arXiv:2008.01808v1 [cs.CV] 4 Aug 2020

# HIGH RESOLUTION NEURAL TEXTURE SYNTHESIS WITH LONG RANGE CONSTRAINTS

NICOLAS GONTHIER$^{†*}$, YANN GOUSSEAU$^{†}$, AND SAÏD LADJAL$^{†}$

**Abstract.** The field of texture synthesis has witnessed important progresses over the last years, most notably through the use of Convolutional Neural Networks. However, neural synthesis methods still struggle to reproduce large scale structures, especially with high resolution textures. To address this issue, we first introduce a simple multi-resolution framework that efficiently accounts for long-range dependency. Then, we show that additional statistical constraints further improve the reproduction of textures with strong regularity. This can be achieved by constraining both the Gram matrices of a neural network and the power spectrum of the image. Alternatively one may constrain only the autocorrelation of the features of the network and drop the Gram matrices constraints. In an experimental part, the proposed methods are then extensively tested and compared to alternative approaches, both in an unsupervised way and through a user study. Experiments show the interest of the multi-scale scheme for high resolution textures and the interest of combining it with additional constraints for regular textures.

**Key words.** Texture Synthesis, Deep Neural Network, High Resolution, Perceptual Evaluation, Multi-scale

**1. Introduction.** Examplar-based texture synthesis consists in automatically generating sample images from a given example texture image. These samples are required to be visually faithful to the example and as diverse as possible. For more than forty years, and despite its inherent ill-posedness, this problem has been a fruitful way to test visually the validity of various mathematical models, ranging from time series [26], Markov random fields [4] to wavelet decompositions [17, 28] or non-parametric Markovian modeling [10]. More recently, Convolutional Neural Networks have permitted impressive progresses in the field, initiated by the work by Gatys et al. [12], itself followed by numerous contributions, e.g. [39, 24, 35].

One challenge that has been faced by all methods since the early days of texture synthesis is the multi-scale nature of texture samples, implying that models should be able to reproduce both small and large scales, possibly over several orders of magnitude. For instance, parametric models for Markov fields are known to be intrinsically badly suited to a multi-scale modeling. Zooming such a model by a given factor implies extremely heavy computations to derive the corresponding parameters [15], inpairing the design of multi-scale such models. Wavelet models are more adapted by nature to multi-scale modeling, but the faithful reproduction of structured textures requires complex interactions between scales to be accounted for. The best such modeling up to date is the second order statistical model proposed in [28], but highly structured textures still represent a challenge to such approaches. Non parametric Markov modeling methods such as those presented in [10] or [9] indeed have the ability to deal simultaneously with several scales, albeit at a high computational cost. However, they are also well known to produce textures with very little variety, often producing verbatim copies, see [1] and the experiments in the present paper. The methods relying on convolutional neural networks, following the seminal work by Gatys et al. [12], are currently the most efficient to capture multi-scale structures. Nevertheless, they still lacks efficiency when large scale regularity is needed, as we will see in detail in this paper. Moreover, they are prone to generate artefacts that

---

$^{*}$Université Paris-Saclay, 91190, Saint-Aubin, France

$^{†}$LTCI, Télécom Paris, Institut polytechnique de Paris, 19 Place Marguerite Perey, 91120 Palaiseau, France (nicolas.gonthier@telecom-paris.fr).

prevent a satisfactory reproduction of small scale structures.

In this work, we present several neural synthesis methods that significantly improves the ability to preserve the large scale organisation of textures. We first propose a simple multi-resolution framework that account for large-scale structures and permits the synthesis of high resolution images. We then show that, in this multi-resolution framework, additional constraints are useful in the case of regular textures. A first approach combines the classical statistical constraints of neural approaches [12] (Gram matrices) with Fourier frequency constraints, similar to those introduced by [11]. A preliminary, mono-scale version of this idea was presented in a conference paper [24]. Alternatively, the multi-resolution framework can be combined with a statistical constraint relying on the full auto-correlation of the features of the network. This approach is closely related to the one introduced in [32], which combines correlations with Gram matrices and various additional constraints. We show that correlation terms alone yield excellent results and therefore that Gram matrices are not necessary in this case.

We then evaluate the proposed methods in an extensive experimental section. The evaluation of texture synthesis results is a challenging task. Some approaches draw on well chosen statistics to estimate the quality of the results (the closest to the examplar, the better), as for instance discussed in [3]. In this paper, we first evaluate results in this manner, relying on Kullback-Leibler divergence between wavelet marginals, following the texture indexing scheme from [6]. Then, we also evaluate the proposed methodology through a perceptual user study. Indeed, it is shown in [3, 7, 8] through extensive experiments that feature-based evaluations do not approach well human-based visual evaluation of texture similarity, especially in the case of long range correlations, which is precisely one of the cases tackled in this paper. We therefore rely on a user study to compare the framework we propose to both the original method from [13] and some of its improvement that focus on the respect of large scale structures [35, 32].

**2. Neuronal texture synthesis.** A complete state-of-the-art on the subject of texture synthesis is out of the scope of this paper. In view of the method that we propose in this work, we focus in this section on the works involving CNNs that have followed the seminal contribution of Gatys et al. [12] and particularly on works proposing new statistical constraints and focusing on long-range structure.

*Accelerations and alternative sampling strategy.* In a first direction, several works have proposed ways to speed-up the synthesis process, notably through feed forward networks [39, 40, 19, 33]. In [18], Generative Adversarial Networks are used to synthesize textures. Such methods enables fast synthesis once the networks have been trained for specific textures, but the quality of results is still inferior to the original approach [12], especially for structured textures. Zhu and other authors have proposed an evolution of the FRAME model [45] in the context of neural networks [25] under the name *DeepFrame*. Textures are synthesized from an exponential model using features from a neural network. In [5], this macrocanonical approach is pushed further and fully analyzed theoretically. It is worth noting that both approaches [25, 5] rely on first order constraints on features and therefore drop the use of the Gram matrices.

*Statistical constraints and losses.* In a different direction, a large body of works has been dedicated to add additional constraints to the synthesis, often relying on new or modified loss functions. In [14], the color of the synthesis is constrained to specified values. In [30], it is proposed to constrain the histograms of some feature maps, in order to reduce halo artefacts. In [19], a total variation term is added

2

in the loss function for perceptual reasons. Other works such as [2, 24, 32] also propose alternative losses to add further statistical constraints. Since they explicitly deal with long range dependency and structure, they will be reviewed in the next paragraph. It should be noted that these approaches propose to combine several statistical constraints by adding them to get the final loss function. Another possibility would be to alternate different projections as it is done in the seminal work of Portilla and Simoncelli [28]. Alternative constraints have also been investigated for the closely related task of style transfer. In [23], it is shown that matching Gram matrices reduces to kernel-based comparison of features, and various kernels are investigated in this setting. Other works investigate alternatives to the original Gram matrices, such as cross-layers (rather than within-layers) Gram matrices as in [43] or [27] (both inspired by [28]).

*Multi-scale neural synthesis.* Neural networks such as the VGG19 used in most texture synthesis methods intrinsically have a multi-scale structure by alternating convolutions, non-linearity and subsampling. However, as we will see in the experimental section, the size of the receptive fields in these networks is not sufficient to synthesize large scale structures, especially when the resolution increases. To the best of our knowledge, only one paper, [35], proposes to rely on a multi-scale strategy to synthesize high resolution textures. The idea is to feed the network with a multi-scale decomposition, in this case a Gaussian pyramid, instead of a single image. In this paper we will propose an alternative approach to the multi-scale neural synthesis and compare our results with [35].

*Incorporating long distance dependency.* In [2], long distance patterns are handled by adding in the loss function a cross-correlation term, made of the correlation between features maps and a shifted version of it. Different sets of shifts are used depending on the layer, up to about a sixth of the image size. In [27], in the context of style transfer, a similar idea is investigated using only one-pixel shifts. In [32], the same idea is pushed further, by considering all cross-correlations at once in order to impose long-range structure for regular textures. Several other terms (smoothness, diversity) are added to the loss function. This approach, to which we will compare our results, indeed yields long-range structure, nevertheless at the price of relatively strong artefacts. Apart from these work dealing with cross-correlation of features and closely related to the present paper, [24] proposed to incorporate the power spectrum in the loss function, thereby enabling the respect of highly structured textures. In a related work, [31], it is proposed to impose the spectrum constraint by using a windowed Fourier Transform, enabling non-stationnary behavior to be accounted for, at the cost of the inherent stationary nature of textures.

**3. Multi-scale spectral control for texture synthesis.** In this section, we detail our method to synthesize high quality texture images. After recalling in subsection 3.1, the classical approach from Gatys et al., we introduce a simple multi-scale framework in subsection 3.2, before presenting in subsection 3.3 the spectral constraint we propose in order to both control artefacts and preserve long-range structures. Finally, we present in subsection 3.4, the use of the autocorrelation of the feature maps as a potential alternative to the Gram matrices.

