



**HAL**  
open science

## Scaling Painting Style Transfer

Bruno Galerne, Lara Raad, José Lezama, Jean-Michel Morel

► **To cite this version:**

Bruno Galerne, Lara Raad, José Lezama, Jean-Michel Morel. Scaling Painting Style Transfer. 2024.  
hal-03897715v2

**HAL Id: hal-03897715**

**<https://hal.science/hal-03897715v2>**

Preprint submitted on 22 Apr 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



# Scaling Painting Style Transfer

Bruno Galerne<sup>1,2</sup>, Lara Raad<sup>3</sup>, José Lezama<sup>3†</sup>, and Jean-Michel Morel<sup>4</sup>

<sup>1</sup>Institut Denis Poisson, Université d’Orléans, Université de Tours, CNRS   <sup>2</sup>Institut Universitaire de France (IUF)  
<sup>3</sup>Instituto de Ingeniería Eléctrica, Facultad de Ingeniería, Universidad de la República   <sup>4</sup>City University of Hong Kong



**Figure 1:** Ultra-high resolution multiscale style transfer. Top row from left to right: style ( $4226 \times 5319$ ), content ( $6048 \times 8064$ ), intermediary transfer at scale 1 ( $756 \times 1008$ ) and final transfer at scale 4 ( $6048 \times 8064$ ) (intermediary results at scale 2 ( $1512 \times 2016$ ) and 3 ( $3024 \times 4032$ ) are not shown). Bottom row: zoomed in detail of size  $1024 \times 1024$  for the three UHR images and  $128 \times 128$  for the low resolution transfer at scale 1. Our method produces a style transfer of unmatched quality for such high resolution. It effectively conveys a pictorial aspect to the output images thanks to fine painting details such as brushstrokes, painting cracks, and canvas texture.

## Abstract

Neural style transfer (NST) is a deep learning technique that produces an unprecedentedly rich style transfer from a painting to an image. It is particularly impressive when it comes to transferring style from a painting to an image. NST was originally achieved by solving an optimization problem to match the global statistics of the style image while preserving the local geometric features of the content image. The two main drawbacks of this original approach is that it is computationally expensive and that the resolution of the output images is limited by high GPU memory requirements. Many solutions have been proposed to both accelerate NST and produce images with larger size. However, our investigation shows that these accelerated methods all compromise the quality of the produced images in the context of painting style transfer. Indeed, transferring the style of a painting is a complex task involving features at different scales, from the color palette and compositional style to the fine brushstrokes and texture of the canvas. This paper provides a solution to solve the original global optimization for ultra-high resolution (UHR) images, enabling multiscale NST at unprecedented image sizes. This is achieved by spatially localizing the computation of each forward and backward passes through the VGG network. Extensive qualitative and quantitative comparisons, as well as a user study, show that our method produces style transfer of unmatched quality for such high-resolution painting styles. By a careful comparison, we show that state of the art fast methods are still prone to artifacts, thus suggesting that fast painting style transfer remains an open problem.

## 1. Introduction

Style transfer is an image editing strategy transferring an image style to a content image. Given style and content, the goal is to extract the style characteristics of the style and merge them to the geometric features of the content. While this problem has a long history in computer vision and computer graphics (e.g. [HJO\*01; APH\*14]), it has seen a remarkable development since the seminal works of Gatys *et al.* [GEB15; GEB16] that introduced *neural style transfer* (NST) [JYF\*20]. These works demonstrate that the Gram matrices of the activation functions of a pre-trained VGG19 network [SZ15] faithfully encode the perceptual style and textures of an input image. NST is performed by optimizing a functional aiming at a compromise between fidelity to VGG19 features of the content image while reproducing the Gram matrices statistics of the style image. Other global statistics have been proven effective for style transfer and texture synthesis [LZW16; SC17; LPSB17; VDKC20; RWB17; HVCB21; DDD\*21; GGL22] and it has been shown that a coarse-to-fine multiscale approach allows one to reproduce different levels of style detail for images of moderate to high-resolution (HR) [GEB\*17; Sne17; GGL22]. The two major drawbacks of such optimization-based NST are the computation time and the limited resolution of images because of large GPU memory requirements. The former limitation is more critical for the present work, since conveying the visual aspects of a painting requires multiple scales of visual detail.

Regarding computation time, several methods have been proposed to generate new stylized images by training feed-forward networks [ULVL16; JAF16; LW16] or by training VGG encoder-decoder networks [CS16; HB17; LFY\*17; LLKY19; CG20]. These models tend to provide images with relatively low style transfer loss and can therefore be considered as approximate solutions to [GEB16]. Despite a remarkable acceleration, these methods suffer from GPU memory limitations due to the large size of the models used for content and style characterization and are therefore limited in terms of resolution (generally limited to  $1024^2$  pixels (px)).

This resolution limitation has received less attention but was recently tackled [ALH\*20; WLW\*20; CWX\*22; WZZ\*22]. Nevertheless, although generating UHR images (larger than 4k images), the approximate results are not able to correctly represent the style resolution. Indeed, for some methods to satisfy the GPU's memory limitations, the transfer is performed locally on small patches of the content image with a zoomed out style image ( $1024^2$  px) [CWX\*22]. In other methods, the multiscale nature of the networks is not fully exploited [WLW\*20].

At the opposite of these machine learning based-approaches, we propose to solve the original NST optimization problem [GEB16] for UHR images by introducing an exact localized algorithm. As illustrated in Figure 1, our UHR multiscale method manages to transfer the different levels of detail contained in the style image from the color palette and compositional style to the fine brushstrokes and canvas texture. The resulting UHR images look like authentic painting as can be seen in the UHR example of Figure 2.

Comparative experiments show that the results of competing methods suffer from brushstroke styles that do not match those of the UHR style image, and that very fine textures are not transferred

well and are subject to local artifacts. To straighten this visual comparison, we also introduce a qualitative and quantitative *identity test* that highlights how well a given texture is being emulated. A user study completes these experiments and confirm the superiority of our approach regarding painting style reproduction.

The main contributions of this work are summarized as follows:

- We introduce a two-step algorithm to compute the style transfer loss gradient for UHR images that do not fit in GPU memory using localized neural feature calculation.
- We show that using this algorithm in a multi-scale procedure leads to a UHR style transfer for images up to  $20k^2$  px with details conveying a natural painting aspect at *every scale*.
- Comparative experiments show that the visual quality of our UHR style transfer is by far richer and more faithful than state of the art fast but approximate solutions, revealing that, in our opinion, fast UHR painting style transfer is still an open problem.

In particular, the superiority of our approach is confirmed by a blind user study. This work provides a new reference method for high-quality style transfer with unequaled multi-resolution depth. It also naturally extends the state of the art for UHR texture synthesis. The main drawback of our approach is that it remains computationally heavy, taking several minutes to produce an image. Nevertheless, it us up to the users to define their speed VS. quality trade-off, and we believe that our algorithm can be viewed as a new gold standard for practitioners wishing to achieve the highest style transfer image quality. Our public implementation (see supplementary material (supp. mat.)) will also allow future research on fast but approximate models to be compared with our method.

