder to perform the alignment task, allowing to jointly align images based on their visual features.

A byproduct is that the encoder returns angular distance to fronto-parallelity, which could be used to guide users during capture.

## 5 RESULTS

### 5.1 Dataset

We created an extended a dataset of flash images for testing and training of our approach. It comprises of 356 images of various types of materials we captured using four different smartphones. We reserve 50 images for testing, augmented by all images from Aittala et al. [2016]. Hence, no image from Aittala et al. [2016] was used for training.

### 5.2 Quantitative Evaluation

For quantitative analysis we compare our approach to a range of alternative methods with respect to different metrics.

*Methods.* We compare to five methods by (i) Aittala et al. [2016], (ii) Deschaintre et al. [2018], (iii) Gao et al. [2019], (iv) Guo et al. [2020b], and (v) Zhao et al. [2020]. All renderings of these methods are done with the material model described in their respective paper. While Gao et al. [2019] and Guo et al. [2020b] were designed to be compatible with multiple image acquisition with known light positions, in our comparisons we provide the same input as to our method: a single input image and an approximate light position.

*Metrics.* We quantify *style*, *diversity*, and *computational speed*. Style is captured by L1 difference of the VGG Gram matrices of rendered images. A good agreement in style has a low number i.e., less is better. We also evaluate XYZ histogram L1 difference and find that all methods have below 1% of difference with Ground Truth renderings, indicating good color matching for all. Histogram difference does not however capture the complex visual difference when comparing materials (as can be seen in Fig. 7). Diversity is captured as the mean pairwise VGG L1 across all realizations [Henzler et al.



Fig. 4. **Infinite spatial extent.** Top result is sampled at high resolution (256×4096) from our BRDF space while the bottom result is a result from Guo et al. [2020b] at 256×256 resolution and horizontally tiled 16 times to achieve high resolution. The absence of repetitiveness in the top result demonstrates that our learned BRDF space can be sampled at any query (x,y) location without producing visible repetition artifact. By construction, our network architecture does not require any special boundary alignment to avoid tiling artifacts.

2020]. Here, more is better. The idea behind this diversity metric is, that a for a diverse method, two realizations should have a high difference. A direct pixel metric would be sensitive to noise which generates small perturbations resulting in false-positive differences. Hence, the choice of VGG features, which detects if realizations are indeed perceivably different. Note that we do not evaluate pixel-wise metrics such as L1 or SSIM as these enforce local coherence, which is, by construction, not targeted by our method.

*Comparisons.* We use the described metrics to compare against multiple state-of-the-art methods in material acquisition and report the results on real (Flash) and synthetic materials (Relit) in Tbl. 2. For real results, we only have access to the frontal flash-illuminated material and therefore compare the picture to a rendering of each method’s result also under frontal illumination.

This, however, does not evaluate well the appearance under novel illumination, which is a key property of svBRDFs. To validate the generalization across light directions, we acquire 30 random stationary synthetic svBRDFs from CC0 Texture and render them to simulate a frontal-flash capture setup using Mitsuba2 [Nimier-David et al. 2019]. All methods are then run with this simulated flash image as input. We report the average of the re-lighting error, against ground truth renderings, for all methods under 10 random point light illuminations.

As shown in Tbl. 2, our approach is the only one to target diverse results, i.e., we produce infinitely many realizations of a texture while all other approaches produce only one. Thus, diversity (Div.) is zero for compared methods, while our approach can generate varied realizations for each material.

In terms of computational speed, Aittala et al. [2016] and Zhao et al. [2020] both require long –between 1 and 3 hours– per-exemplar optimization to produce a stationary texture. Our approach requires around 500 ms to generate a material and a few minutes to fine-tune it to a given input. This is in the same order of speed than Deschaintre et al. [2018] for generation and Gao et al. [2019] and Guo et al. [2020b] for the fine-tuning. Once fine-tuned, our method can generate new realizations and high resolutions versions of the targeted material in around 500 ms.

Table 2. We compare to recent material acquisition approaches on the L1 difference between VGG Gram matrices (VGG Style, lower is better) on both real and synthetic results as described in Sec. 5.2. Additionally we evaluate each method’s capacity to generate diverse realizations of a material with the mean pairwise VGG L1 across all realizations (Div, higher is better). We see that ours outperforms others on perceived similarity with the VGG style metric. Additionally, ours is the only one generating diverse material variations from a single image.

Method	Style err. ↓		Div. ↑
	Flash	Relit	
Aittala et al. [2016]	0.922	0.512	0.00
Deschaintre et al. [2018]	0.943	0.653	0.00
Gao et al. [2019]	0.738	0.556	0.00
Zhao et al. [2020]	<b>0.545</b>	0.618	0.00
Guo et al. [2020b]	0.843	0.582	0.00
Ours	0.597	<b>0.439</b>	<b>2.08</b>

### 5.3 Qualitative Evaluation

*Decomposition.* A qualitative example of our svBRDF decomposition (Normal, Diffuse Roughness and Specular maps) and re-renderings under different lights are depicted in Fig. 5. Please see our supplemental material for all results decomposition and comparison. We see that our method captures best the material behaviour and does not suffer from artefact in the over-exposed area of the input image which can be seen in previous work. As our method uses materials statistics rather than direct pixel aligned image to material transformation, it is immune to such artefacts.