**3.1. Reminder on the work from [12].** The seminal work [12] is based on the idea that a network trained for classification purpose, in this case a VGG network as introduced in [34], can be repurposed for a synthesis task. Roughly speaking, the synthesis is achieved by backpropagation of texture-adapted statistical constraints from the inner layers of the network up to pixels of the synthesized image.

3

More precisely, the method works as follows. We consider a given convolutional neural network [1] consisting of $l$ layers. For a given color texture exemplar $I \in \mathbb{R}^{h \times w \times 3}$, where $h, w$ are the dimensions of the image, we write $f_l$ for the output (that is, the activations) of layer $l$. This output will be called a *feature map* from now on. Each feature map $f^l$ belongs to $\mathbb{R}^{h_l \times w_l \times m_l}$, where $w_l, h_l$ are the spatial dimension of layer $l$ and $m_l$ the number of channels of the feature. We further write $N = h \times w$ and $N_l = h_l \times w_l$ for the spatial dimensions of respectively the image and the feature maps.

To synthesis a new texture, some statistics are imposed on a subset[2] $\mathcal{S}$ of the layers of the CNN. The statistics considered in [13] are strongly inspired by the work from Portilla and Simoncelli [28] and rely on the so-called Gram matrices $G^l \in \mathbb{R}^{m_l \times m_l}$, defined for each couple $p, q \in \{1, \cdots, m_l\}$ as

$$(3.1) \qquad G_{p,q}^l = \frac{1}{N_l^2} \sum_{i=1}^{N_l} f_p^l(i) \cdot f_q^l(i) = \frac{1}{N_l^2} \langle f_q^l, f_p^l \rangle,$$

where $f_p^l \in \mathbb{R}^{N_l}$, for $p \in \{1, \cdots, m_l\}$, is the vectorized $p^{\text{th}}$ channel of feature $l$.

To generate a new texture image $\tilde{I}$ on the basis of a reference one $I$, a gradient descent is used, starting from a white noise image, to find an image that matches the reference statistics. Usually the L-BFGS-B [44] second-order optimization method is chosen.

The corresponding loss function on the features is defined as :

$$(3.2) \qquad \mathcal{L}_{Gram} = \sum_{l=1}^{L} \omega_l \| G^l - \tilde{G}^l \|_2^2,$$

with $\omega_l \in \mathbb{R}$ the weight of the layer $l$.

The loss function (3.2) can be seen as multi-objective cost functions agglomerated into a single-objective cost function. Although comparing different objectives is generally difficult, choosing identical weights, i. e. $\omega_l = 1 \; \forall l \in [1, L]$, yields perceptually acceptable results.

A central questions of texture synthesis is to identify the best sets of statistics to incorporate in this loss function and possibly the irreducible set of those statistics ([20]). Although the method from [12] yields synthesis results of unprecedented quality, a strong limitation is its inability to respect long range dependency, particularly when large scale structures have some regularity. This can be seen in first row of Figure 2. Neural networks such as VGG-19 have a multiscale structure, through alternating convolution and subsampling, that allow some large scale structures to be accounted for. Nevertheless, the size of the filters used in CNNs such as VGG-19, and therefore the size of the corresponding receptive fields, are small with respect to the size of the image especially when synthesizing high resolution images (here $1024 \times 1024$). As we have mentioned in the introduction, several works have addressed this limitation [24, 2, 27, 31, 32], but, as we will see in the experimental section, none is fully satisfactory. In the following sections, we propose several improvement of the original neural texture synthesis method in order to address this limitation.

---

[1]In this work, as in [12], we consider the VGG19 network [34] but other choices are possible, including networks with random weights [16].

[2]For simplicity's sake, we will consider layers from 1 to $L$ in the rest of this document but the user can choose non-consecutive layers in the network.

**3.2. Multi-scale synthesis.** The first modification we introduce to the method from [12] is a straightforward multi-scale framework that will help preserving the large scale organisation of images. This strategy is relatively classical for texture synthesis methods and has been used in the past in different settings [22, 38]. This approach is much simpler than the related method introduced in [35] and, as we will see in the experimental section, yields better results.

The idea is simply to first synthesize a coarse resolution image, which is then upsampled and given as initialization for a synthesis at the next scale. This process is repeated $K$ times until the desired resolution is reached. As illustrated in Figure 1, we first build an image pyramid from the examplar image $I$, iteratively down-sampling it by factors $2^1, 2^2, \cdots, 2^K$, resulting in images $I^{(1)}, I^{(2)}, \cdots, I^{(K)}$. A first synthesis result is obtained by using the smallest image as the examplar and white noise as initialization. Then, for step $k \in K, K-1, \cdots, 1$, we generate a new result using $I^{(k)}$ as the exemplar and the obtained synthesis result $\tilde{I}^{(k-1)}$ as the initialization instead of white noise. The upsampling of $\tilde{I}^{(k-1)}$ is performed using bilinear interpolation. The only parameter of this generic multi-scale framework is the number of scales $K$.

As can be seen in Figure 15, this strategy can yield strong improvements in some cases but is not enough to allow the reproduction of highly structured textures. In the next section, we show how the result can be improved by adding a careful control of the Fourier spectrum into the multi-scale scheme.
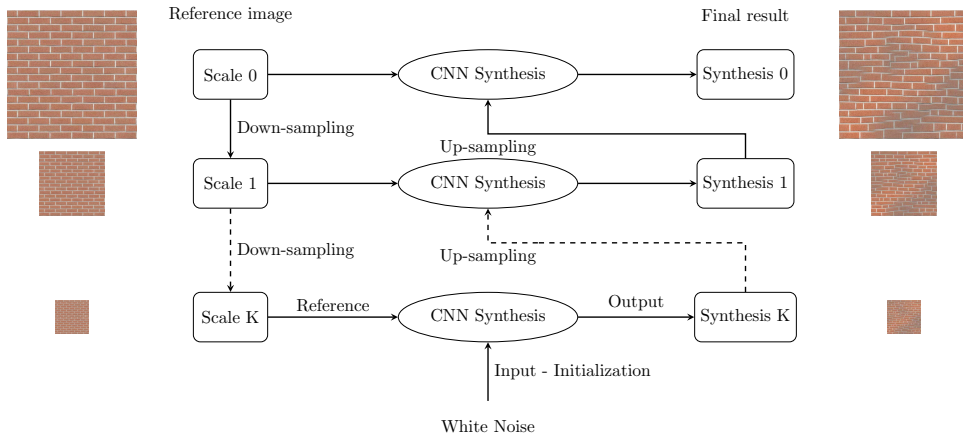


FIG. 1. *Multiscale strategy. The examplar is down-sampled by factors $2^{-1}, 2^{-2}, \cdots, 2^{-k}$ to build a pyramid $I^{(1)}, I^{(2)}, \cdots, I^{(k)}$. At scale $K$, a new texture is synthesized by using $I^{(k)}$ as the exemplar and the upsampled result of the synthesis at scale $K-1$ as initialization (instead of white noise). We repeat this step until we reach the top of the image pyramid.*

**3.3. Spectrum Constraint.** We propose to include in the synthesis a new constraint based on the Fourier spectrum of the image. It is known that such a constraint alone is an efficient way to reproduce the so-called *micro-textures* [11] made of uniformly distributed small details. This constraint has also been used in combination with more structured synthesis methods in [38].

Let us write $\mathcal{F}(I)$ for the Discrete Fourier Transform (DFT) of an image $I$, $\mathcal{F}^{-1}$ for the inverse DFT and $|.|$ for the complex modulus. The idea is to constrain the synthesized image $\tilde{I}$ to have a Fourier spectrum $|\mathcal{F}(\tilde{I})|$ as similar as possible to $|\mathcal{F}(I)|$, the spectrum of $I$. A simple way to do this is to first perform the multi-scale neural

5

synthesis described above, and then to replace the phases of the Fourier transform of the synthesized image with random phases, before applying the inverse Fourier transform to the result [11]. Now, this sequential strategy is not satisfactory, since the randomization of phases would destroy the effect of both the statistical constraint on the VGG features and the effect of the multi-scale strategy. Therefore, we propose to incorporate the Fourier constraint into the multi-scale synthesis process. A preliminary, mono-scale version of this idea was presented in [24].