## 2. Related work

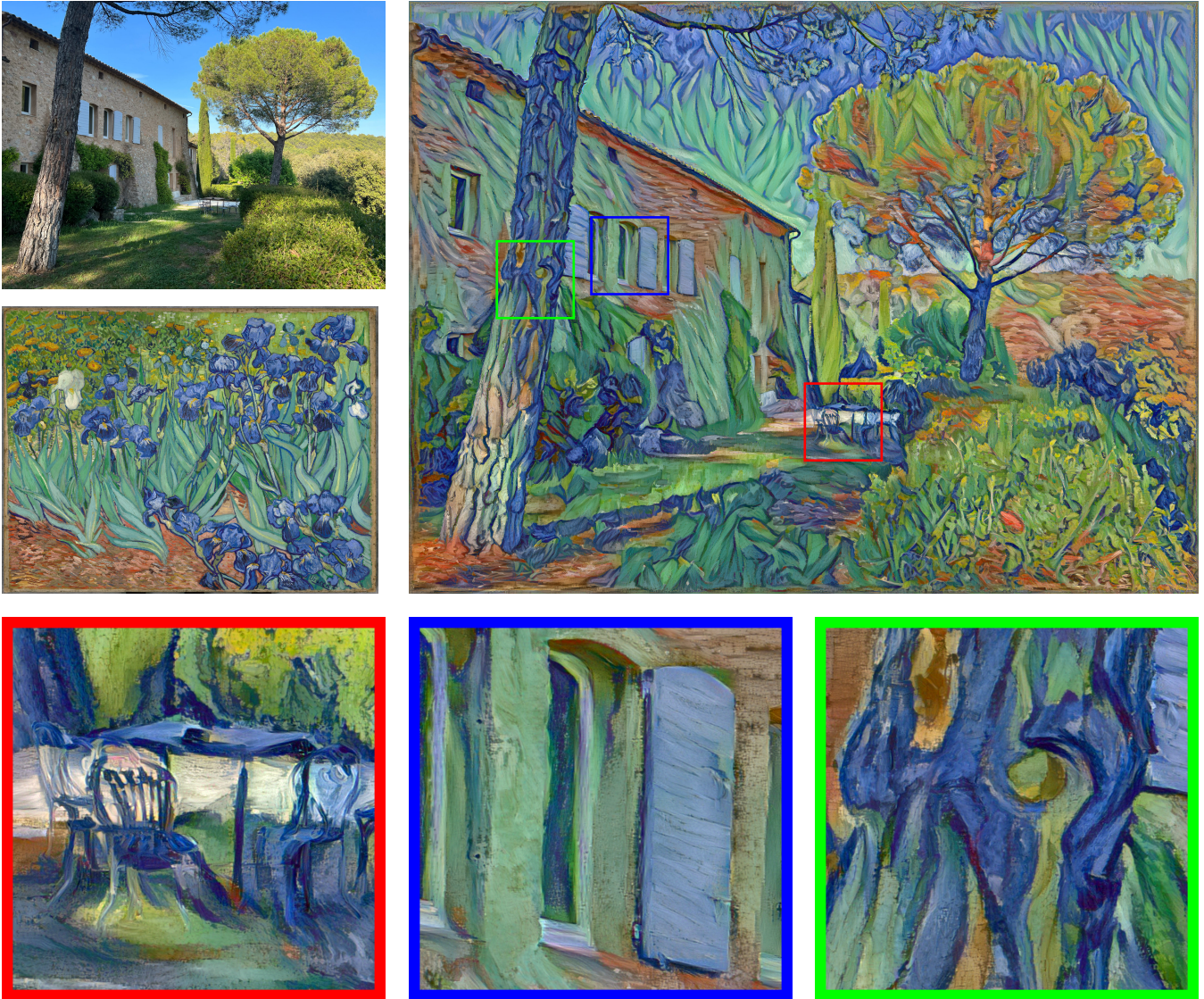
### 2.1. Style transfer by optimization

As recalled in the introduction, the seminal work of Gatys *et al.* formulated style transfer as an optimization problem minimizing the distances between Gram matrices of VGG features [GEB16]. Other global statistics have been proven effective for style transfer and texture synthesis such as deep correlations [SC17; GGL22], Bures metric [VDKC20], spatial mean of features [LZW16; DDD\*21], feature histograms [RWB17], or even the full feature distributions [HVCB21]. Specific cost function corrections have also been proposed for photorealistic style transfer [LPSB17]. When dealing with HR images, a coarse-to-fine multiscale strategy has been proven efficient to capture the different levels of details present in style images [GEB\*17; Sne17; GGL22]. Style transfer by optimization has also been extended for video style transfer [RDB16] and style transfer for neural fields [ZKB\*22]. The original optimization approach [GEB16] was considered unfitted for UHR style transfer due to high memory requirements (limited to  $1k^2$  px images, see e.g. [TFF\*20]). This paper presents an algorithm that solves this very problem for UHR images.

### 2.2. Universal style transfer (UST)

[ULVL16; UVL17] and [JAF16] showed that feed-forward networks could be trained to approximately solve style transfer. Although these models produce a very fast style transfer, they require





**Figure 2:** UHR style transfer. Top row, content image (top-left,  $6048 \times 8064$ ), style image (bottom left,  $6048 \times 7914$ ), result (right,  $6048 \times 8064$ ) (the three UHR images are downscaled  $\times 4$  for visualization). Bottom row: three zoomed in details of the result image ( $800^2$ , true resolution). Observe how very fine details such as the chairs look as if painted.

learning a new model for each style type, making them slower than the original optimization approach when training time is included.

Style limitation was addressed by training a VGG autoencoder that attempts to reverse VGG feature computations after normalizing them at the autoencoder bottleneck. Chen *et al.* [CS16] introduce the encoder-decoder framework with a style swap layer replacing content features with the closest style features on overlapping patches. Huang *et al.* [HB17] propose to use an Adaptive Instance Normalization (AdaIN) that adjusts the mean and variance of the content image features to match those of the style image. [LFY\*17] match the covariance matrices of the content image features to those of the style image by applying whitening and color-

ing transforms. These operations are performed layer by layer and involve specific reconstruction decoders at each step. [SLSW18] use one encoder-decoder block combining the transformations of [LFY\*17] and [CS16]. [PL19] introduce an attention-based transformation module to integrate the local style patterns according to the spatial distribution of the content image. [LLKY19] train a symmetric encoder-decoder image reconstruction module and a transformation learning module. [CG20] extend [LFY\*17] by embedding a new transformation that iteratively updates features in a cascade of four autoencoder modules. Despite the many improvements in fast UST strategies, we remark that: (a) they rely on matching

VGG statistics as introduced by [GEB16] (b) they are limited in resolution due to GPU memory required for large size models.

### 2.3. UST for high-resolution images

Some methods attempt to reduce the size of the network in order to perform high resolution style transfer. [ALH\*20] propose ArtNet which is a channel-wise pruned version of GoogLeNet [SLJ\*15]. [WLW\*20] propose a collaborative distillation approach in order to compress the model by transferring the knowledge of a large network (VGG19) to a smaller one, hence reducing the number of convolutional filters involved in [LFY\*17] and [HB17]. [CWX\*22] proposed an UHR style transfer framework where the content image is divided into patches and a patch-wise style transfer is performed from a zoomed out version of the style image of size  $1024^2$  px. [WZZ\*22] recently proposed to avoid using pre-trained convolutional deep neural networks for inference and instead train three very lightweight models, a content encoder, a style encoder, and a decoder, resulting in a ultra-high resolution UST with very low inference time. However, as will be shown below, the UHR style transfer results generally suffer from visual artifacts and do not faithfully convey the complexity of the style painting at all scales.

[TFF\*20] present a hybrid approach that combines neural networks and patch-based synthesis. They first perform NST between the low-resolution versions of the content and the style images, then refine the style details using patch-based transfer at a medium resolution followed by an upscaling. By design, this approach only consider a low-resolution version of the content image and suffers from a loss of details in comparison to our method (see supp. mat.).

## 3. Global optimization for neural style transfer

### 3.1. Single scale style transfer

Let us recall the algorithm of [GEB16]. It solely relies on optimizing some VGG19 second-order statistics for changing the image style while maintaining some VGG19 features to preserve the content image’s geometric features. Style is encoded through Gram matrices of several VGG19 layers, namely the set  $\mathcal{L}_s = \{\text{ReLU}_k, k \in \{1, 2, 3, 4, 5\}\}$  while the content is encoded with a single feature layer  $L_c = \text{ReLU}_{4_2}$ .

Given a content image  $u$  and a style image  $v$ , one optimizes the loss function

$$E_{\text{transfer}}(x; (u, v)) = E_{\text{content}}(x; u) + E_{\text{style}}(x; v) \quad (1)$$

where  $E_{\text{content}}(x; u) = \lambda_c \left\| V^{L_c}(x) - V^{L_c}(u) \right\|_F^2$ , with  $\lambda_c > 0$ , and

$$E_{\text{style}}(x; v) = \sum_{L \in \mathcal{L}_s} E_{\text{style}}^L(x; v). \quad (2)$$

The style loss for a layer  $L \in \mathcal{L}_s$  is the Gram loss

$$E_{\text{style}}^L(x; v) = w_L \left\| G^L(x) - G^L(v) \right\|_F^2, \quad w_L > 0, \quad (3)$$

where  $\|\cdot\|_F$  is the Frobenius norm, and, for an image  $w$  and a layer index  $L$ ,  $G^L(w)$  denotes the Gram matrix of the VGG19 features at layer  $L$ : if  $V^L(w)$  is the feature response of  $w$  at layer  $L$  that has spatial size  $n_h^L \times n_w^L$  and  $n_c^L$  channels, one first reshapes  $V^L(w)$  as a

matrix of size  $n_p^L \times n_c^L$  with  $n_p^L = n_h^L n_w^L$  the number of feature pixels, its associated Gram matrix is the  $n_c^L \times n_c^L$  matrix

$$G^L(w) = \frac{1}{n_p^L} V^L(w)^\top V^L(w) = \frac{1}{n_p^L} \sum_{k=0}^{n_p^L} V^L(w)_k (V^L(w)_k)^\top, \quad (4)$$

where  $V^L(w)_k \in \mathbb{R}^{n_c^L}$  is the column vector corresponding to the  $k$ -th line of  $V^L(w)$ .  $E_{\text{style}}^L(x; v)$  is a fourth-degree polynomial and non convex with respect to (wrt) the VGG features  $V^L(x)$ . [GEB15] propose to use the L-BFGS algorithm [Noc80] to minimize this loss, after initializing  $x$  with the content image  $u$ . L-BFGS is an iterative quasi-Newton procedure that approximates the inverse of the Hessian using a fixed size history of the gradient vectors computed during the last iterations.

### 3.2. Gram loss correction

Previous works [SC17; RWB17; HVCB21] have shown that optimizing the Gram loss alone may introduce some loss of contrast artifacts. A proposed explanation is that Gram matrices encompass information regarding both the mean values and correlation of features. While it has been shown that reproducing the full histogram of the features [RWB17; HVCB21] permits to avoid this artefact, we found that simply correcting for the mean and standard deviation (std) of each feature produced visually satisfying results and is computationally simpler.

Given some (reshaped) features  $V \in \mathbb{R}^{n_p \times n_c}$ , define  $\text{mean}(V)$  and  $\text{std}(V) \in \mathbb{R}^{n_c}$  as the spatial mean and standard deviation vectors of each feature channel. Throughout the paper, the Gram loss  $w_L \left\| G^L(u) - G^L(v) \right\|_F^2$  of Eq. (3) is replaced by the following augmented style loss

$$\begin{aligned} \tilde{E}_{\text{style}}^L(x; v) = & w_L \left\| G^L(u) - G^L(v) \right\|_F^2 \\ & + w'_L \left\| \text{mean}(V^L(x)) - \text{mean}(V^L(v)) \right\|^2 \\ & + w''_L \left\| \text{std}(V^L(x)) - \text{std}(V^L(v)) \right\|^2 \end{aligned} \quad (5)$$

for a better reproduction of the feature distribution. Note that using the “mean plus std loss” alone was proposed in [LWLH17] as an alternate loss for NST (see also [HB17]). The values of all the weights  $\lambda_c, w_L, w'_L, w''_L, L \in \mathcal{L}_s$ , have been fixed for all images (see the provided source code for the exact values).