*Relighting.* In Fig. 7 we show qualitative rendering comparisons on real materials with illumination coming from the top. In this more challenging setting, it is clear that existing work struggle to remove the highlight from the center of the flash image, which does not affect our method. As Aittala et al. [2016] reconstruct a small (representative) patch of the large input picture, their method is also immune to flash artefacts, but result in a very zoomed representation of the material. To compensate for this "zoom factor", we tile the results in each direction. We empirically found that 3 times works best for most materials.

*Seeds.* In Fig. 8 we show the variation of our results when changing the seed. The overall appearance of the material remains the

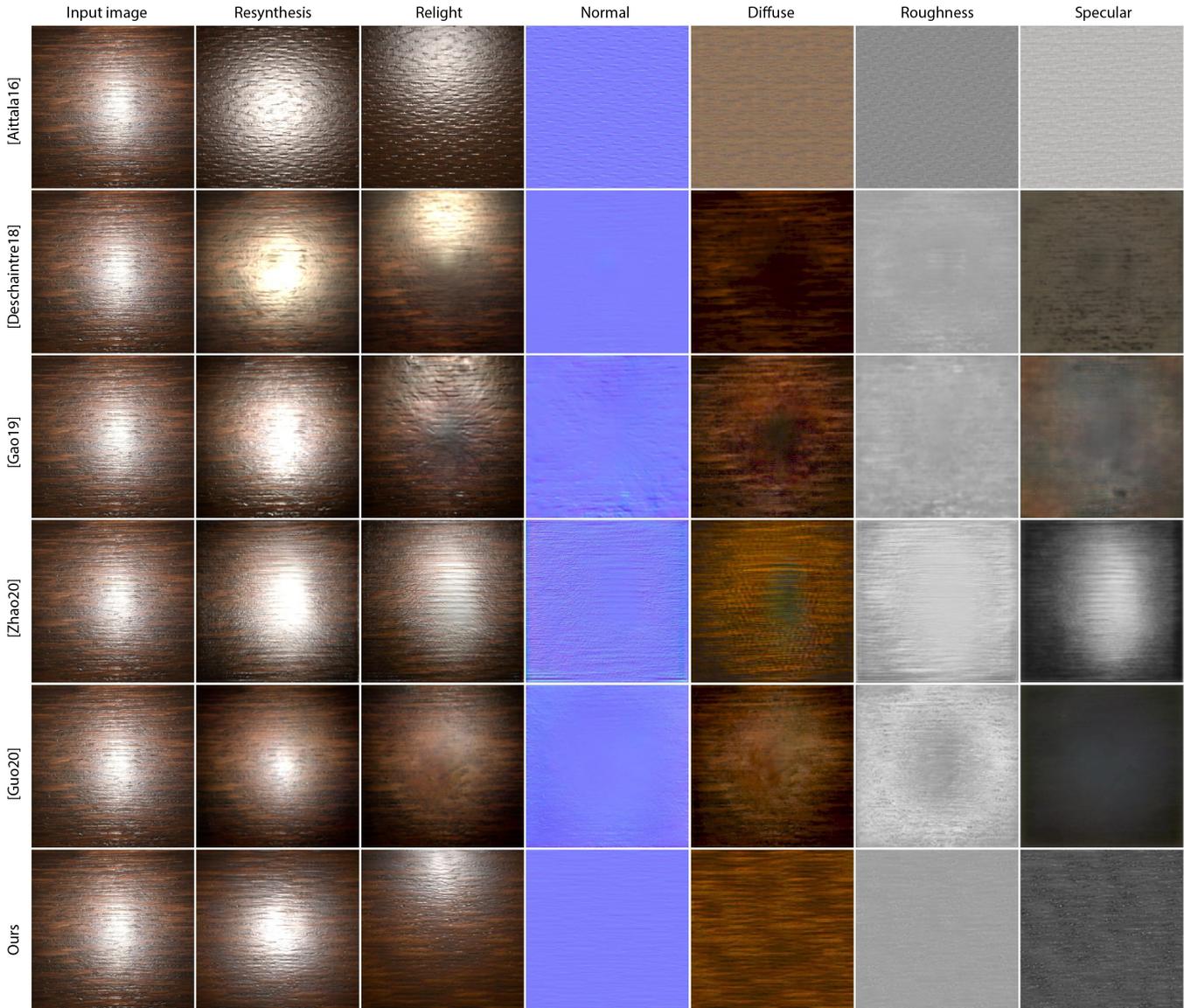


Fig. 5. **Comparison with other methods.** Each method (rows) decomposes an image into svBRDF parameters (columns). The first column shows the flash image input and the second column the rendering of the results under a similar fronto-parallel lighting. The third column is the material relit from the top, showing the generalization capacity across light. Our method’s quality is particularly visible under a novel illumination (see also Fig. 7). This is because other methods leave a trace of the flash in the svBRDF maps, as can be seen in the decomposed channels (four right-most columns). These results are obtained with our single image setting, compared methods Gao et al. [2019] and Guo et al. [2020b] could benefit from additional aligned images or accurate light calibration when available. Please see the supplemental material for similar results on many more materials.

same but the details (such as the rust or the leather normals and color variation) vary.