In order to include the Fourier constraint into the loss function used for synthesizing images, we first introduce $\mathcal{E}_I$, the set of images having the same spectrum as $I$ the examplar image. In the case of color images, this is defined as

$$\mathcal{E}_I = \left\{ J \in \mathbb{R}^{h \times w \times 3} | \exists \phi \in \mathbb{R}^{h \times w} : \mathcal{F}(J) = e^{i\phi} \mathcal{F}(I) \right\}.$$

Next, we define the Fourier loss associated to the image $\tilde{I}$ as the normalized Euclidean distance between $\tilde{I}$ and $\mathcal{E}_I$,

$$(3.3) \qquad \mathcal{L}_{spe} = \frac{1}{2N} d(\tilde{I}, \mathcal{E}_I)^2 = \frac{1}{2N} \|\tilde{I} - \mathcal{P}(\tilde{I})\|^2,$$

and the total loss as

$$(3.4) \qquad \mathcal{L} = \mathcal{L}_{Gram} + \beta \mathcal{L}_{spe},$$

where $\beta$ is a weighting parameter. Since the Fourier loss is the distance to $\mathcal{E}_I$, its gradient is given by

$$\Delta_{spe} = N^{-1}(\tilde{I} - \mathcal{P}(\tilde{I})),$$

where $\mathcal{P}$ is the projection operator on $\mathcal{E}_I$. This projection is given by (see [38], Appendix A)

$$(3.5) \qquad \mathcal{P}(\tilde{I}_c) = \mathcal{F}^{-1} \left( \frac{\mathcal{F}(\tilde{I}) \cdot \mathcal{F}(I)}{|\mathcal{F}(\tilde{I}) \cdot \mathcal{F}(I)|} \cdot \mathcal{F}(I_c) \right), c \in \{r, g, b\}$$

where $\cdot$ is the scalar product in $\mathbb{C}^3$, that is

$$\mathcal{F}(\tilde{I}) \cdot \mathcal{F}(I) = \sum_{c=r,g,b} \mathcal{F}(\tilde{I}_c) \mathcal{F}(I_c)^*,$$

$I_c$, for $c = r, g, b$, being the color channels of $I$ and $a^*$ the conjugate of complex number $a$. This spectrum constraint can be seen as a regularization to the Ill-posed example based synthesis problem.

**3.4. Autocorrelation of the feature maps.** In this section, we consider an alternative way to impose long-range consistency, based on the autocorrelation of the features maps. This is motivated by the fact that the autocorrelation is a proxy of repeating patterns, such as the presence of periodic elements in the signal. As explained in section 2, this idea has been explored with different modality in [2, 27, 32].

The autocorrelation function of an image is defined as the convolution of the image with itself. Let $I \in \mathbb{R}^{h \times w}$, the autocorrelation $C(I) = \in \mathbb{R}^{h \times w}$ is defined, for

$\forall k \in \{1, \cdots, h\}$ and $\forall l \in \{1, \cdots, w\}$, as

$$(3.6) \qquad C(I)(k,l) = \frac{1}{N^2} \sum_{i=1}^{h} \sum_{j=1}^{w} I(i,j) I(\mid i+k \mid_h, \mid j+l \mid_w)$$

$$(3.7) \qquad = \frac{1}{N^2} I * I$$

$\mid \bullet \mid_h$ being the modulo operation with divisor h.

And efficient way to compute the autocorrelation is to use the Discrete Fourier Transform (DFT). According to the Wiener-Khintchin theorem we have :

$$C(I) = \mathcal{F}^{-1}(\mid \mathcal{F}(I) \mid^2).$$

Then, we define the autocorrelation constraint at the layer $l$ as $A^l \in \mathbb{R}^{h_l \times w_l \times m_l}$ the tensor of the squared modulus of the Fourier transform of the features maps, ie :

$$(3.8) \qquad A_p^l = \frac{1}{N_l^2} \mid \mathcal{F}(f_p^l) \mid^2$$

with $p \in \{1, \cdots, m_l\}$ the corresponding indexes of the feature map $p$. Using this toric representation allows one to consider all possible shifts between pixels.

This constraint is similar to the one in [32], except that it is dealt with in the Fourier domain and there is no weighting of the elements of the autocorrelation matrix.

**4. Experiments.** In this section, we perform experiments to illustrate both the multi-scale framework and the additional constraints we propose for neural texture synthesis. After briefly introducing the methods we compare ourselves to, we first show some visual results. Then, we propose a method to evaluate the innovation capacity of algorithms, and more precisely their tendency to produce verbatim copy of the input. Further, we evaluate the methods quantitatively using the Kullback-Leibler divergence between wavelet statistics. Despite the interest of such quantitative evaluations, it is known that they have severe limitations, in particular to evaluate results at large scales [8]. Therefore, we also have conducted a medium scale perceptual evaluation from human observers, the results of which we analyze in Section 4.3.3. These different evaluations have been conducted on the 20 texture images visible in Figure 7. These high resolution ($1024 \times 1024$) textures have been chosen to include both structured and irregular textures. Some of them display strong long-range dependency. All results can be found in Supplementary Materials. Eventually, we study the effects of various parameters and briefly illustrate the ability of our method to produce higher resolution textures.

**4.1. Architecture and parameters.** We use a VGG-19 network pre-trained on ImageNet with rescaled weights[3] as in [12] and we also use the same layers i.e. : 'Conv1_1', 'Pooling1', 'Pooling2', 'Pooling3', 'Pooling4'. The corresponding weights[4] are set to be $w_1 = w_2 = w_3 = w_4 = w_5 = 10^9$. When the spectrum constraint is considered, we use a weighting parameter $\beta = 10^5$ unless otherwise specified. Synthesis are performed using 2000 iterations. We use Tensorflow as a deep learning framework and Scipy as an optimization package. Synthesizing one texture of size $1024 \times 1024$ takes 60 minutes with a GeForce 1080 Ti for the method multi-scale "Gram + MSInit". The overhead compared to Gatys [12] at the same scale is limited because the synthesis at lower resolutions are faster.

---

[3]The rescaled VGG-19 network can be found at http://github.com/leongatys/DeepTextures
[4]Due to the numerical sensitivity of the LBFGS-B optimization algorithm.

**4.2. Other texture synthesis methods.** The first method we compare ourselves to is the original synthesis method from Gatys et al. [12], that from now we refer to as "Gatys". We also consider the method "Deep Corr", introduced in [32], using the code from the authors[5], using a maximum of 2000 iterations. We also consider the multi-scale texture algorithm from [35], using the code from the author [6], using layers 3 and 8 and 5 octaves in the Gaussian pyramid as in the original paper. We use a maximum of 2000 iterations. From now on, we refer to this method as "Snelgrove". Those last two methods have been chosen because they explicitly address the problem of reproducing large scale structures. We also consider the Feed Forward approach proposed in [39] using a PyTorch implementation by Jorge Gutierrez[7]. We refer to this method as "Ulyanov". Finally, we consider two patch-based methods, from the works [10] and [9], using implementations from the online journal IPOL [29, 1], with default parameters settings. We refer to these two methods respectively as "Efros Leung" and "Efros Freeman".

**4.3. Visual comparisons.** In Figures 2 to 5 we display synthesis results using our methods and those presented in the previous paragraph. For space reason, we only consider 4 textures, all exhibiting some kind of long-range dependency. Their resolution is $1024 \times 1024$. Some details can be seen on Figure 6. All results can be seen in Supplementary Materials.

We first notice that patch-based methods are very faithful to the reference image. However, they have the tendency to produce regions that are exact copy of the input, a verbatim effect already noticed in [1] and investigated in the next section. They also at times yield images with constant or repetitive patterns.

Among neural methods, the original "Gatys" method is still competitive, but struggles to reproduce large scales on these high resolution textures. This is due to the size of the receptive fields, which is clearly not sufficient in this case. The method from "Ulyanov" is worse in this respect. The method "Deep Corr" improves the preservation of large scale structures, but results are not satisfying, some structures are lacking and artefacts are visible. In contrast, the plain use of the auto-correlation term as an additional constraint, as we propose in "Autocorr", yields better results, even though no use of the Gram matrices is made. The regularization and innovation terms present in the method from [32] may also be harmful in these cases. Next, we observe that adding the Fourier spectrum constraint alone (at a single scale) yields interesting results, but is not enough to get fully satisfying results. The multi-scale methods, be it "Snelgrove" or the one we propose, "Gram+MSInit", "Gram+Spectrum+MSInit", "Autocorr+MSInit", all improve the original synthesis method "Gatys". In the case of very regular textures, as in Figures 3 to 5, our multi-scale methods "XXX+MSInit" yields better results, as will be confirmed by the user study in subsection 4.3.3. The method "Autocorr+MSInit" sometimes yields results that are clearly better than others, especially for very structured textures, as can be seen in Figure 4 or on the last line of Figure 6. Nevertheless, it also sometimes fails as in Figure 2 and may produce artefacts on some examples. For this reason and for human resources constraints, we choose, among our methods, to only include "Gram+MSInit" and "Gram+Spectrum+MSInit" in the user study presented in Section 4.3.3.