Limiting our style loss  $\tilde{E}_{\text{style}}^L(x; v)$  to second-order statistics is capital for our localized algorithm described in Section 4. Indeed, using more involved techniques such as slice Wasserstein distance minimization [HVCB21] is not feasible for UHR images due to prohibitive memory requirement. The visual improvement when replacing  $E_{\text{style}}^L$  by  $\tilde{E}_{\text{style}}^L$  is illustrated in the supp. mat.

### 3.3. Limited resolution

Unfortunately, applying this Gatys *et al.* algorithm off-the-shelf with UHR images is not possible in practice for images of size larger than 4000 px, even with a high-end GPU. The main limitation comes from the fact that differentiating the loss  $E_{\text{transfer}}(x; (u, v))$  wrt  $x$  requires fitting into memory  $x$  and all its intermediate VGG19 features. While this requires a moderate 2.61



GB for a  $1024^2$  px image, it requires 10.2 GB for a  $2048^2$  while scaling up to  $4096^2$  is not feasible with a 40 GB GPU. In the next section we describe a practical solution to overcome this limitation.

#### 4. Localized style transfer loss gradient

As mentioned in the introduction, our main contribution is to emulate the computation of

$$\nabla_x E_{\text{transfer}}(x; (u, v)) \quad (6)$$

even for images larger than  $4000^2$  px for which evaluation and automatic differentiation of the loss is not feasible due to large memory requirements.

We first discuss how one can compute neural features in a localized way and straightforwardly compute the style transfer loss using a spatial partition of the image. Then, we demonstrate that this approach allows for the exact computation of the loss gradient using a two-pass procedure.

##### 4.1. Localized computation of neural features

First suppose one wants to compute the feature maps  $V^L(x)$ ,  $L \in \mathcal{L}_s \cup \{L_c\}$ , of an UHR image  $x$ . The natural idea developed here is to compute the feature maps piece by piece, by partitioning the input image  $x$  into small images of size  $512^2$ , that we will call blocks. This approach will work up to boundary issues. Indeed, to compute exactly the feature maps of  $x$  one needs the complete receptive field centered at the pixel of interest. Hence, each block of the partition must be extracted with a margin area, except on the sides that are actual borders for the image  $x$ . In all our experiments we use a margin of width 256 px in the image domain.

This localized way to compute features allows one to compute global feature statistics such as Gram matrices and means and stds vectors. Indeed, these statistics are all spatial averages that can be aggregated block by block by adding sequentially the contribution of each block. Hence, this easy to implement procedure allows one to compute the value of the loss  $E_{\text{transfer}}(x; (u, v))$  (see Equation (1)). However, in contrast with standard practice, it is *not* possible to automatically differentiate this loss wrt  $x$ , because the computation graph linking back to  $x$  has been lost.

##### 4.2. Localized gradient given global statistics

A close inspection of the different style losses wrt the neural features shows that they all have the same form: For each style layer  $L \in \mathcal{L}_s$ , the gradient of the layer style loss  $\tilde{E}_{\text{style}}^L(x; v)$  wrt the layer feature  $V^L(x)_k \in \mathbb{R}^{n_c^L}$  at some pixel location  $k$  only depends on the local value  $V^L(x)_k$  and on some difference between the global statistics (Gram matrix, spatial mean, std) of  $V^L(x)$  and the corresponding ones from the style layer  $V^L(v)$ .

**Proposition 4.1 (Locality of style loss gradient)** Given the layer global statistics values, the gradient of the layer style loss  $\tilde{E}_{\text{style}}^L(x; v)$  wrt the layer feature  $V^L(x) \in \mathbb{R}^{n_c^L}$  is local: The gradient value at location  $k$  only depends on the feature  $V^L(x)_k$  at the same location  $k$ .

*Proof* Recall from Equation (5) that  $\tilde{E}_{\text{style}}^L(x; v)$  is a linear combination of the Gram, mean, and std losses. As shown in Table 1, given the global statistics, each of these losses satisfies the local property.  $\square$

Exploiting this locality of the gradient, it is also possible to *exactly* compute the gradient vector  $\nabla_x E_{\text{transfer}}(x; (u, v))$  block by block using a two-pass procedure: The first pass is used to compute the global VGG19 statistics of each style layer and the second pass is used to locally backpropagate the gradient wrt the local neural features. The whole procedure is described by Algorithm 1 and illustrated by Figure 3. As illustrated by Figure 4, Algorithm 1 enables to exactly compute the global gradient of the loss in a localized way. The used block margin of size 256 is necessary to avoid visual discontinuities at block boundaries (see Figure 4).

---

**Algorithm 1** Localized computation of the style transfer loss and its gradient wrt  $x$

---

**Input:** Current image  $x$ , content image layer  $V^{L_c}(u)$ , and list of feature statistics of  $v$   $\{(G^L(v), \text{mean}(V^L(v)), \text{std}(V^L(v)))\}$ ,  $L \in \mathcal{L}_s$  (computed block by block)

**Output:**  $E_{\text{transfer}}(x; (u, v))$  and  $\nabla_x E_{\text{transfer}}(x; (u, v))$

**Step 1: Compute the global style statistics of  $x$  block by block:**  
**for** each block in the partition of  $x$  **do**

Extract the block  $b$  with margin and compute  $\text{VGG}(b)$  without computation graph

For each style layer  $L \in \mathcal{L}_s$ : Extract the features of the block by properly removing the margin and add their contribution to  $G^L(x)$  and  $\text{mean}(V^L(x))$ .

**end for**

For each style layer  $L \in \mathcal{L}_s$ : compute  $\text{std}(V^L(x))$  as a function of  $G^L(x)$  and  $\text{mean}(V^L(x))$ .

**Step 2: Compute the transfer loss and its gradient wrt  $x$  block by block:**

Initialize the loss and its gradient:  $E_{\text{transfer}}(x; (u, v)) \leftarrow \tilde{E}_{\text{style}}(x; v)$ ;  $\nabla_x E_{\text{transfer}}(x; (u, v)) \leftarrow 0$

**for** each block in the partition of  $x$  **do**

Extract the block  $b$  with margin and compute  $\text{VGG}(b)$  with computation graph

For each style layer  $L \in \mathcal{L}_s$ : Compute the gradient of the style loss wrt the local features using the global statistics of  $x$  from Step 1 and the style statistics of  $v$  as reference (Table 1)

For the content layer  $L_c$ , add the contribution of  $V^{L_c}(b)$  to the loss  $E_{\text{transfer}}(x; (u, v))$  and compute the gradient of the content loss wrt the local features (first row of Table 1)

Use automatic differentiation to backpropagate all the feature gradients to the level of the input block image  $b$ .

Populate the corresponding block of  $\nabla_x E_{\text{transfer}}(x; (u, v))$  with the inner part of the gradient obtained by backpropagation.

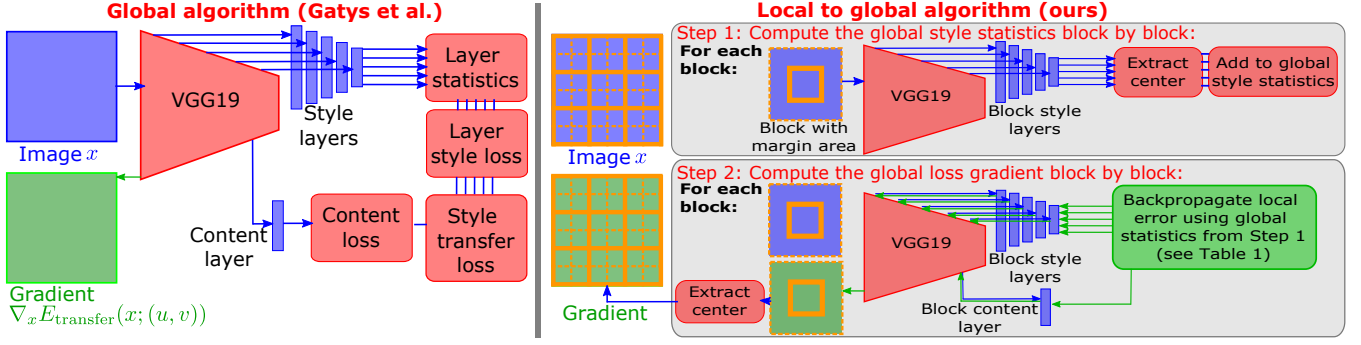
**end for**

---

## 5. Multiscale high-resolution painting style transfer

### 5.1. Coarse-to-fine style transfer

Thanks to Algorithm 1, we can apply style transfer to unprecedented scales. However, applying a direct style transfer to UHR images generally does not produce the desired effects due to the



**Figure 3:** Algorithm overview: Our localized algorithm (right part) allows to compute the global style transfer loss and its gradient wrt  $x$  for images that are too large for the original algorithm of [GEB16] (left part). See Algorithm 1 for the fully detailed procedure.