Overall, we see in Fig. 5, Fig. 7 and Fig. 8 that our approach can capture a large range of different stationary materials, reproducing their style, yet being diverse. This enables different properties described next.

*Infinite.* We show in Fig. 4 the “infinite” resolution capacity of our approach against the common approach of tiling. Our result (top image) shows no sign of repetitiveness even for very large resolution ( $4096 \times 256$ ).

*Interpolation.* We show results of interpolation between materials, as described in Sec. 4.5, in Fig. 6 and Supplemental Material. We compare against the linear interpolation baseline and Guo et al.

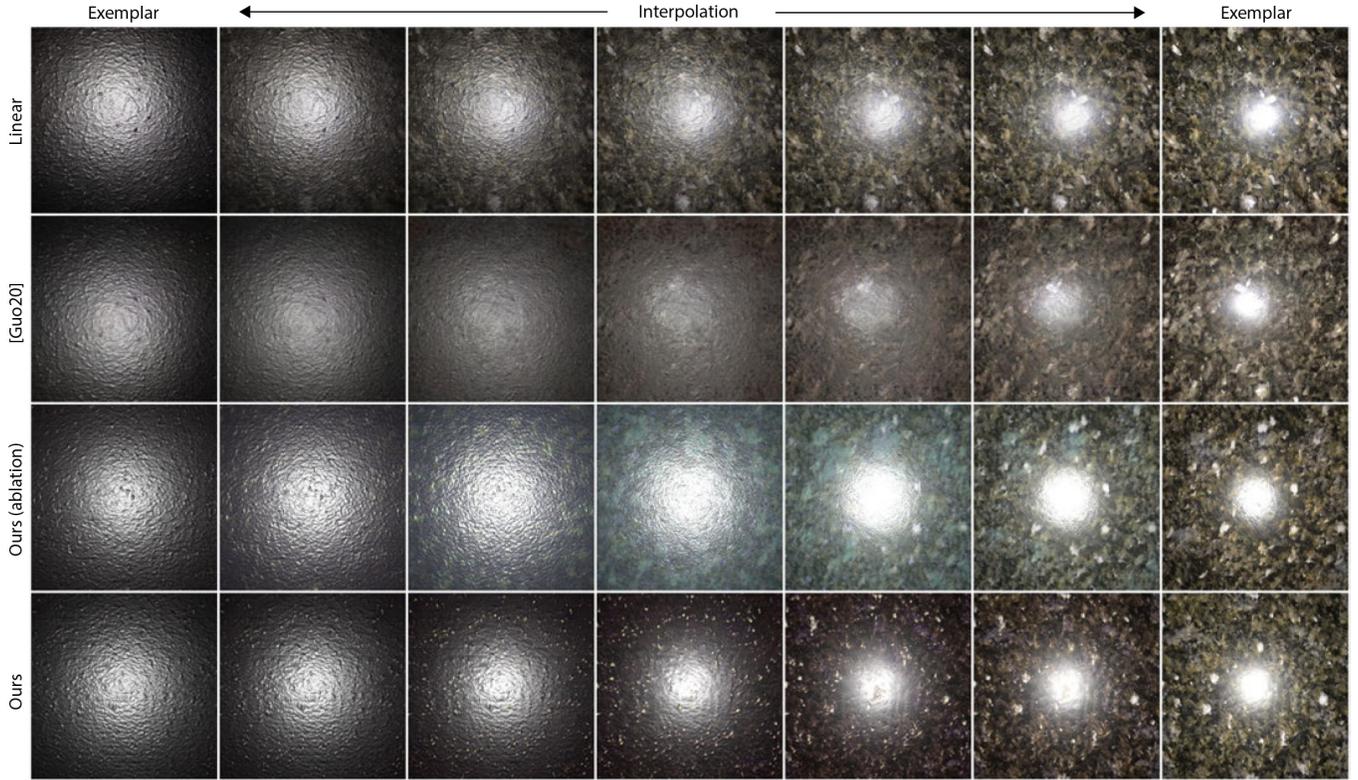


Fig. 6. **Interpolation of latent BRDF texture codes.** In each row, a left and a right latent code and generator weights  $z_1, z_2, g_1, g_2$  are obtained by encoding two flash images, respectively. The intermediate, continuous field of BRDF parameters is computed by interpolating, in the learned BRDF space, from  $z_1$  &  $g_1$  to  $z_2$  &  $g_2$  and conditioning the decoder Convolutional Neural Network (CNN) with the intermediate codes. The result is lit with a fronto-parallel light source to demonstrate the changes in appearance. For comparison, the first row shows image space linear interpolation, the second compares to Guo et al. [2020b]. The third row shows an ablation of our approach trained on a single material (without previous full dataset training). This lack of training prevents it from creating a cohesive space in which to interpolate. Overall our approach allows for interpolation, progressively changing both structure and reflectance.