---

[5]The code of [32] can be found on Github : https://github.com/omrysendik/DCor

[6]The code of [35] https://github.com/wxs/subjective-functions

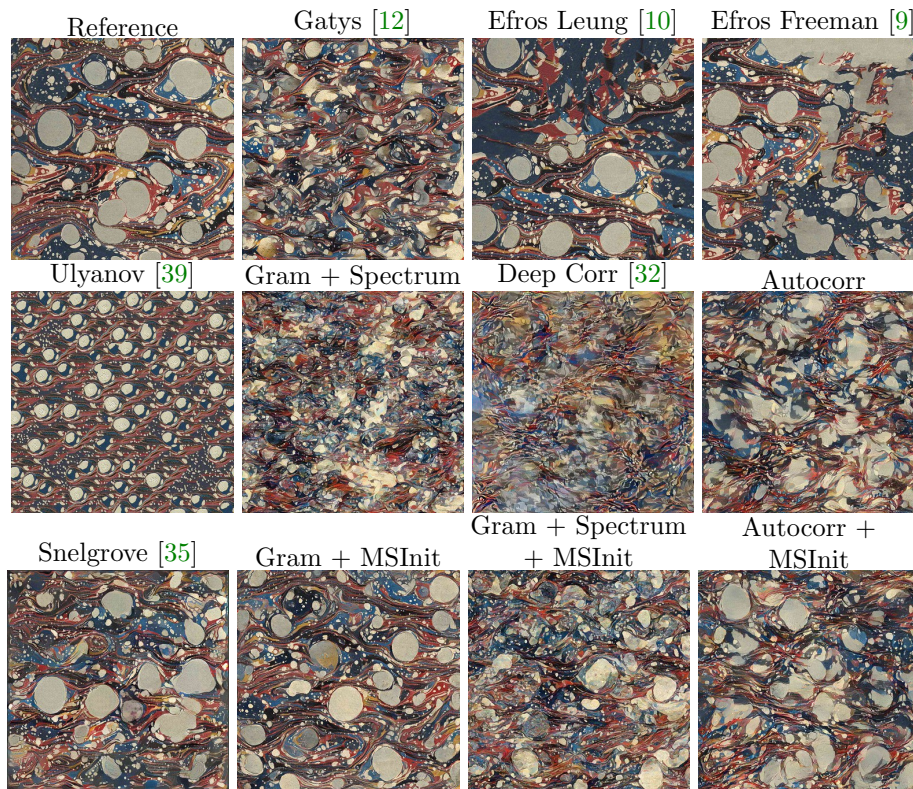[7]https://github.com/JorgeGtz/TextureNets_implementation

FIG. 2. *Synthesis results using different methods for a given reference of size* $1048 \times 1048$.

**4.3.1. Verbatim copy.** Texture synthesis methods should have the capacity to produce new images that are as diverse as possible. In the pioneering work FRAME [44], this is achieved by maximizing the entropy. Similar ideas have recently been explored in [25, 5]. Following these ideas, texture synthesis methods could be evaluated based on their capacity to maximize the entropy under some given constraints. Such a quantitative evaluation, however, is far from being trivial and probably not tractable. In this section, we take a pragmatic and much more modest way. We propose a simple way to evaluate the tendency of methods to locally produce verbatim copy of the input. This is a known default of patch-based methods, see e.g. [1, 29].

For each pixel of a given synthesis result, we look for its nearest neighbor in the input image. The notion of proximity is defined by comparing small square neighborhoods (patches) around each pixel. In Figure 8, we display the corresponding displacement map. The used color scale is obtained by assigning the $x$ coordinate of the displacement map to red, and the $y$ coordinate to blue. Verbatim copy of the input appear as constant regions in these displacement maps. As expected, the only two methods displaying large such regions are patch-based methods. All others seem to produce a reasonable amount of innovation, even though the multi-scale method from [35] can very occasionally produce small verbatim copies, probably due to the strong constraints it puts on the Gaussian pyramid.

In order to quantify the visual effect of the displacement maps, we propose to measure the flat regions corresponding to locally constant displacements. For each

9

| Reference | Gatys [12] | Efros Leung [10] | Efros Freeman [9] |
|---|---|---|---|

| Ulyanov [39] | Gram + Spectrum | Deep Corr [32] | Autocorr |
|---|---|---|---|

| Snelgrove [35] | Gram + MSInit | Gram + Spectrum + MSInit | Autocorr + MSInit |
|---|---|---|---|

FIG. 3. *Synthesis results using different methods for a given reference of size* $1048 \times 1048$.

pixel of the displacement map, we count how many of its neighbors (in 4-connexity) share its color value. Denoting $n$ this number, a score is defined as $DS = 1 - n/N$, where $N$ is the total number of investigated neighbors. The more verbatim copy there are in the synthesis, the closest the score is to 0.

The box plots of this score for the different methods and the twenty reference images can be seen in Figure 9. They confirm the impression given by the displacement map that the patch-based methods yields significantly more verbatim copy than neural methods. It should be noted however that the proposed methodology is relatively rough and does not account neither for small perturbations on the pixel positions nor for noisy pixel values.

**4.3.2. Feature-based evaluation.** Feature-based evaluation of textures is not straightforward, because no existing feature is considered as the reference one. More-over, such evaluations are inherently biased. In the most extreme case, one could even try to optimize the chosen features to synthesize new textures. In this work, we choose to rely on wavelet filters, that both are classical texture features and are not used in any of the considered methods. More precisely, we rely on the texture features proposed in [6]. In this paper, two textures are compared by computing the Kullback-Leibler divergence between parametric estimation (using generalized Gaussians) of the marginal distributions of wavelet coefficients.

In order to quantify the proximity of a synthesized texture to the reference image, we propose to :

10

FIG. 4. *Synthesis results using different methods for a given reference of size* $1048 \times 1048$.

1. Compute the wavelets coefficients of the reference image and the synthesized one (in our case we choose a Daubechies 4 wavelets as in [6] with 8 scales instead of 3, in order to account for large scale structures).
2. For each scale and orientations, estimate the parameters of a generalized Gaussian from the empirical distribution of wavelets coefficients
3. For each scale and orientations, compute the Kullback-Leibler divergence between the estimated generalized Gaussians (using a closed-form formula)

We display in Figure 10 the boxplots of the log KL scores over the 20 considered images, for the different methods. For each box, the horizontal orange line corresponds to the average result and the star to the median. On the average, the best method for this evaluation scheme appears to be "Gram+MSInit". Then follow the two patch-based methods. This is in agreement with results from the previous paragraph, since indeed a verbatim copy will have a perfect score. The next method is "Gram+Spectrum+MSInit", followed by "Snelgrove" and "Autocorr + MSInit". This evaluation confirms the good quality of results produced by the proposed "XXX+MSInit" methods, as well as "Snelgrove", at least on this image dataset containing a relatively high proportion of structured textures.

**4.3.3. Perceptual evaluation of texture synthesis methods.** Next, we further evaluate the proposed methods by performing a medium scale user study. Indeed, as shown in [8], feature-based methods such as the one of the previous section may not correlate very well with human observations, especially for long range structures.

11

FIG. 5. *Synthesis results using different methods for a given reference of size* $1048 \times 1048$.

For ethical reasons, we decided not to rely on micro-work platforms. Most users involved are volunteer PhD students or researchers, which certainly induces some bias. The total number of persons involved was 93, each having the possibility to answer up to 40 questions.

*Methodology.* Each question aims at comparing two methods on a given texture. In order to evaluate results at different scales, both the complete synthesis and a detail are presented to the user, see Figure 11. The evaluation is performed on the twenty $1024 \times 1024$ images considered in this paper. In order to get further insight on the methods, we have split the textures in two groups : regular and irregular, see Figure 7.

Following the results of the previous sections, we chose to include in the study the five following methods : "Gatys" [12], "Gram+MSInit", "Gram+Spectrum+MSInit", "Snelgrove" [35] and "Deep Coor" [32]. The first four correspond to the best feature-based score and visual impression. The last one appears to us as the most directly related to the present work in the literature, since it explicitly aims at preserving large scale structures through additional statistical constraints.

For each couple of methods (out of five) and each image, we build up two setups corresponding to the two possible respective position of methods (right and left) to avoid a possible lateral bias. This results in 400 different questions, for which we got 3170 answers.

For each question, four images are presented corresponding to the two methods at two different scales (global and local). There are 4 possible answers (method 1 is

12

the best for the global and the local scale, method 1 is the best for the local scale and method 2 the best for the global scale, etc.). Even though this is presented as a single question to the user, we treat its answers as two answers, one for the local scale and one for the global scale. This survey has been made with PsyToolkit servers [37, 36].

It should be noted that asking a question such as "which result is most similar to the reference" is not trivial. Users were indicated that by "the most similar", it should be understood "which gives the most similar visual impression". Images are not expected to correspond pixel by pixel. Ideally, a synthesized image should give the impression to correspond to a different region of the same material as the reference.

*Bradley-Terry model.* In order to quantify the results of this study, we rely on the Bradley-Terry model, as used in other perceptual study, see [41].