Global statistics	Feature loss Expression	Gradient wrt the feature
Raw features $V$	MSE: $E(V) = \ V - V_{\text{ref}}\ ^2$	$\nabla_V E(V) = 2(V - V_{\text{ref}})$
Gram matrix: $G = \frac{1}{n_p} VV^T$	Gram loss: $E(V) = \ G - G_{\text{ref}}\ _{\text{F}}^2$	$\nabla_V E(V) = \frac{4}{n_p} V(G - G_{\text{ref}})$
Feature mean: $\text{mean}(V)$	Mean loss: $E(V) = \ \text{mean}(V) - \mu_{\text{ref}}\ ^2$	$(\nabla_V E(V))_k = \frac{2}{n_p} (\text{mean}(V) - \mu_{\text{ref}})$
Feature std: $\text{std}(V)$	Std loss: $E(V) = \ \text{std}(V) - \sigma_{\text{ref}}\ ^2$	$\nabla_V E(V)_{k,j} = \frac{2}{n_p} (V_{k,j} - (\text{mean}(V))_j) \frac{(\text{std}(V))_j - \sigma_{\text{ref},j}}{(\text{std}(V))_j}$

**Table 1:** Expression of the feature loss gradient wrt a generic feature  $V$  having  $n_p$  pixels and  $n_c$  channels (matrix size  $n_p \times n_c$ ).

### Algorithm 2 Multiscale style transfer

**Input:** Content image  $u$ , a style image  $v$ , number of scales  $n_{\text{scales}}$

**Output:** Style transferred image  $x$

**for** scale  $s = 1$  to  $n_{\text{scales}}$  **do**

**Downscale**  $u$  and  $v$  by a factor  $2^{n_{\text{scales}} - 1}$  to obtain the low-resolution couple  $(u^\downarrow, v^\downarrow)$

**Initialization:** If  $s = 1$  let  $x = u^\downarrow$ , otherwise upscale current  $x$  by a factor 2

**Style transfer at current scale:**

$x^\downarrow \leftarrow \text{StyleTransfer}((u^\downarrow, v^\downarrow), x^\downarrow)$  using  $n_{\text{it}}^s$  iterations of L-BFGS with gradient computed with Algorithm 1

**end for**

fixed size of VGG19 receptive fields. For images larger than  $500^2$  px, visually richer results are obtained by adopting a multiscale approach [GEB\*17] corresponding to the standard coarse-to-fine texture synthesis [WL00] that we recall in Algorithm 2.

Our two step localized computation approach allows to apply style transfer through up to 6 scales (e.g. from  $512^2$  px to  $16384^2$  px). Except for the first step, all subsequent style transfers are well-initialized, allowing for a faster optimization [GEB\*17]. For our baseline implementation, we use L-BFGS with 600 iterations for the first scale and 300 iterations for the subsequent scales. Due to the large memory needed to store UHR images, the L-BFGS history is limited to the 10 last gradients for all scales except the first one that uses the standard history size of 100.

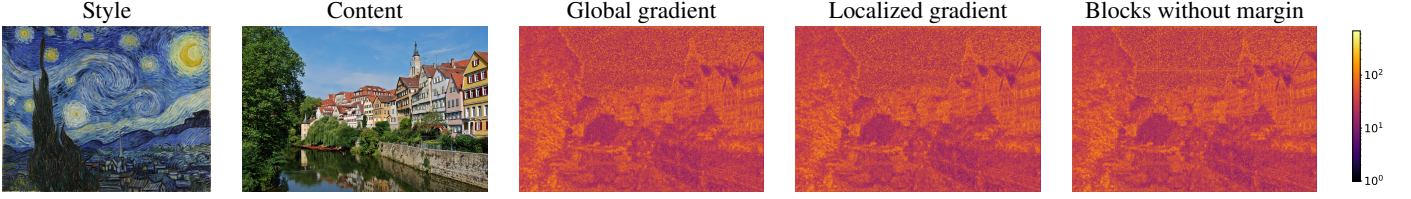
Finally, in order to avoid GPU memory saturation, for very large images we perform the L-BFGS update procedure and gradient history storing on the CPU for the last scale. This allows to increase the maximal number of pixels by 190% (+70% in square image side), as reported in the left column of Table 2. In particular this allows to apply style transfer on images with the unprecedented size of  $20k^2$  using a GPU with 80 GB of memory.

The coarse-to-fine procedure is revealed to be essential to convey the visual complexity of UHR digital photograph of a painting: the first scale encompasses color and large strokes while subsequent scales refine the stroke details up to the painting texture, bristle brushes, fine painting cracks and canvas texture, as illustrated in Figure 1. Surprisingly, fast methods for universal style transfer are not based on a coarse-to-fine approach, which is probably the main reason for their lack of fidelity to fine details (see Section 7.1).

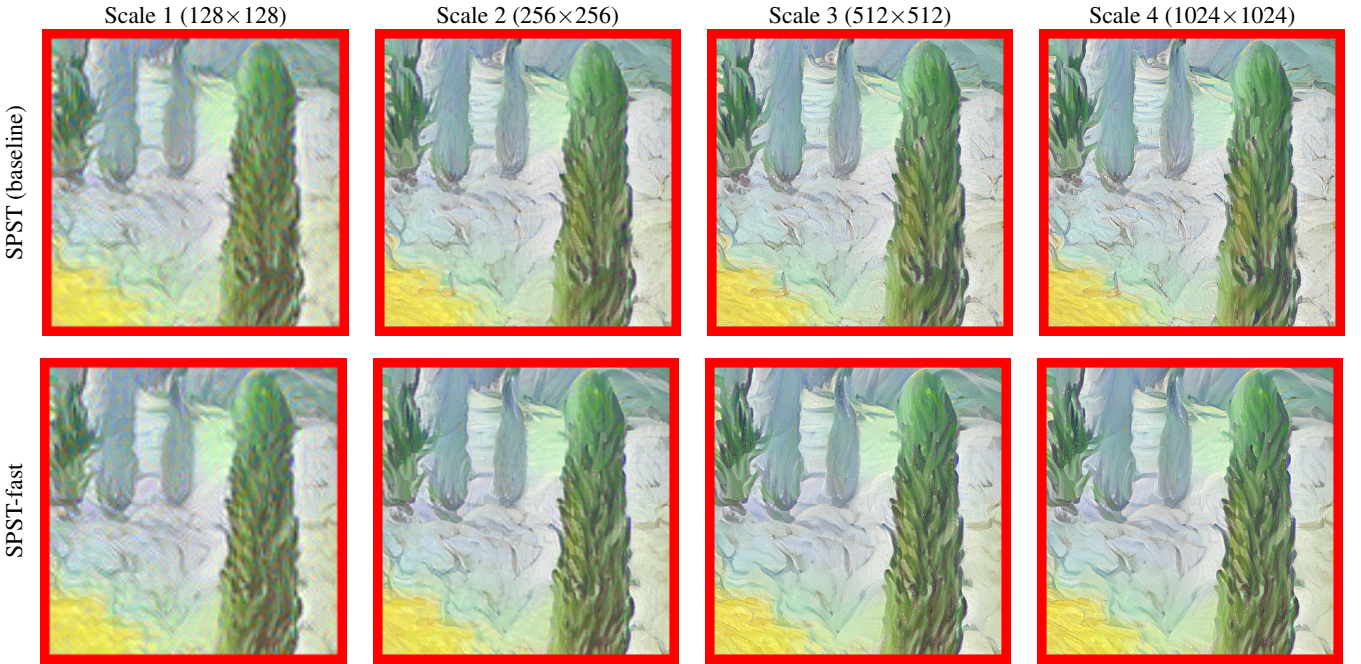
## 5.2. Accelerated multiscale style transfer

The main drawback of our baseline approach is the computational cost. Indeed, the complexity is linear in the number of pixels, making each upscaling step four times longer than the previous one. Nevertheless, we experimentally observed that the style transfer is remarkably stable from one step to the next, as can be observed in the top row of Figure 5. To the best of our knowledge, this property has never been reported, probably because style transfer involving several scales was not reachable without our localized algorithm for gradient computation.

The role of the last steps is to refine local texture in accor-



**Figure 4:** Localized gradient computation: From left to right: Style image (size  $1953 \times 2466$ ), content image (size  $1953 \times 2900$ ), norm of the RGB gradient at each pixel computed with three different approaches: reference global gradient [GEB16], localized gradient using Algorithm 1 (blocks of size  $512 \times 512$ , block margin is 256), and localized gradient with block margin set to zero. Algorithm 1 allows for the exact computation of the gradient up to numerical errors (relative error is  $1.03e-2$ ). Using a block margin of size zero instead of 256 produces a gradient with visible seams at block boundaries (relative error is  $5.23e-1$ ).



**Figure 5:** UHR multiscale style transfer: Stability of upscaling and SPST-fast. The top row shows intermediary steps for the experiment of Figure 1 displaying details of size  $128^2$  to  $1024^2$ . While the transfer is globally stable from one scale to the next, each upscaling enables the addition of fine pictorial details which give an authentic painting aspect to the final output image. Bottom row: Same details for the output of the SPST-fast alternative that uses less and less L-BFGS iterations after each upscaling. Computation times for the full scale image outputs of size  $6048 \times 8064$  are 74 minutes for SPST and 13 minutes for SPST-fast ( $\times 5.6$  speed up). Observe that both methods produce very close results but that the very fine details of the SPST-fast output are slightly less complex.

dance to the style image at the current resolution. To allow for a faster alternative, we found that these last steps can be alleviated by reducing the number of iterations. We thus propose an alternative procedure, called *SPST-fast* in what follows, that reduces the number of iterations by a factor 3 from one scale to the other, while ensuring a minimal number of 30 iterations, e.g. for 4 scales one uses  $(n_{it}^s)_{1 \leq s \leq 4} = (600, 200, 66, 30)$  instead of  $(n_{it}^s)_{1 \leq s \leq 4} = (600, 300, 300, 300)$  for the baseline implementation. Computation times for both SPST and SPST-fast are reported in Table 2 for three different GPU hardware. They show that SPST-fast

is about five times faster than SPST. Note that our algorithm allows for multiscale style transfer of UHR images up to  $20k^2$  px. Even on a moderate GPU with 11 GB of memory, our algorithm can deal with images of size  $5k^2$  px, while the original implementation of [GEB16] does not run on a 40GB GPU for an image of size  $4k^2$  px.