[2020b] which also allows interpolation. We find our method to provide smoother interpolation than the Linear approach and to better preserve intermediate material color than Guo et al. [2020b]. We additionally evaluate interpolation if we directly train on material individually (without the training step described in Sec. 4.4). This confirms that this pre-training forms a coherent latent space in which we can navigate.

*Texturing.* Fig. 1 shows examples of applying maps produced by our approach to a complex 3D shape. Thanks to our generative model, we can easily texture many sneakers, without spatial or material repetition. At any point, a user can randomize the generated material, generate new materials from pictures or interpolate between new materials and old ones.

*Generation.* Our  $z$  space can be sampled to generate new materials as shown in in Fig. 9 with a variety of examples

*Interactive demo.* The visual quality is best inspected from our interactive WebGL demo in the supplemental material. It allows exploring the space by relighting, changing the random seed and visualizing individual BRDF model channels and their combinations.

The same package contains all channels of all materials as images as well as compared methods. See the accompanying video for a demonstration of our interactive interface.

*Fine-tuning.* We show the results quality improvement when using the proposed fine-tuning approach in Fig. 10. We can see that the structure and details better match the input picture.

#### 5.4 Ablation Experiments

We study several variants of our approach to evaluate the relevance of individual contributions to our FULL method.

We report the results of these evaluation in Tbl. 3 with VGG Style error in Sec. 5.2. We did not find the diversity of our method to be affected by these ablations.

Table 3. VGG style error for ablations relative to our FULL. For reference, our full method has an absolute score of 0.44.

Ablation	Error ↓
SINGLE	-0.5 %
NONTUNED	+24.0 %
DECODERONLY	+2.0 %
FOURIER	+0.9 %
LIGHT	+1.7 %