Let $\beta_i \in \mathbb{R}$ represent the strength of method $i$ (also called performance score), and let the outcome of a duel between methods $i$ and $j$ be determined by $\beta_i - \beta_j$. The Bradley-Terry model treats these outcomes as independent Bernoulli random variables with parameter $p_{ij}$, where the log-odds corresponding to the probability $p_{ij}$ that method $i$ beats method $j$ is modeled as :

$$(4.1) \qquad \log \frac{p_{ij}}{1 - p_{ij}} = \beta_i - \beta_j$$

Equivalently, solving for $p_{ij}$ yields

$$(4.2) \qquad p_{ij} = \frac{e^{\beta_i - \beta_j}}{1 + e^{\beta_i - \beta_j}} = \frac{e^{\beta_i}}{e^{\beta_i} + e^{\beta_j}}$$

This model is over-parameterized in the sense that it is exactly the same if we add a fixed constant to all values . The Bradley-Terry model assigns scores to a fixed set of items based on pairwise comparisons of these items, where the log-odds of item "beating" item is given by the difference of their scores. The strength is estimated by second order optimization of the maximum likelihood and the standard deviation of the difference is approximated with the Hessian of this likelihood.

*Duel results.* First, we can consider all the duels between all pairs of methods and all reference images, either from the complete set (20 images) or from the subsets of regular and irregular images separately. Results can be averaged for the global and local scale or treated separately. The results can be found on Tables 1 to 3.

Overall, the two best methods for this evaluation appear to be "Gram+MSInit" and "Gram+Spectrum+MSInit".

For the global scale, there is a draw for the complete dataset and for the irregular images, while "Gram+Spectrum+MSInit" wins on the regular images. For the local scale, "Gram+MInit" always win.

From this, we may deduce that the spectrum constraint may be useful for preserving large scale structure of regular texture, possibly at the price of a slight degradation at a more local scale. For more irregular textures, method "Gram + MSInit" should be preferred. When we consider all images and both scales (Table 1) we can extract a full ranking : "Gram+MSInit" > "Gram+Spectrum+MSInit" > Snelgrov > Gatys > Deep Cor.

*Winning probability.* An alternative evaluation consists in calculating the probability that a method $i$ is chosen among all candidates. This "winning probability" is given by the average over $j$ of the probability $p_{ij}$ that a participant chooses the

## Global case

| All images | Regular images | Irregular images |

**Table 1**

*Difference between the methods strengths $(\beta_i - \beta_j)$ (eq. (4.1)) Index $i$ corresponds to rows and index $j$ to columns. When $|\beta_i - \beta_j| > 1.96\hat{se}_{ij}$ the method $i$ is considered as beatting the method $j$ and the cell is displayed in green. In the opposite case, the cell is red. When the cell is white, the difference is not significant.*

## Local case

| All images | Regular images | Irregular images |

**Table 2**

*Difference between the methods strengths $(\beta_i - \beta_j)$ (eq. (4.1)) Index $i$ corresponds to rows and index $j$ to columns. When $|\beta_i - \beta_j| > 1.96\hat{se}_{ij}$ the method $i$ is considered as beatting the method $j$ and the cell is displayed in green. In the opposite case, the cell is red. When the cell is white, the difference is not significant.*

## Gloabl and local case

| All images | Regular images | Irregular images |

**Table 3**

*Difference between the methods strengths $(\beta_i - \beta_j)$ (eq. (4.1)) Index $i$ corresponds to rows and index $j$ to columns. When $|\beta_i - \beta_j| > 1.96\hat{se}_{ij}$ the method $i$ is considered as beatting the method $j$ and the cell is displayed in green. In the opposite case, the cell is red. When the cell is white, the difference is not significant.*

candidate $i$ over $j$:

$$(4.3) \qquad W_i = \frac{1}{N-1}\sum_{j\neq i}^{N} p_{ij} = \frac{1}{N-1}\sum_{j\neq i}^{N} \frac{e^{\beta_i - \beta_j}}{1 + e^{\beta_i - \beta_j}}$$

In contrast to the pairwise probability $p_{ij}$, $W_i$ represents the probability that a candidate $i$ was preferred over all other candidates.

We can estimate the standard error of $W_i$ as :

$$(4.4) \qquad \Sigma_i = \frac{1}{N-1}\sqrt{\sum_{j\neq i}^{N} \hat{\sigma}_{ij}^2}$$

under the hypothesis that the $p_{ij}$ are independent.

These winning probabilities are displayed on Figures 12 to 14 and confirm the duel results.

14

**4.4. Influence of parameters.** In this section, we display experiments illustrating the effects of two parameters of the proposed method : $K$, the number of considered scales, and $\beta$, the weighting of the spectrum term when using the method "Gram+Spectrum+MSInit".

**4.4.1. Multi-scale strategy.** In Figure 15, we display synthesis results with $K$ ranging from 0 (original method from[12]) to 4. The quality of results increases up to $K = 2$. This confirms the fact that the size of the filters in the VGG19 network is too small to describe large scales. Its also illustrates the fact than the VGG filters are versatile and provide good features at different scales, since the network has been trained on $224 \times 224$ input images. An interesting experiment in this respect would be to synthesize textures using the scale-invariant features from [42].

From $K = 3$, the method starts to produce results that are very similar to the reference, the case $K = 4$ being almost a copy of the reference. This may be due to the fact that in these cases, the number of parameters of the synthesis model is up to two orders of magnitude larger than the number of pixels of the coarse image. In other words, the multi-scale strategy reduces too much the solution space for this optimization problem. In practice, $K = 2$ appears a good choice for synthesizing $1024 \times 1024$ images.

**4.4.2. Weighting of the spectrum constraint.** In Figure 16, we display the result of the synthesis for different values of $\beta$, the parameter weighting the Spectrum constraint, using the method "Gram+Spectrum+MSInit". For the structured textures for which the spectrum term is useful, the best results are obtained for a relatively large $\beta$, of the order of $10^5$ for the brick image (second column). For more irregular textures, such high values may deteriorate results. This is in agreement with the results from the previous evaluations, where a value $\beta = 10^5$ was used. The problem of automatically setting this parameter is open.

**4.5. High resolution synthesis..** We conclude this experimental section by showing synthesis results of higher resolution ($1024 \times 1024$). We consider methods "Gatys", "Gram+MSInit", "Gram+Spectrum+MSInit" (both using $K = 3$). The results can be seen Figures 17 and 18. More results can be seen in Supplementary Materials.Unsurprisingly the interest of the multi-scale schemes is even stronger in this case and the mono-scale method fails. Figure 18 shows the ability of the spectrum constraint to enforce large scale regularity at this resolution.

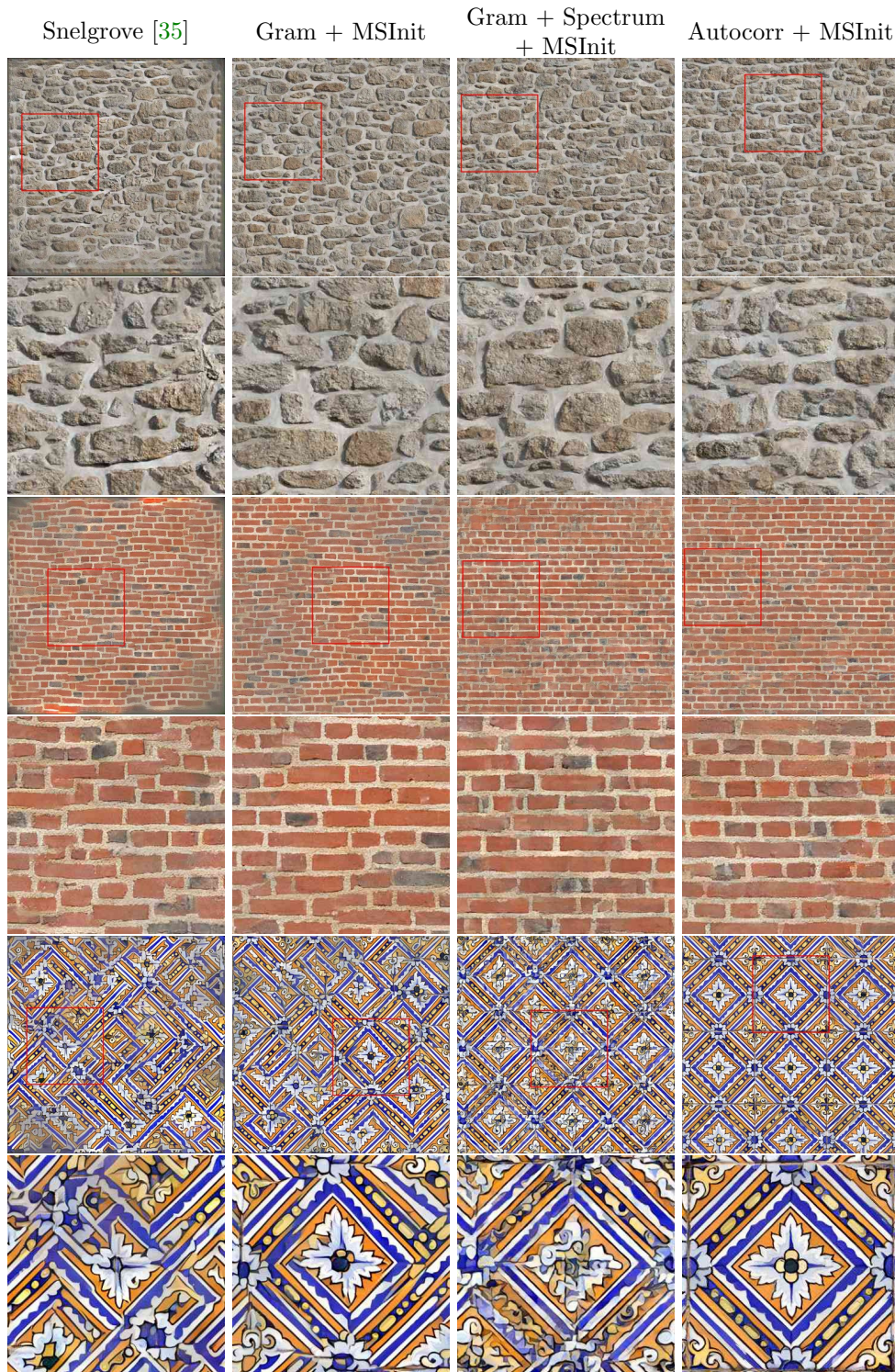| Snelgrove [35] | Gram + MSInit | Gram + Spectrum + MSInit | Autocorr + MSInit |
|---|---|---|---|