As shown in Figure 5, SPST-fast produces visually satisfying results but with small texture details that are slightly less aligned with the UHR content image compared to the SPST baseline approach.



GPU (VRAM)		Computation time (in min.)			
w/ max. res.	GPU / GPU+CPU	2k	4k	8k	16k
RTX 2080 (11GB)	SPST	12.8	70.3*	-	-
max. res. 3k / 5k*	SPST-fast	4.3	13.1*	-	-
A100 (40GB)	SPST	4.0	20.6	96.6	-
max. res. 8k / 14k*	SPST-fast	1.7	4.1	15.2	-
A100 (80GB)	SPST	3.8	19.0	89.7	406*
max. res. 12k / 20k*	SPST-fast	1.5	3.9	14.4	61.4*

**Table 2:** Resolution and computation time of SPST and SPST-fast depending on GPU hardware. Below each hardware name, we give the maximum resolution achievable with full computation on the GPU and the maximum resolution achievable when using the CPU for L-BFGS steps for the last scale (denoted by \*). For the computation times, \* indicates that L-BFGS optimization had to be moved to the CPU to avoid GPU memory saturation. All images are square.

## 6. Numerical Results

### 6.1. Ultra-high resolution style transfer

An example of UHR style transfer is displayed in Figure 2 with several highlighted details. Figure 1 illustrates intermediary steps of our high resolution multiscale algorithm. The result for the first scale (third column) corresponds to the ones of the original paper [GEB16] (except for our slightly modified style loss) and suffers from poor image resolution and grid artifacts. As already discussed with Figure 5, while progressing to the last scale, the texture of the painting gets refined and stroke details gain a natural aspect. This process is remarkably stable; the successive global style transfers results remain consistent with the one of the first scale.

### 6.2. Ultra-high resolution texture synthesis

Although we focus our discussion on style transfer, our approach also allows for UHR texture synthesis. Following the original paper on texture synthesis [GEB15], given a texture exemplar  $v$ , texture synthesis is performed by minimizing  $E_{\text{style}}(x; v)$  (2), starting from a random white noise image  $x_0$ . From a practical point of view, it consists in minimizing the style transfer loss with the three following differences: a) The style image is replaced by the texture image. b) There is no content image and no content loss (set  $\lambda_c = 0$ ). c) The image  $x$  is initialized as a random white noise  $x_0$ . We perform texture synthesis following the same multiscale approach and using the augmented style loss  $\tilde{E}_{\text{style}}^L(x; v)$  defined in Equation (5).

Our experiments show that for texture synthesis, one should use a number of scales as high as possible, that is, the multiscale process starts with images of moderate size (about 200 pixels). To illustrate this point we show two different UHR texture synthesis in Figures 6 (six additional results are displayed in the supp. mat.). For each example, the synthesis using three scales (same setting as for style transfer) and five scales is shown. Starting with a first scale with small size is critical for a satisfying synthesis quality. Indeed, using only three scales yields textures that are spatially homogeneous due to the white noise initialization.

Let us recall that our approach enables to reach up to  $20k^2$  px (see Table 2 for maximal resolutions), which pushes by far the maximal resolution for neural texture synthesis. Indeed, to the best of our knowledge the highest resolution reported in the neural texture synthesis literature was limited to  $2048^2$  px [GGL22] for the multiscale version of Gatys *et al.* algorithm [GEB15].

## 7. Comparison with very fast alternatives

### 7.1. Visual comparison

We compare our method with two fast alternatives for UHR style transfer, namely collaborative distillation (CD) [WLW\*20] and URST [CWX\*22] (based on [LFY\*17]) using their official implementations. To improve readability of Figure 7, results for SPST-fast, which are really close to the ones of SPST but have slightly less details, are only reproduced in supp. mat..

As already discussed in Section 2, URST decreases the resolution of the style image to  $1024^2$  px, so the style transfer is not performed at the proper scale and fine details cannot be transferred (e.g. the algorithm is not aware of the brushstroke style). As in UST methods, CD does not take into account details at different scales but simply proposes to reduce the number of filters in the auto-encoder network through collaborative distillation, to process larger images. Unsurprisingly, one observes in Figure 7 that our method is the only one capable of conveying the aspect of the painting strokes to the content image. CD suffers from halos around objects (e.g., tress in the first example), saturated color, and high-frequency artifacts (see fourth column of Figure 7). URST presents visible patch boundaries, a detail frequency mismatch due to improper scaling, loss of structure (e.g., buildings in the second example) and sometimes critical shrinking of the color palette (see fifth column of Figure 7).

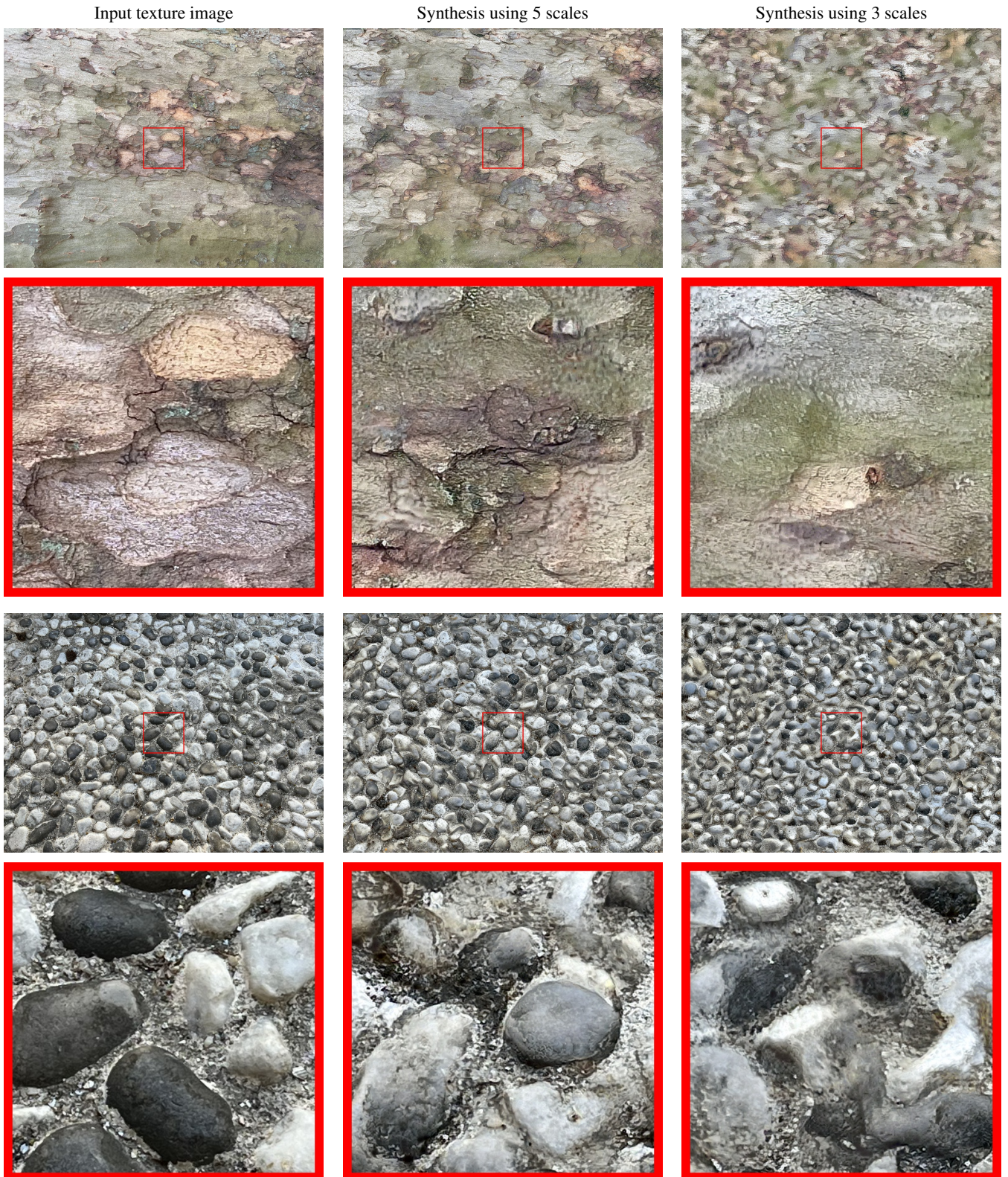
All in all, even though CD and URST produce UHR images, one can argue that the effective resolution of the output does not match their size due to the many visual artifacts. In comparison, our iterative SPST algorithm produces images for which every image part is in accordance with UHR painting style, up to the pixel level.

Finally, let us observe that the style transfer results are in general better when the geometric content of the style image and the content image are close, regardless of the method. See supp. mat. for an illustration of this limitation.