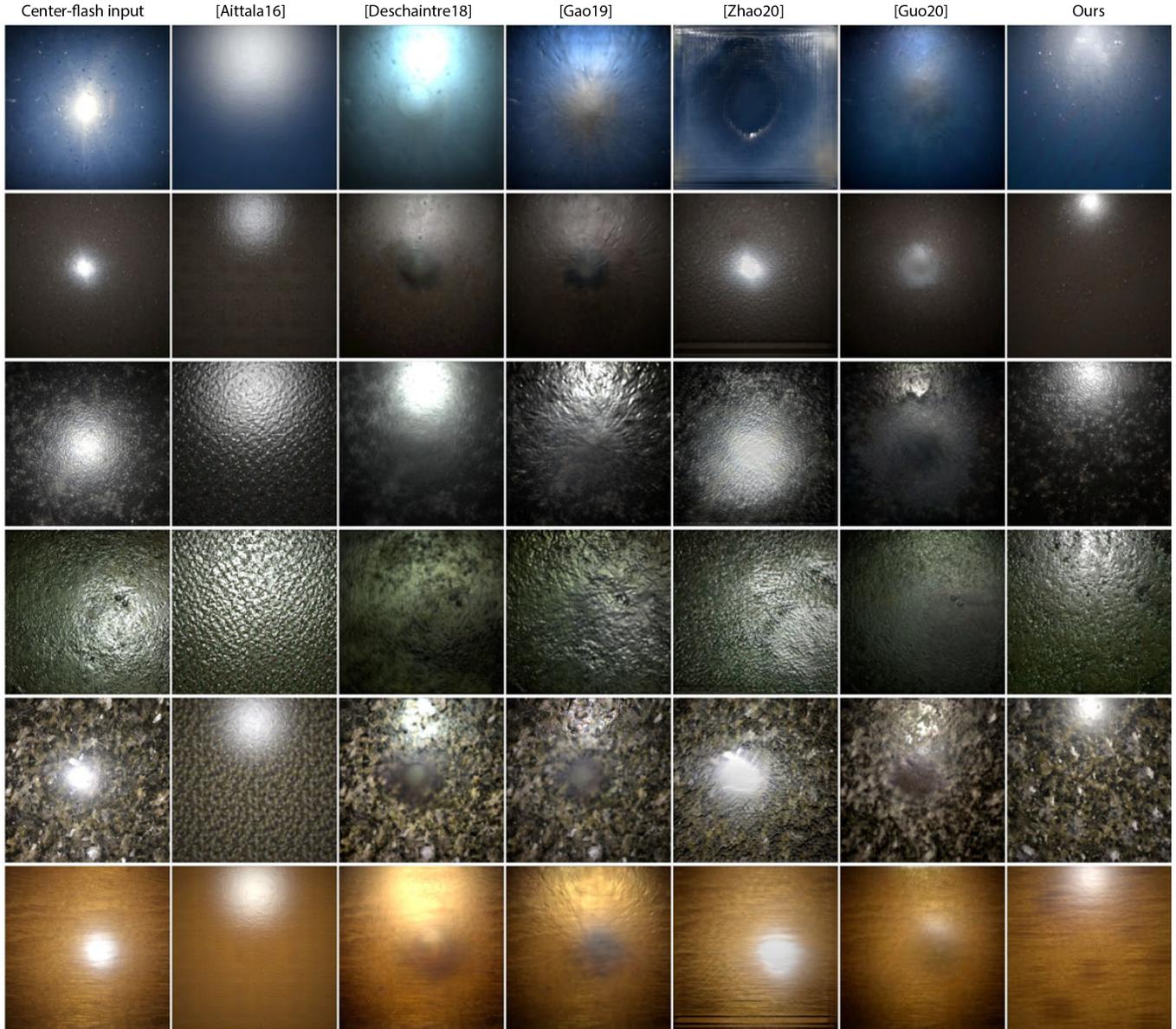


Fig. 7. **Relighting** of different materials (rows) using material maps extracted by different methods (columns). The first column shows the input flash image where light is fronto-parallel. The light in all other images comes from the top. While no reference is available for this task, it is apparent that all the methods except ours struggle to generalize to novel light conditions. Note that Deschaintre et al. [2018], Gao et al. [2019] and Guo et al. [2020b] leave a dark residual of the flash in the material maps. Zhao et al. [2020] and Aittala et al. [2016] fare slightly better and avoid the residual, but the structures do not match. These results are obtained with our single image setting, compared methods Gao et al. [2019] and Guo et al. [2020b] could benefit from additional aligned images or accurate light calibration when available.

SINGLE describes our method trained on a single example, without the previous training step. The results are slightly better than our FULL method, but requires twice longer per material training and does not generalize to a space, preventing interpolation and generation of materials.

NONTUNED is our method without the fine-tuning step from Sec. 4.5, confirming that it significantly improves the match to the

acquired material. DECODERONLY describes the change of our generator to a decoder only architecture. We show that removing the encoder part of the generator slightly degrades the results. FOURIER and LIGHT respectively result from the removal of the Fourier component (power spectrum) of our loss (Sec. 4.3) and the removal of the light alignment branch of our encoder (described in Sec. 4.7), which both lead to slightly worse results.

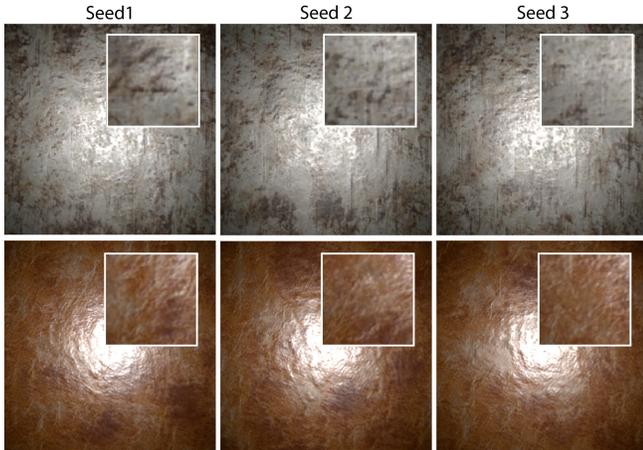


Fig. 8. **Seeds variation.** We vary the seed for the generation of different realizations for acquired materials, while preserving its overall appearance. The zoomed-in insets all show the same region of the material, allowing to better appreciate the variations.



Fig. 9. **Generation.** Random samples from our space. We generate new materials by sampling the  $z$  space and render them with a frontal flash. See supplemental materials for more generated materials.

## 6 USER EXPERIMENT

We perform a user study to better understand the capabilities of different methods. Our main aim is to provide material maps that robustly generalize to all light conditions so they can be deployed in production rendering. Hence we study a *relighting task*: given a the input image in one light condition, we ask humans to pick the method that looks most plausible “in a different light”.

*Methods.* Subjects anonymously completed an online form with-out time limit. At the start of the user study, participants were shown two photos of a marble material taken under two different lighting conditions to exemplify what a valid relighting could look like. They performed 10 trials, each corresponding to one material. In each trial, they were presented a reference image rendered in one

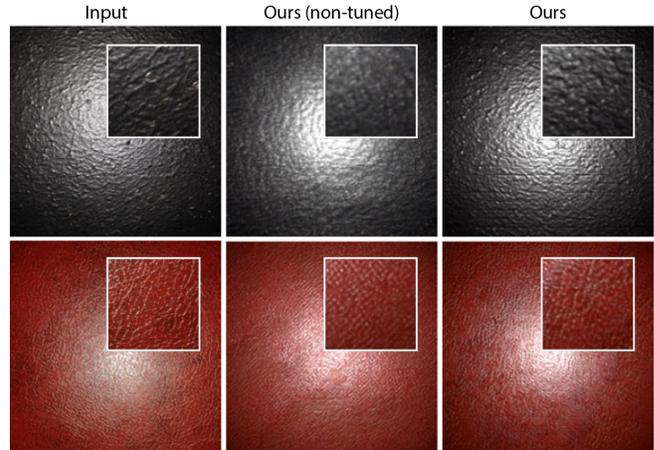


Fig. 10. **Fine-tuning.** We show results on two results of real materials reproduced using our pre-trained network (ours non-tuned) and the same material using our fine-tuning approach. We can see that our fine-tuned results match the input material appearance significantly better. Note that fine-tuning is only with image supervision and does not have access to any underlying BRDF supervision.

light conditions (“flash”) and six relit images in another light condition (“top”). Relit images were displayed in a randomized spatial 2D layout. We consider six different methods: Aittala et al. [2016]; Deschaintre et al. [2018]; Gao et al. [2019]; Guo et al. [2020b]; Zhao et al. [2020] and ours. Samples of those stimuli are seen in Fig. 7. Participants were asked to pick the image (images were not named) that, according to them, was the best faithful relighting of the source (flash) image. Note that no relit reference was shown.