FIG. 6. *Zoom in some of the texture synthesis results. For the MSInit cases, we use $K = 2$. The region of each image framed by a red square is shown in the row below.*

16

Regular images



Irregular images



FIG. 7. *Reference images used in the different evaluation methods.*

Gatys [12]    Efros Leung [10]    Efros Freeman [9]



Ulyanov [39]    Gram + Spectrum    Deep Corr [32]    Autocorr



Snelgrove [35]    Gram + MSInit    Gram + Spectrum + MSInit    Autocorr + MSInit



FIG. 8. *Displacement map for results using different methods, for a given reference image. An area with constant color indicates a verbatm copy of the input. The synthesis can be found in 2.*

FIG. 9. *Boxplots of the displacement score for the different methods on the twenty reference images of size* $1024 \times 1024$.



FIG. 10. *Boxplots of the displacement score for the different methods on the twenty reference images of size* $1024 \times 1024$.
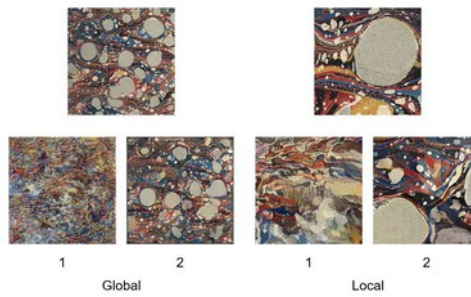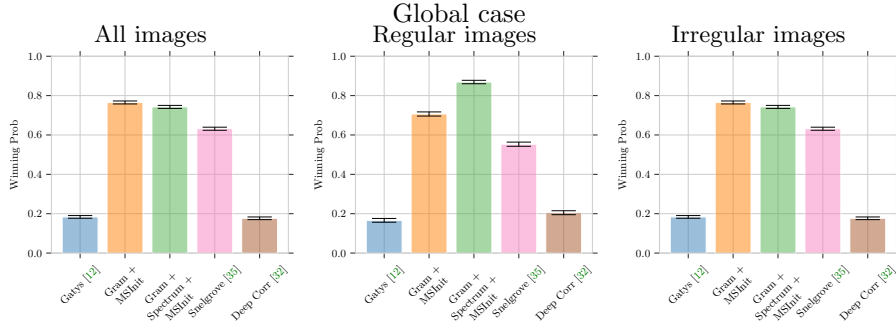


FIG. 11. *Example of the layout for one question.*

18

FIG. 12. *Winning probabilities $W_i$ with standard error $\Sigma_i$ for the different methods for the global case.*
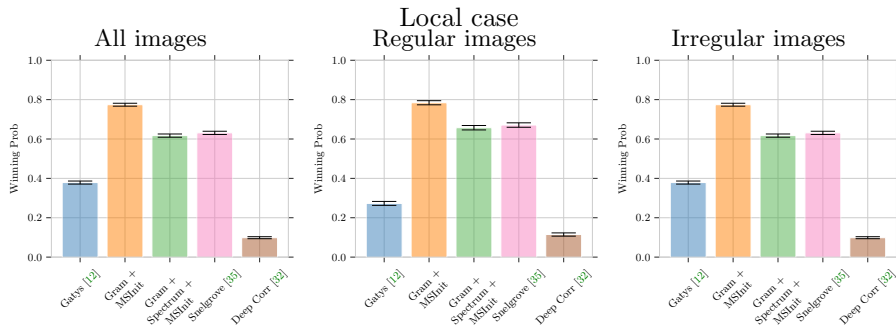


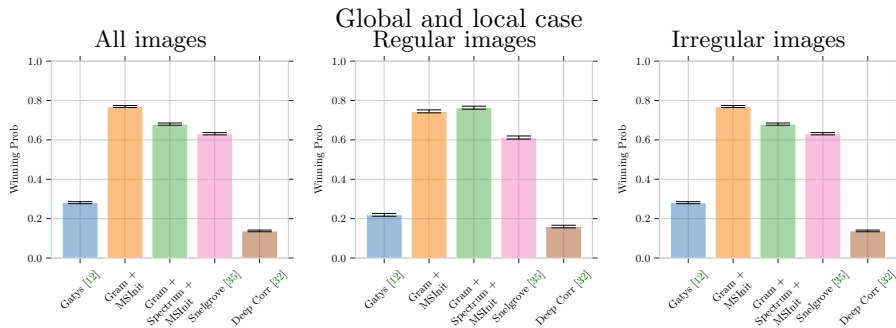FIG. 13. *Winning probabilities $W_i$ with standard error $\Sigma_i$ for the different methods for the local case.*



FIG. 14. *Winning probabilities $W_i$ with standard error $\Sigma_i$ for the different methods for both global and local cases.*
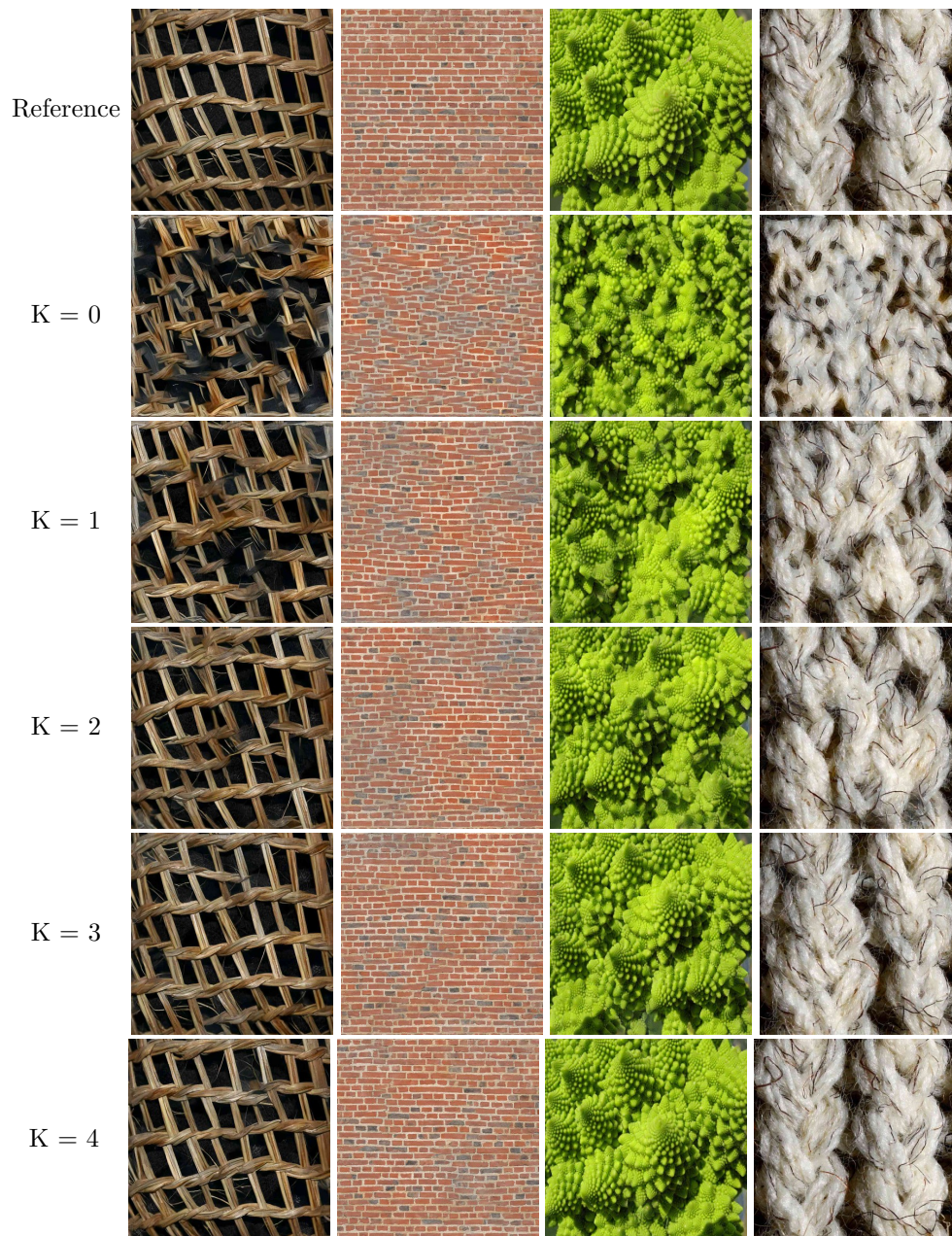
19