### 7.2. Identity test for style transfer quality assessment

Style transfer is an ill-posed problem by nature. We introduce here an *identity test* to evaluate if a method is able to reproduce a painting when using the same image for both content and style. Two examples of this sanity check test are shown in Figure 8. We observe that our iterative algorithm is slightly less sharp than the original style image, yet high-resolution details from the paint texture are faithfully conveyed. In comparison, the results of [WLW\*20] suffer from color deviation and frequency artifacts while the results of [CWX\*22] apply a style transfer that is too homogeneous and present color and scale issues as already discussed. Again corresponding results for Figure 8 for SPST-fast are only reproduced in supp. mat. for readability.





**Figure 6:** UHR texture synthesis (same as Figure 6): From left to right: Input texture image, synthesis using 5 scales, synthesis using 3 scales. Image have size  $(3024 \times 4032)$  (downscaled by a factor 4 for inclusion in the .pdf) and true resolution details have size  $(512 \times 512)$ .





**Figure 7:** Comparison of UHR style transfers. For each example, top row, left to right: style, content, our result (SPST), CD [WLW\*20], URST [CWX\*22]. Bottom row: zoom in of the corresponding top row. First row: content ( $3168 \times 4752$ ), style ( $2606 \times 3176$ ), SPST uses three scales; third row: content ( $3024 \times 4032$ ), style ( $3024 \times 3787$ ), SPST uses three scales; fifth row: content ( $4480 \times 5973$ ), style ( $6000 \times 4747$ ), SPST uses four scales. In comparison to our results, state of the art very fast methods produce images with many defects: halo effect, neural artifacts, blending, unfaithful color palette, ... This result in images that do not look like painting contrary to SPST outputs.

Method	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Gram $\downarrow$	Time $\downarrow$
SPST	<b>24.6</b>	<b>0.454</b>	<b>0.352</b>	<b>1.99e5</b>	25.1
SPST-fast	<b>24.6</b>	<u>0.438</u>	<u>0.446</u>	<u>4.08e5</u>	4.96
CD	21.8	0.413	0.500	4.28e7	<u>0.373</u>
URST	19.0	0.413	0.546	6.77e7	<b>0.232</b>

**Table 3:** Quantitative evaluation of identity test for UHR style transfer. Results include PSNR, SSIM [WBSS04], LPIPS [ZIE\*18], the Gram (style distance) metrics, and computation time in minutes for our results (SPST), its faster alternative (SPST-fast), CD [WLW\*20] and URST [CWX\*22]. All metrics are averages using 79 HR paintings images used as both content and style. Best results are in bold, second best underlined. Our iterative procedures SPST and SPST-fast are the best for all the image fidelity metrics but are respectively 100 and 20 times slower.

Some previous works introduced a *style distance* [WLW\*20] that corresponds to the Gram loss for some VGG19 layers, showing again that fast approximate methods try to reproduce the algorithm of Gatys *et al.* which we extend to UHR images. Since we explicitly minimize this quantity, it is not fair to only consider this criterion for a quantitative evaluation. For this reason, we also calculated PSNR, SSIM [WBSS04] and LPIPS [ZIE\*18] metrics on a set of 79 paint styles (see supp. mat.) to quantitatively evaluate our results. We further report the ‘‘Gram’’ metric, that is, the style loss of Equation (2) using the original Gram loss of Equation (3), computed on UHR results using our localized approach. The average scores reported in Table 3 confirm the good qualitative behavior discussed earlier: SPST and SPST-fast are by far the best for all the scores. However, SPST and SPST-fast are respectively 100 and 20 times slower than the fastest method.

### 7.3. User study

To further compare our results, we performed a user study comparing the fast version of our algorithm (SPST-fast) to CD [WLW\*20] and URST [CWX\*22].

The user study consisted of several evaluation instances, each of which compared four images: the style used for the transfer and the results of the three methods (SPSt-fast, URST, and CD), which were displayed at random positions for each evaluation instance. Each participant was asked to select the result closest to the style of the style image among the three displayed results.

Participants were presented three types of experiments, each of which had five instances to evaluate, thus yielding a total of 15 instances to evaluate per candidate. The results were saved only if the participant conducted the whole test. The first two experiments aim to compare the results of our identity test. In one case, the overall performance of the methods is evaluated by displaying the complete results at a resolution of  $1280 \times 720$ , and in the other case, the performance of the methods on fine details is evaluated by displaying a close-up of the results at a size of  $512 \times 512$ . For the identity test, 79 painting styles were used and each participant was shown five random instances for the global evaluation and another five for the detail evaluation. The third experiment aims to compare the re-

	Voting results (%)		
	Id global	Id detail	Style transfer
CD	<u>6.56</u>	<u>22.95</u>	<u>4.92</u>
URST	0.33	2.29	4.26
SPST-fast	<b>93.11</b>	<b>74.75</b>	<b>90.82</b>

**Table 4:** This user study results shows the percentage of times each method was selected out of the 305 comparisons for each experiment. Best results are in bold, second best underlined.

sults of the three methods when transferring a painting style image to a generic content image. Only the overall performance of the methods is compared displaying the whole results at a resolution of  $1280 \times 720$ . In this case, 13 pairs of style/content images were used, and five instances were randomly shown to each participant.

A total of 61 participants took the test, yielding a total of 305 evaluations for each type of experiment. All invited participants were image processing experts in academy and industry. The results of the study are shown in Table 4. They confirm that our approach, both for the identity test (global and close-up) and the transfer of a painting style to any image, is by far superior to CD and URST in terms of visual quality: Our method is considered better more than 90% of the times as the one that better reproduces the style of the painting (Table 4 third column).

## 8. Discussion

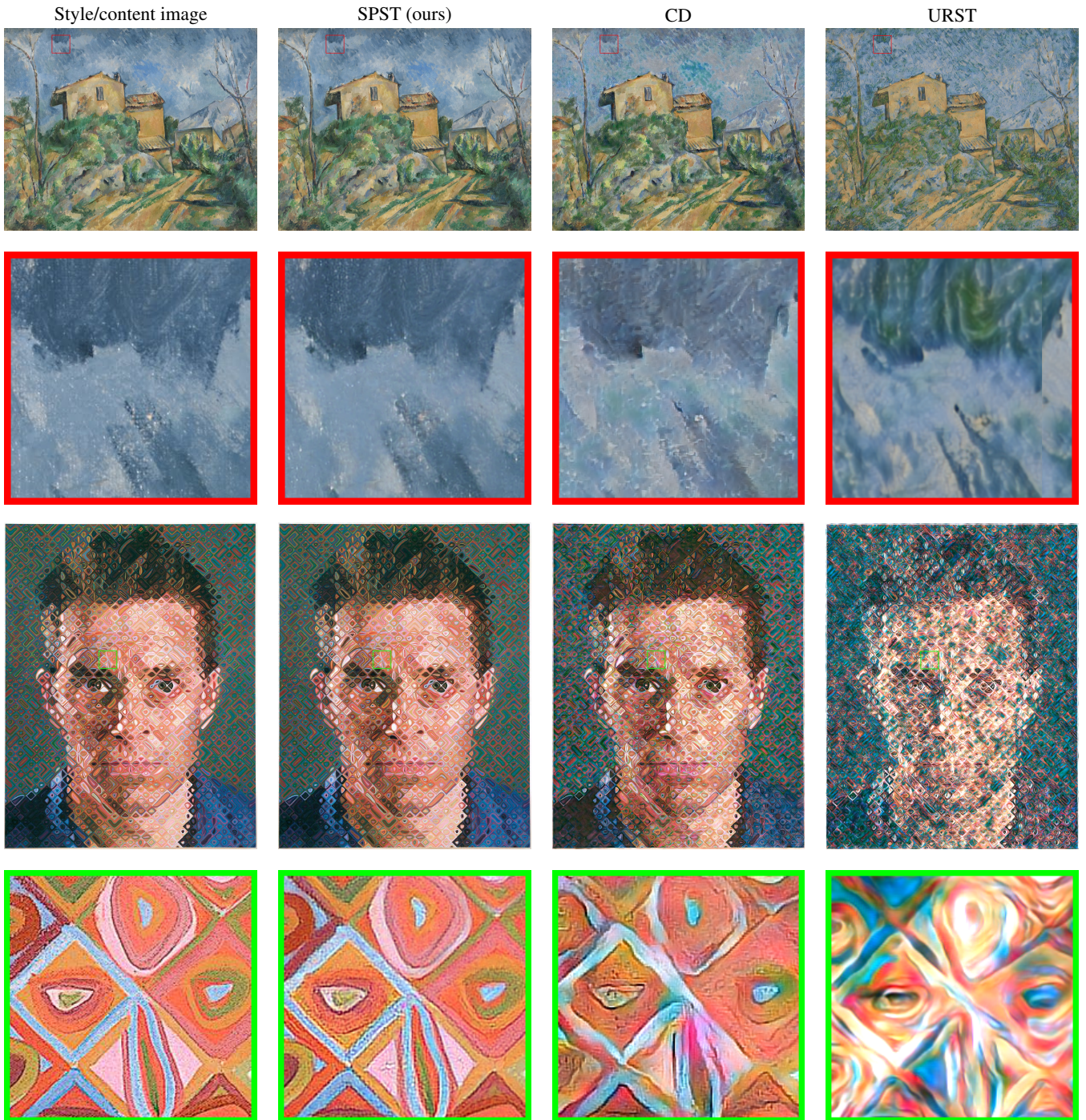
This work presented the SPST algorithm, a provably correct extension of the Gatys *et al.* style transfer algorithm to UHR images. Regarding visual quality, our algorithm outperforms competing UHR methods by conveying a true painting feel thanks to faithful HR details such as strokes, paint cracks, and canvas texture. This is clearly supported by our user study and our proposed quantitative *identity test*. SPST also allows for the synthesis by example of high-quality UHR textures. While the baseline SPST method can become prohibitively slow, even though its complexity scales linearly with image size, we proposed a faster alternative SPST-fast that limits computations as the scale grows by exploiting the stability of multiscale style transfer.