*Analysis.* A total of  $N = 39$  participants completed the experiment as summarized in Tbl. 4. A  $\chi^2$  test rejects ( $p < 0.0001$ ) the hypothesis that choices were random. Pair-wise binomial post-hoc tests further show that our method is different from any other method, at the same significance level. Most importantly, subjects choose our method in 314 out of 390 total answers (80.5 %). We did not analyze the relation of other methods relative to each other.

Table 4. User preferences per method.

Method	Freq.
Aittala et al. [2016]	21
Deschaintre et al. [2018]	10
Gao et al. [2019]	4
Guo et al. [2020b]	11
Zhao et al. [2020]	30
Ours	314

## 7 LIMITATIONS

Our method relies on fronto-parallel flash acquisition. While we propose a mitigation solution in Sec. 4.7, we show in Fig. 11 that we are not completely invariant to large light and plane rotations. Our approach is also limited to stationary isotropic materials and relies on the planarity of the captured surface.

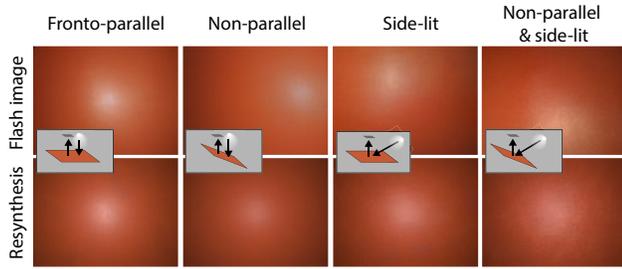


Fig. 11. **Flash acquisition assumption.** We show example of how our results degrade when fronto-parallel, collocated flash assumptions are broken. The recovered material appearance varies (roughness, high frequency normal) but maintains the overall appearance of the input image.

## 8 CONCLUSION

We have presented an approach to generate a space of BRDF textures using a small set of flash images in an unsupervised way. Comparing this approach to the literature shows competitive metrics for re-renderings with the unique advantage of being able to generate an infinite and diverse field of BRDF parameters.

In the future, it would be interesting to increase the complexity of supported material whether in term of shading or non stationarity. Also, not relying on fine-tuning to increase the network expressiveness would allow to create an even more cohesive space. Further, more refined differentiable rendering material models could be used to derive stochastic textures, including shadows, displacement, or scattering as well as volumetric or time-varying textures. We believe that our framework will represent a stepping stone for more complex infinite and diverse BRDFs acquisition.

## ACKNOWLEDGEMENT

This work was supported by the ERC Starting Grant SmartGeometry, a Google AR/VR Research Award, Dr. Abhijeet Ghosh and his EPSRC Early Career Fellowship (EP/N006259/1) and a GPU donation by NVIDIA Corporation. We thank the authors of Aittala et al. [2016], Gao et al. [2019], Guo et al. [2020b], Zhao et al. [2020] for providing their implementations and helping with the comparisons.

## REFERENCES

Miika Aittala, Timo Aila, and Jaakko Lehtinen. 2016. Reflectance modeling by neural texture synthesis. *ACM Trans Graph (Proc. SIGGRAPH)* 35, 4 (2016).

Brett Allen, Brian Curless, and Zoran Popović. 2003. The space of human body shapes: reconstruction and parameterization from range scans. *ACM Trans Graph (Proc. SIGGRAPH)* 22, 3 (2003), 587–94.

Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. 2009. Patch-Match: A randomized correspondence algorithm for structural image editing. *ACM Trans Graph (Proc. SIGGRAPH)* 28, 3 (2009).

Urs Bergmann, Nikolay Jetchev, and Roland Vollgraf. 2017. Learning texture manifolds with the periodic spatial GAN. In *J MLR*. 469–77.

Volker Blanz, Thomas Vetter, et al. 1999. A morphable model for the synthesis of 3D faces. In *Siggraph*, Vol. 99. 187–194.

Mark Boss, Varun Jampani, Kihwan Kim, Hendrik P.A. Lensch, and Jan Kautz. 2020. Two-shot Spatially-varying BRDF and Shape Estimation. In *CVPR*.

Robert L Cook and Tony DeRose. 2005. Wavelet noise. *ACM Trans Graph* 24, 3 (2005), 803–11.

Robert L Cook and Kenneth E Torrance. 1982. A reflectance model for computer graphics. *ACM Trans Graph* 1, 1 (1982), 7–24.

Valentin Deschaintre, Miika Aittala, Frédo Durand, George Drettakis, and Adrien Bousseau. 2018. Single-image svbrdf capture with a rendering-aware deep network.

*ACM Trans Graph (Proc. SIGGRAPH)* 37, 4 (2018).

Valentin Deschaintre, Miika Aittala, Frédo Durand, George Drettakis, and Adrien Bousseau. 2019. Flexible SVBRDF Capture with a Multi-Image Deep Network. *Comp Graph Forum* 38, 4 (2019).