FIG. 15. *Synthesis results using different numbers of scales K in the multi-scale strategy. The case K = 0 corresponds to the original method from [12].*

$\beta = 0$

$\beta = 0.1$

$\beta = 10^2$

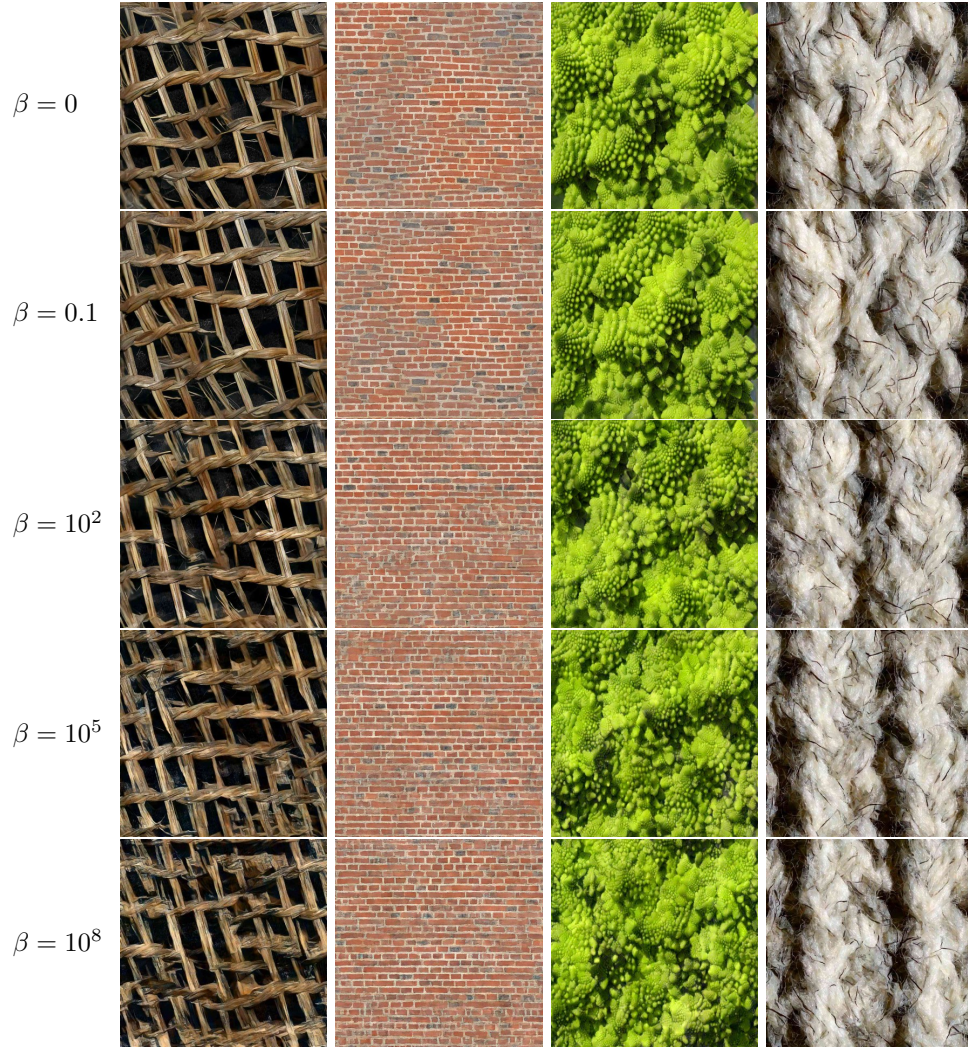$\beta = 10^5$

$\beta = 10^8$



FIG. 16. *Synthesis results using different $\beta$ in Formula* (3.4) *(original can be seen in Fig.* 15), $\beta = 0, 10^{-1}, 10^2, 10^5, 10^8$ *with the multi-scale strategy and $K = 2$.*
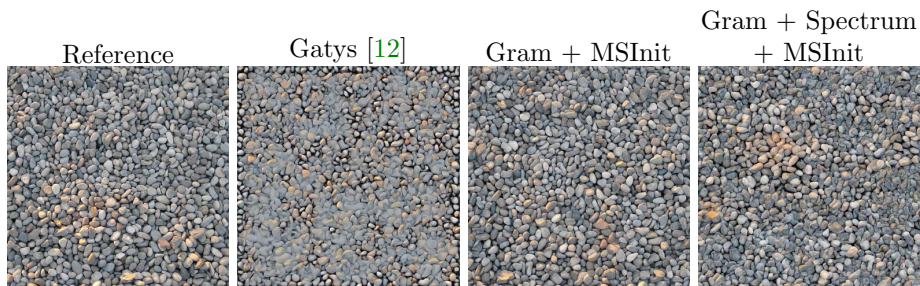
Reference     Gatys [12]     Gram + MSInit     Gram + Spectrum + MSInit



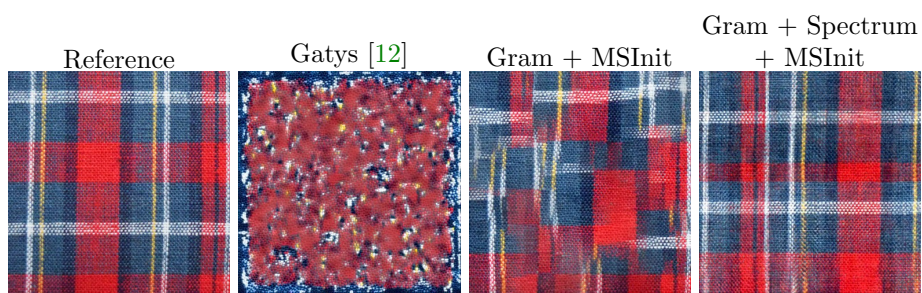FIG. 17. *Synthesis results using different methods for one given reference of size* $2048 \times 2048$.

FIG. 18. *Synthesis results using different methods for a given reference of size* $2048 \times 2048$.

**5. Conclusion.** In this paper, we have shown how a multi-resolution framework and additional statistical constraints related to long-range dependency enables one to significantly improve texture synthesis results in comparison to the seminal work [12], especially for high resolution and possibly regular textures. A natural extension would be to investigate the use of such multi-resolution strategies for style transfer for high-resolution images, following [13]. More generally, most generative methods dealing with high resolution images incorporate more or less explicitly some multi-resolution steps in their synthesis process. This is for instance the case for the very efficient StyleGan [21] approach to face synthesis. In this context, it is of great interest to investigate generic procedures to develop multi-resolution frameworks for such generative approaches.

A strong limitation of the neural methods investigated in this work is the unreasonably large number of parameters of the models. In this respect, the next question is not "what set of statistical constraint is sufficient", but "what is the minimal set of statistical constraints" needed to produce realistic synthesis. Some works [5] have shown that second order statistics between features are not necessary to get satisfying results. This, combined with the highly redundant nature of networks such as VGG19, trained for recognition, suggests that much room is available to reduce the number of parameters in these models.

## REFERENCES

[1] C. AGUERREBERE, Y. GOUSSEAU, AND G. TARTAVEL, *Exemplar-based Texture Synthesis: The Efros-Leung Algorithm*, Image Processing On Line, 3 (2013), pp. 223–241, https://doi.org/10.5201/ipol.2013.59.

[2] G. BERGER AND R. MEMISEVIC, *Incorporating long-range consistency in CNN-based texture generation*, arXiv preprint arXiv:1606.01286, (2016).

[3] A. D. CLARKE, F. HALLEY, A. J. NEWELL, L. D. GRIFFIN, AND M. J. CHANTLER, *Perceptual similarity: A texture challenge.*, in BMVC, 2011, pp. 1–10.

[4] G. R. CROSS AND A. K. JAIN, *Markov Random Field Texture Models*, IEEE Transactions on Pattern Analysis and Machine Intelligence, PAMI-5 (1983), pp. 25–39, https://doi.org/10.1109/TPAMI.1983.4767341.