As we have demonstrated, very fast methods do not reach a satisfying quality. They fail our proposed identity test due to the presence of many artifacts, and our results are considered more faithful to the style image by a vast majority of users. This work also leads to conclude that very fast high-quality style transfer remains an open problem and that our results provide a new standard to assess the overall quality of such algorithms.

This work opens the way for several future research directions, from allowing local control for UHR style transfer [GEB\*17] to training fast CNN-based models to reproduce our results. Another promising direction is to extend our framework to video or radiance fields style transfer for which reaching ultra-high resolution would be beneficial.

**Acknowledgements:** B. Galerne and L. Raad acknowledge the support of the project MISTIC (ANR-19-CE40-005).





**Figure 8:** Identity test: a style image is transferred to itself. We compare three style transfer strategies. From left to right: ground truth style, our result (SPST), CD [WLW\*20], URST [CWX\*22]. First row: The style image has resolution  $3375 \times 4201$ ; Third row: The style image has resolution  $3095 \times 4000$  (UHR images have been downscaled by  $\times 4$  factor to save memory). Second and fourth row: Close-up view with true resolution. Observe that our results are the more faithful to the input painting and do not suffer from color blending.



## References

- [ALH\*20] AN, JIE, LI, TAO, HUANG, HAOZHI, et al. “Real-time universal style transfer on high-resolution images via zero-channel pruning”. *arXiv preprint arXiv:2006.09029* (2020) **2**, **4**.
- [APH\*14] AUBRY, MATHIEU, PARIS, SYLVAIN, HASINOFF, SAMUEL W., et al. “Fast Local Laplacian Filters: Theory and Applications”. *ACM Trans. Graph.* 33.5 (Sept. 2014). ISSN: 0730-0301. DOI: [10.1145/2629645](https://doi.org/10.1145/2629645). URL: <https://doi.org/10.1145/2629645>.
- [CG20] CHIU, TAI-YIN and GURARI, DANNA. “Iterative Feature Transformation for Fast and Versatile Universal Style Transfer”. *European Conference on Computer Vision*. Springer. 2020, 169–184 **2**, **3**.
- [CS16] CHEN, TIAN QI and SCHMIDT, MARK. “Fast patch-based style transfer of arbitrary style”. *arXiv preprint arXiv:1612.04337* (2016) **2**, **3**.
- [CWX\*22] CHEN, ZHE, WANG, WENHAI, XIE, ENZE, et al. “Towards Ultra-Resolution Neural Style Transfer via Thumbnail Instance Normalization”. *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 36. 1. June 2022, 393–400. DOI: [10.1609/aaai.v36i1.19916](https://doi.org/10.1609/aaai.v36i1.19916). URL: <https://ojs.aaai.org/index.php/AAAI/article/view/19916> **2**, **4**, **8**, **10–12**, **15**, **20**, **21**.
- [DDD\*21] DE BORTOLI, VALENTIN, DESOLNEUX, AGNÈS, DURMUS, ALAIN, et al. “Maximum Entropy Methods for Texture Synthesis: Theory and Practice”. *SIAM Journal on Mathematics of Data Science* 3.1 (2021), 52–82. DOI: [10.1137/19M1307731](https://doi.org/10.1137/19M1307731) **2**.
- [GEB\*17] GATYS, LEON A., ECKER, ALEXANDER S., BETHGE, MATTHIAS, et al. “Controlling Perceptual Factors in Neural Style Transfer”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. July 2017 **2**, **6**, **11**, **15**.
- [GEB15] GATYS, L., ECKER, A. S., and BETHGE, M. “Texture Synthesis Using Convolutional Neural Networks”. *Advances in Neural Information Processing Systems* 28. 2015, 262–270. URL: <http://papers.nips.cc/paper/5633-texture-synthesis-using-convolutional-neural-networks.pdf> **2**, **4**, **8**.
- [GEB16] GATYS, L. A., ECKER, A. S., and BETHGE, M. “Image Style Transfer Using Convolutional Neural Networks”. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2016, 2414–2423. DOI: [10.1109/CVPR.2016.265](https://doi.org/10.1109/CVPR.2016.265) **2**, **4**, **6–8**.
- [GGL22] GONTHIER, NICOLAS, GOUSSEAU, YANN, and LADJAL, SAÏD. “High-Resolution Neural Texture Synthesis with Long-Range Constraints”. *Journal of Mathematical Imaging and Vision* 64.5 (2022), 478–492 **2**, **8**.
- [HB17] HUANG, XUN and BELONGIE, SERGE. “Arbitrary Style Transfer in Real-Time With Adaptive Instance Normalization”. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. Oct. 2017 **2–4**.
- [HJO\*01] HERTZMANN, AARON, JACOBS, CHARLES E., OLIVER, NURIA, et al. “Image Analogies”. *Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques*. SIGGRAPH '01. New York, NY, USA: Association for Computing Machinery, 2001, 327–340. ISBN: 158113374X. DOI: [10.1145/383259.383295](https://doi.org/10.1145/383259.383295). URL: <https://doi.org/10.1145/383259.383295> **2**.
- [HVCB21] HEITZ, ERIC, VANHOEY, KENNETH, CHAMBON, THOMAS, and BELCOUR, LAURENT. “A Sliced Wasserstein Loss for Neural Texture Synthesis”. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2021, 9412–9420 **2**, **4**, **15**.
- [JAF16] JOHNSON, JUSTIN, ALAHI, ALEXANDRE, and FEI-FEI, LI. “Perceptual losses for real-time style transfer and super-resolution”. *European Conference on Computer Vision*. 2016, 694–711. DOI: [10.1007/978-3-319-46475-6\\_43](https://doi.org/10.1007/978-3-319-46475-6_43) **2**.
- [JYF\*20] JING, YONGCHENG, YANG, YEZHOU, FENG, ZUNLEI, et al. “Neural Style Transfer: A Review”. *IEEE Transactions on Visualization and Computer Graphics* 26.11 (2020), 3365–3385. DOI: [10.1109/TVCG.2019.2921336](https://doi.org/10.1109/TVCG.2019.2921336) **2**.
- [LFY\*17] LI, YIJUN, FANG, CHEN, YANG, JIMEI, et al. “Universal style transfer via feature transforms”. *Advances in neural information processing systems* 30 (2017), 386–396 **2–4**, **8**.
- [LLKY19] LI, XUETING, LIU, SIFEI, KAUTZ, JAN, and YANG, MING-HSUAN. “Learning linear transformations for fast image and video style transfer”. *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2019, 3809–3817 **2**, **3**.
- [LPSB17] LUAN, FUJUN, PARIS, SYLVAIN, SHECHTMAN, ELI, and BALA, KAVITA. “Deep Photo Style Transfer”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. July 2017 **2**.
- [LW16] LI, CHUAN and WAND, MICHAEL. “Precomputed real-time texture synthesis with markovian generative adversarial networks”. *European conference on computer vision*. Springer. 2016, 702–716 **2**.
- [LWLH17] LI, YANGHAO, WANG, NAIYAN, LIU, JIAYING, and HOU, XI-AODI. “Demystifying Neural Style Transfer”. *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*. 2017, 2230–2236. DOI: [10.24963/ijcai.2017/310](https://doi.org/10.24963/ijcai.2017/310). URL: <https://doi.org/10.24963/ijcai.2017/310> **4**.
- [LZW16] LU, YANG, ZHU, SONG-CHUN, and WU, YING NIAN. “Learning FRAME Models Using CNN Filters”. *Thirtieth AAAI Conference on Artificial Intelligence*. 2016 **2**.
- [Noc80] NOCEDAL, JORGE. “Updating Quasi-Newton Matrices with Limited Storage”. *Mathematics of Computation* 35.151 (1980), 773–782. ISSN: 00255718, 10886842. DOI: [10.1090/S0025-5718-1980-0572855-7](https://doi.org/10.1090/S0025-5718-1980-0572855-7) **4**.
- [PL19] PARK, DAE YOUNG and LEE, KWANG HEE. “Arbitrary style transfer with style-attentional networks”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2019, 5880–5888 **3**.
- [RDB16] RUDER, MANUEL, DOSOVITSKIY, ALEXEY, and BROX, THOMAS. “Artistic style transfer for videos”. *Pattern Recognition: 38th German Conference, GCPR 2016, Hannover, Germany, September 12–15, 2016, Proceedings* 38. Springer. 2016, 26–36 **2**.
- [RWB17] RISSER, ERIC, WILMOT, PIERRE, and BARNES, CONNELLY. *Stable and Controllable Neural Texture Synthesis and Style Transfer Using Histogram Losses*. 2017. DOI: [10.48550/ARXIV.1701.08893](https://arxiv.org/abs/1701.08893). URL: <https://arxiv.org/abs/1701.08893> **2**, **4**, **15**.
- [SC17] SENDIK, OMRY and COHEN-OR, DANIEL. “Deep Correlations for Texture Synthesis”. *ACM Trans. Graph.* 36.5 (July 2017). ISSN: 0730-0301. DOI: [10.1145/3015461](https://doi.org/10.1145/3015461). URL: <https://doi.org/10.1145/3015461> **2**, **4**, **15**.
- [SLJ\*15] SZEGEDY, CHRISTIAN, LIU, WEI, JIA, YANGQING, et al. “Going Deeper With Convolutions”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2015 **4**.
- [SLSW18] SHENG, LU, LIN, ZIYI, SHAO, JING, and WANG, XIAOGANG. “Avatar-net: Multi-scale zero-shot style transfer by feature decoration”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2018, 8242–8250 **3**.
- [Sne17] SNELGROVE, XAVIER. “High-resolution multi-scale neural texture synthesis”. *SIGGRAPH Asia 2017 Technical Briefs*. 2017, 1–4 **2**.
- [SZ15] SIMONYAN, KAREN and ZISSERMAN, ANDREW. “Very Deep Convolutional Networks for Large-Scale Image Recognition”. *International Conference on Learning Representations*. 2015 **2**.
- [TFF\*20] TEXLER, ONDŘEJ, FUTSCHIK, DAVID, FIŠER, JAKUB, et al. “Arbitrary style transfer using neurally-guided patch-based synthesis”. *Computers & Graphics* 87 (2020), 62–71. ISSN: 0097-8493. DOI: <https://doi.org/10.1016/j.cag.2020.01.002>. URL: <https://www.sciencedirect.com/science/article/pii/S0097849320300042> **2**, **4**, **15**, **16**, **23**.
- [ULVL16] ULYANOV, D., LEBEDEV, V., VEDALDI, A., and LEMPITSKY, V. “Texture Networks: Feed-forward Synthesis of Textures and Stylized Images”. *ICML*. New York, NY, USA, 2016, 1349–1357. URL: <http://dl.acm.org/citation.cfm?id=3045390.3045533> **2**.