Valentin Deschaintre, George Drettakis, and Adrien Bousseau. 2020. Guided Fine-Tuning for Large-Scale Material Transfer. *Comp Graph Forum (Proc. EGSR)* 39, 4 (2020).

Valentin Deschaintre, Yiming Lin, and Abhijeet Ghosh. 2021. Deep polarization imaging for 3D shape and SVBRDF acquisition. In *CVPR*.

Alexei A Efros and Thomas K Leung. 1999. Texture synthesis by non-parametric sampling. In *ICCV*, Vol. 2.

Duan Gao, Xiao Li, Yue Dong, Pieter Peers, Kun Xu, and Xin Tong. 2019. Deep inverse rendering for high-resolution SVBRDF estimation from an arbitrary number of images. *ACM Trans Graph (Proc. SIGGRAPH Asia)* 38, 4 (2019).

Leon Gatys, Alexander S Ecker, and Matthias Bethge. 2015. Texture synthesis using convolutional neural networks. In *NIPS*.

Stamatios Georgoulis, Konstantinos Rematas, Tobias Ritschel, Efstratios Gavves, Mario Fritz, Luc Van Gool, and Tinne Tuytelaars. 2017. Reflectance and natural illumination from single-material specular objects using deep learning. *PAMI* 40, 8 (2017), 1932–47.

Darya Guarnera, Giuseppe Claudio Guarnera, Abhijeet Ghosh, Cornelia Denk, and Mashhuda Glencross. 2016. BRDF representation and acquisition. *Comp Graph Forum* 35, 2 (2016), 625–50.

Jie Guo, Shuichang Lai, Chengzhi Tao, Yuelong Cai, Lei Wang, Yanwen Guo, and Ling-Qi Yan. 2021. Highlight-aware two-stream network for single-image SVBRDF acquisition. *ACM Trans Graph (Proc. SIGGRAPH)* 40, 4 (2021).

Yu Guo, Miloš Hašan, Lingqi Yan, and Shuang Zhao. 2020a. A Bayesian Inference Framework for Procedural Material Parameter Estimation. *Comp Graph Forum* 39, 7 (2020).

Yu Guo, Cameron Smith, Miloš Hašan, Kalyan Sunkavalli, and Shuang Zhao. 2020b. MaterialGAN: Reflectance Capture using a Generative SVBRDF Model. *ACM Trans Graph* 39, 6 (2020).

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *CVPR*. 770–8.

David J Heeger and James R Bergen. 1995. Pyramid-based texture analysis/synthesis. In *Proc. SIGGRAPH*. 229–38.

Eric Heitz. 2014. Understanding the Masking-Shadowing Function in Microfacet-Based BRDFs. *Journal of Computer Graphics Techniques (JCGT)* 3, 2 (30 June 2014), 48–107. <http://jcg.org/published/0003/02/03/>

Philipp Henzler, Niloy J Mitra, and Tobias Ritschel. 2020. Learning a Neural 3D Texture Space from 2D Exemplars. *CVPR* (2020).

Yiwei Hu, Julie Dorsey, and Holly Rushmeier. 2019. A Novel Framework for Inverse Procedural Texture Modeling. *ACM Trans. Graph.* 38, 6 (2019).

Xun Huang and Serge Belongie. 2017. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*. 1501–10.

Justin Johnson, Alexandre Alahi, and Li Fei-Fei. 2016. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*.

Bela Julesz. 1965. Texture and visual perception. *Scientific American* 212, 2 (1965).

Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *CVPR*. 4401–10.

Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv:1312.6114* (2013).

Vivek Kwatra, Irfan Essa, Aaron Bobick, and Nipun Kwatra. 2005. Texture optimization for example-based synthesis. *ACM Trans Graph* 24, 3 (2005).

Xiao Li, Yue Dong, Pieter Peers, and Xin Tong. 2017. Modeling surface appearance from a single photograph using self-augmented convolutional neural networks. *ACM Trans Graph* 36, 4 (2017).

Xiao Li, Yue Dong, Pieter Peers, and Xin Tong. 2019. Synthesizing 3D Shapes From Silhouette Image Collections Using Multi-Projection Generative Adversarial Networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Zhengqin Li, Kalyan Sunkavalli, and Manmohan Chandraker. 2018a. Materials for masses: SVBRDF acquisition with a single mobile phone image. In *ECCV*. 72–87.

Zhengqin Li, Zexiang Xu, Ravi Ramamoorthi, Kalyan Sunkavalli, and Manmohan Chandraker. 2018b. Learning to reconstruct shape and spatially-varying reflectance from a single image. *ACM Trans Graph (Proc. SIGGRAPH Asia)* 37, 4 (2018).

Guilin Liu, Duygu Ceylan, Ersin Yumer, Jimei Yang, and Jyh-Ming Lien. 2017. Material editing using a physically based rendering network. In *ICCV*. 2261–9.

Gang Liu, Yann Gousseau, and Gui-Song Xia. 2016. Texture synthesis through convolutional neural networks and spectrum constraints. In *ICPR*. 3234–9.