[5] V. DE BORTOLI, A. DESOLNEUX, B. GALERNE, AND A. LECLAIRE, *Macrocanonical models for texture synthesis*, in International Conference on Scale Space and Variational Methods in Computer Vision, Springer, 2019, pp. 13–24.

[6] M. N. DO AND M. VETTERLI, *Wavelet-based texture retrieval using generalized Gaussian density and Kullback-Leibler distance*, IEEE Transactions on Image Processing, 11 (2002), pp. 146–158, https://doi.org/10.1109/83.982822.

[7] X. DONG AND M. J. CHANTLER, *The importance of long-range interactions to texture similarity*, in International Conference on Computer Analysis of Images and Patterns, Springer, 2013, pp. 425–432.

[8] X. DONG, J. DONG, AND M. CHANTLER, *Perceptual Texture Similarity Estimation: An Evaluation of Computational Features*, IEEE Transactions on Pattern Analysis and Machine Intelligence, (2020), pp. 1–1, https://doi.org/10.1109/TPAMI.2020.2964533.

[9] A. A. EFROS AND W. T. FREEMAN, *Image quilting for texture synthesis and transfer*, in Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques, ACM, 2001, pp. 341–346.

[10] A. A. EFROS AND T. K. LEUNG, *Texture synthesis by non-parametric sampling*, in Computer Vision, 1999. The Proceedings of the Seventh IEEE International Conference On, vol. 2, IEEE, 1999, pp. 1033–1038.

[11] B. GALERNE, Y. GOUSSEAU, AND J.-M. MOREL, *Micro-Texture Synthesis by Phase Randomization*, Image Processing On Line, 1 (2011), pp. 213–237, https://doi.org/10.5201/ipol.2011.ggm_rpn.

[12] L. GATYS, A. S. ECKER, AND M. BETHGE, *Texture Synthesis Using Convolutional Neural Networks*, (2015), pp. 262–270.

[13] L. A. GATYS, A. S. ECKER, AND M. BETHGE, *A Neural Algorithm of Artistic Style*, arXiv:1508.06576 [cs, q-bio], (2015), https://arxiv.org/abs/1508.06576.

[14] L. A. GATYS, A. S. ECKER, AND M. BETHGE, *Image style transfer using convolutional neural networks*, (2016), pp. 2414–2423.

[15] B. GIDAS, *A renormalization group approach to image processing problems*, IEEE Transactions on Pattern Analysis and Machine Intelligence, 11 (1989), pp. 164–180.

[16] K. HE, Y. WANG, AND J. HOPCROFT, *A powerful generative model using random weights for the deep image representation*, (2016), pp. 631–639.

[17] D. J. HEEGER AND J. R. BERGEN, *Pyramid-based texture analysis/synthesis*, in Proceedings of the 22nd Annual Conference on Computer Graphics and Interactive Techniques, ACM, 1995, pp. 229–238.

[18] N. JETCHEV, U. BERGMANN, AND R. VOLLGRAF, *Texture Synthesis with Spatial Generative Adversarial Networks*, arXiv:1611.08207 [cs, stat], (2016), https://arxiv.org/abs/1611.08207.

[19] J. JOHNSON, A. ALAHI, AND L. FEI-FEI, *Perceptual losses for real-time style transfer and super-resolution*, (2016), pp. 694–711.

[20] B. JULESZ, *Visual Pattern Discrimination*, IRE Transactions on Information Theory, 8 (1962), pp. 84–92, https://doi.org/10.1109/TIT.1962.1057698.

[21] T. KARRAS, S. LAINE, AND T. AILA, *A style-based generator architecture for generative adversarial networks*, in Proceedings of the IEEE conference on computer vision and pattern recognition, 2019, pp. 4401–4410.

[22] V. KWATRA, I. ESSA, A. BOBICK, AND N. KWATRA, *Texture optimization for example-based synthesis*, in ACM Transactions on Graphics (ToG), vol. 24, ACM, 2005, pp. 795–802.

[23] Y. LI, N. WANG, J. LIU, AND X. HOU, *Demystifying Neural Style Transfer*, in IJCAI, Jan. 2017, https://arxiv.org/abs/1701.01036.

[24] G. LIU, Y. GOUSSEAU, AND G.-S. XIA, *Texture synthesis through convolutional neural networks and spectrum constraints*, in 2016 23rd International Conference on Pattern Recognition (ICPR), IEEE, 2016, pp. 3234–3239.

[25] Y. LU, S.-C. ZHU, AND Y. N. WU, *Learning frame models using CNN filters*, in Thirtieth AAAI conference on artificial intelligence, 2016.

[26] B. H. MCCORMICK AND S. N. JAYARAMAMURTHY, *Time series model for texture synthesis*, International Journal of Computer & Information Sciences, 3 (1974), pp. 329–343.

[27] R. NOVAK AND Y. NIKULIN, *Improving the Neural Algorithm of Artistic Style*, arXiv:1605.04603 [cs], (2016), https://arxiv.org/abs/1605.04603.

[28] J. PORTILLA AND E. P. SIMONCELLI, *A parametric texture model based on joint statistics of complex wavelet coefficients*, International journal of computer vision, 40 (2000), pp. 49–70.

[29] L. RAAD AND B. GALERNE, *Efros and Freeman Image Quilting Algorithm for Texture Synthesis*, Image Processing On Line, 7 (2017), pp. 1–22, https://doi.org/10.5201/ipol.2017.171.

[30] E. RISSER, P. WILMOT, AND C. BARNES, *Stable and Controllable Neural Texture Synthesis and Style Transfer Using Histogram Losses*, arXiv:1701.08893 [cs], (2017), https://arxiv.org/abs/1701.08893.

[31] S. SCHREIBER, J. GELDENHUYS, AND H. DE VILLIERS, *Texture synthesis using convolutional neural networks with long-range consistency and spectral constraints*, in 2016 Pattern Recognition Association of South Africa and Robotics and Mechatronics International Conference (PRASA-RobMech), Nov. 2016, pp. 1–6, https://doi.org/10.1109/RoboMech.2016.7813173.

[32] O. SENDIK AND D. COHEN-OR, *Deep correlations for texture synthesis*, ACM Transactions on Graphics (TOG), 36 (2017), p. 161.

[33] W. SHI AND Y. QIAO, *Fast texture synthesis via pseudo optimizer*, in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 5498–5507.

[34] K. SIMONYAN AND A. ZISSERMAN, *Very Deep Convolutional Networks for Large-Scale Image Recognition*, arXiv:1409.1556 [cs], (2014), https://arxiv.org/abs/1409.1556.

[35] X. SNELGROVE, *High-resolution multi-scale neural texture synthesis*, in SIGGRAPH Asia, ACM Press, 2017, pp. 1–4, https://doi.org/10.1145/3145749.3149449.

[36] G. STOET, *PsyToolkit: A software package for programming psychological experiments using Linux*, Behavior Research Methods, 42 (2010), pp. 1096–1104, https://doi.org/10.3758/BRM.42.4.1096.

[37] G. Stoet, *PsyToolkit: A Novel Web-Based Method for Running Online Questionnaires and Reaction-Time Experiments*, Teaching of Psychology, 44 (2017), pp. 24–31, https://doi.org/10.1177/0098628316677643.

[38] G. Tartavel, Y. Gousseau, and G. Peyré, *Variational Texture Synthesis with Sparsity and Spectrum Constraints*, Journal of Mathematical Imaging and Vision, 52 (2015), pp. 124–144, https://doi.org/10.1007/s10851-014-0547-7.

[39] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. S. Lempitsky, *Texture networks: Feed-forward synthesis of textures and stylized images.*, in ICML, vol. 1, 2016, p. 4.

[40] D. Ulyanov, A. Vedaldi, and V. Lempitsky, *Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis*, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 6924–6932.

[41] K. Um, X. Hu, B. Wang, and N. Thuerey, *Spot the Difference: Accuracy of Numerical Simulations via the Human Visual System*, Journal of Computational Physics, (2019), https://arxiv.org/abs/1907.04179.

[42] N. van Noord and E. Postma, *Learning scale-variant and scale-invariant features for deep image classification*, Pattern Recognition, 61 (2017), pp. 583–592, https://doi.org/10.1016/j.patcog.2016.06.005.

[43] M.-C. Yeh and S. Tang, *Improved Style Transfer by Respecting Inter-layer Correlations*, arXiv:1801.01933 [cs], (2018), https://arxiv.org/abs/1801.01933.

[44] C. Zhu, R. H. Byrd, P. Lu, and J. Nocedal, *Algorithm 778: L-BFGS-B: Fortran Subroutines for Large-scale Bound-constrained Optimization*, ACM Trans. Math. Softw., 23 (1997), pp. 550–560, https://doi.org/10.1145/279232.279236.

[45] S. C. Zhu, Y. Wu, and D. Mumford, *Filters, Random Fields and Maximum Entropy (FRAME): Towards a Unified Theory for Texture Modeling*, International Journal of Computer Vision, 27 (1998), pp. 107–126.