Benoit B Mandelbrot. 1983. *The fractal geometry of nature*. Vol. 173. WH Freeman New York.

W Matusik. 2003. A data-driven reflectance model. *ACM Trans Graph* 22, 3 (2003), 759–69.

Wojciech Matusik, Matthias Zwicker, and Frédo Durand. 2005. Texture design using a simplicial complex of morphable textures. *ACM Trans Graph (Proc. SIGGRAPH)* 24,

- 3 (2005).
- Giljoo Nam, Joo Ho Lee, Diego Gutierrez, and Min H. Kim. 2018. Practical SVBRDF Acquisition of 3D Objects with Unstructured Flash Photography. *ACM Trans. Graph.* 37, 6 (2018).
- Chuong H Nguyen, Tobias Ritschel, and Hans-Peter Seidel. 2015. Data-driven color manifolds. *ACM Trans Graph* 34, 2 (2015).
- Merlin Nimier-David, Delio Vicini, Tizian Zeltner, and Wenzel Jakob. 2019. Mitsuba 2: A retargetable forward and inverse renderer. *ACM Trans Graph (Proc. SIGGRAPH Asia)* 38, 6 (2019).
- Keunhong Park, Konstantinos Rematas, Ali Farhadi, and Steven M Seitz. 2019. Photoshape: photorealistic materials for large-scale shape collections. *ACM Trans Graph (Proc. SIGGRAPH Asia)* 37, 6 (2019).
- Ken Perlin. 1985. An Image Synthesizer. *SIGGRAPH Comput Graph* 19, 3 (1985).
- Bui Tuong Phong. 1975. Illumination for computer generated pictures. *Commun ACM* 18, 6 (1975), 311–17.
- Javier Portilla and Eero P Simoncelli. 2000. A parametric texture model based on joint statistics of complex wavelet coefficients. *Int J Comp Vis* 40, 1 (2000), 49–70.
- K. Rematas, T. Ritschel, M. Fritz, E. Gavves, and T. Tuytelaars. 2016. Deep Reflectance Maps. In *CVPR*.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*. 234–41.
- Christophe Schlick. 1994. An inexpensive BRDF model for physically-based rendering. *Comp Graph Forum* 13, 3 (1994), 233–46.
- Christopher Schwartz, Ralf Sarlette, Michael Weinmann, and Reinhard Klein. 2013. DOME II: A Parallelized BTF Acquisition System. In *MAM*. 25–31.
- Omyr Sendik and Daniel Cohen-Or. 2017. Deep correlations for texture synthesis. *ACM Trans Graph* 36, 5 (2017).
- Ana Serrano, Diego Gutierrez, Karol Myszkowski, Hans-Peter Seidel, and Belen Masia. 2016. An intuitive control space for material appearance. *ACM Trans Graph (Proc. SIGGRAPH Asia)* 35, 6 (2016).
- Tamar Rott Shaham, Tali Dekel, and Tomer Michaeli. 2019. Singan: Learning a generative model from a single natural image. In *ICCV*.
- Liang Shi, Beichen Li, Miloš Hašan, Kalyan Sunkavalli, Tamy Boubekeur, Radomir Mech, and Wojciech Matusik. 2020. Match: Differentiable Material Graphs for Procedural Material Capture. *ACM Trans. Graph.* 39, 6 (2020).
- Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv:1409.1556* (2014).
- Dmitry Ulyanov, Vadim Lebedev, Andrea Vedaldi, and Victor S Lempitsky. 2016. Texture Networks: Feed-forward Synthesis of Textures and Stylized Images.. In *ICML*.
- Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. 2017. Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis. In *CVPR*.
- Bruce Walter, Stephen R Marschner, Hongsong Li, and Kenneth E Torrance. 2007. Microfacet models for refraction through rough surfaces. In *Proc. EGSR*. 195–206.
- Li-Yi Wei and Marc Levoy. 2000. Fast texture synthesis using tree-structured vector quantization. In *Proc. SIGGRAPH*.
- Wenjie Ye, Xiao Li, Yue Dong, Pieter Peers, and Xin Tong. 2018. Single Image Surface Appearance Modeling with Self-augmented CNNs and Inexact Supervision. *Comp Graph Forum* 37, 7 (2018), 201–11.
- Yezi Zhao, Beibei Wang, Yanning Xu, Zheng Zeng, Lu Wang, and Nicolas Holzschuch. 2020. Joint SVBRDF Recovery and Synthesis From a Single Image using an Unsupervised Generative Adversarial Network. In *EGSR*.
- Xilong Zhou and Nima Khademi Kalantari. 2021. Adversarial Single-Image SVBRDF Estimation with Hybrid Training. *Computer Graphics Forum (Proc. Eurographics)* (2021).
- Yang Zhou, Zhen Zhu, Xiang Bai, Dani Lischinski, Daniel Cohen-Or, and Hui Huang. 2018. Non-stationary texture synthesis by adversarial expansion. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 1–13.
- Károly Zsolnai-Fehér, Peter Wonka, and Michael Wimmer. 2018. Gaussian material synthesis. *ACM Trans. Graph.* 37, 4 (2018).