- [UVL17] ULYANOV, DMITRY, VEDALDI, ANDREA, and LEMPITSKY, VICTOR. “Improved Texture Networks: Maximizing Quality and Diversity in Feed-Forward Stylization and Texture Synthesis”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. July 2017 2.
- [VDKC20] VACHER, JONATHAN, DAVILA, AIDA, KOHN, ADAM, and COEN-CAGLI, RUBEN. “Texture Interpolation for Probing Visual Perception”. *Advances in Neural Information Processing Systems*. Ed. by LAROCHELLE, H., RANZATO, M., HADSELL, R., et al. Vol. 33. Curran Associates, Inc., 2020, 22146–22157. URL: <https://proceedings.neurips.cc/paper/2020/file/fba9d88164f3e2d9109ee770223212a0-Paper.pdf> 2.
- [WBSS04] WANG, ZHOU, BOVIK, A.C., SHEIKH, H.R., and SIMONCELLI, E.P. “Image quality assessment: from error visibility to structural similarity”. *IEEE Transactions on Image Processing* 13.4 (2004), 600–612. DOI: [10.1109/TIP.2003.819861](https://doi.org/10.1109/TIP.2003.819861) 11.
- [WL00] WEI, L. Y. and LEVOY, M. “Fast texture synthesis using tree-structured vector quantization”. *SIGGRAPH '00*. ACM Press/Addison-Wesley Publishing Co., 2000, 479–488. DOI: [10.1145/344779.345009](https://doi.org/10.1145/344779.345009) 6.
- [WLW\*20] WANG, HUAN, LI, YIJUN, WANG, YUEHAI, et al. “Collaborative Distillation for Ultra-Resolution Universal Style Transfer”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2020 2, 4, 8, 10–12, 15, 20, 21.
- [WZZ\*22] WANG, ZHIZHONG, ZHAO, LEI, ZUO, ZHIWEN, et al. “MicroAST: Towards Super-Fast Ultra-Resolution Arbitrary Style Transfer”. *arXiv* (2022). DOI: [10.48550/ARXIV.2211.15313](https://doi.org/10.48550/ARXIV.2211.15313). URL: <https://arxiv.org/abs/2211.15313> 2, 4.
- [ZIE\*18] ZHANG, RICHARD, ISOLA, PHILLIP, EFROS, ALEXEI A., et al. “The Unreasonable Effectiveness of Deep Features as a Perceptual Metric”. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2018 11.
- [ZKB\*22] ZHANG, KAI, KOLKIN, NICK, BI, SAI, et al. “ARF: Artistic Radiance Fields”. *ECCV 2022*. 2022 2.

# Supplementary material

## A. Implementation details and reproducibility

Our implementation is based on PyTorch and our source code is provided as supplemental material (see `scaling_painting_style_transfer.zip` that will be made a public github repository in case of acceptance). It mainly uses a class dealing with tensor spatial partitioning to compute the localized gradient block by block.

In most figures the UHR images have been downscaled and saved in jpeg format. UHR images of all figures will be made available for detailed inspection and future comparison on a dedicated project website.

As mentioned in the paper, the values of all the weights  $\lambda_c, w_L, w'_L, w''_L, L \in \mathcal{L}_s$ , have been fixed for all images, hence all figures are reproducible.

## B. Visualization of Gram loss correction

A “mean plus std” corrective term is added to the Gram loss to avoid loss of contrast artifacts and grayish color alteration that may occur when minimizing this loss, as previously documented [SC17; RWB17; HVCB21]. Figure 10 shows comparative experiment when using the original Gram loss  $E_{\text{style}}^L(x; v)$  instead of our proposed augmented style loss  $\tilde{E}_{\text{style}}^L(x; v)$  (see the main paper for equations).

## C. Additional style transfer results

### C.1. Multiscale example

Figure 9 presents a second example of ultra-high resolution (UHR) style transfer with each intermediate results of the multiscale algorithm. One can see how the texture of the painting and the stroke details are gradually refined.

### C.2. Additional comparison between SPST and SPST-fast

We introduce a faster version of SPST, called SPST-fast that uses less and less iterations as the scale grows. SPST-fast is five times faster than SPST. It produces results that are really close to our SPST baseline. The main difference is that some textural details at the highest scale are less aligned with the UHR content since we only use only a few iterations for the last scale (eg 66 or 30 iterations instead of 300). Figure 11 shows the three examples of style transfer used as comparison in the main paper. One can see that both results are visually very close. A close inspection shows some slight variations such as for the eye detail in the last example.

For the sake of completeness, the results of SPST-fast for the two examples of identity tests are shown in Figure 12. One can observe on the top part of the first close-up that some brush strokes texture are better reproduced with SPST.

### C.3. Additional comparison experiments

Figure 13 presents some additional comparison examples. These results show that the competing methods suffer from color imbalance and do not match the fine texture and strokes of the style painting. This is due to the fact that neither URST [CWX\*22] nor CD [WLW\*20] take into account details at different scales. On the other hand, we are able to convey the appearance of fine details and painting strokes to the content image. In the first example, neither CD nor URST is consistent with the size of the brushstroke visible in the style image. In the second example, CD doubles the frequency of details (e.g., branches), resulting in structural inconsistency, while URST loses the branches completely. In the third example, CD has a halo effect around the bell tower and URST has inconsistent color compared to the input style and diffused edges.

Figure 14 shows an additional examples of style transfer on a portrait. Once again, we can observe that our method is the only one able to transfer the texture and strokes of the painting to the portrait content. In this example, URST has inconsistent colors compared to the style image and neither URST nor CD is consistent with the fine scale details of the style brush strokes.

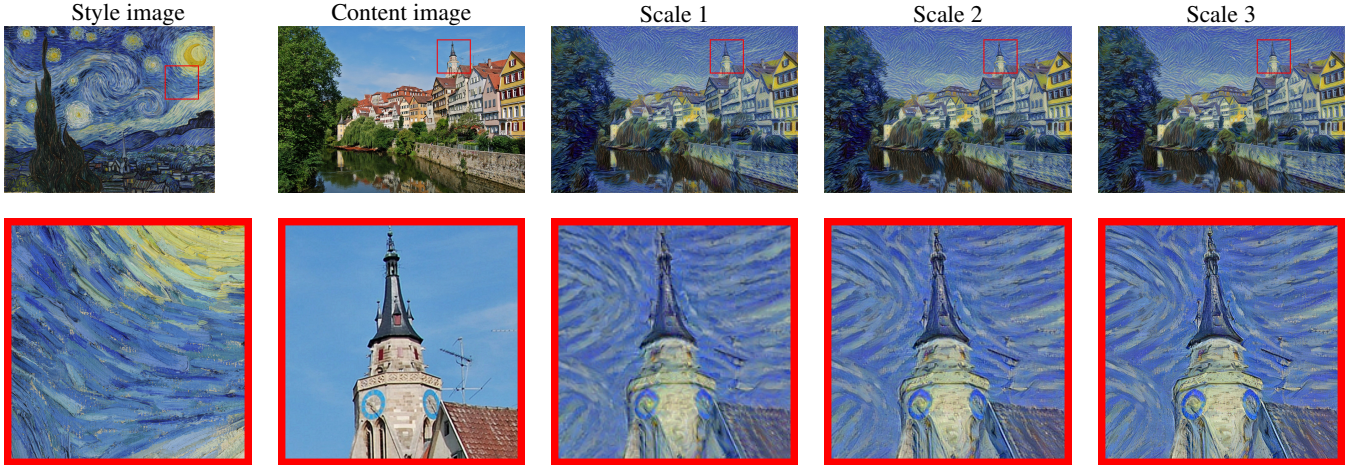
### C.4. Failure cases due to content-style mismatch

Figure 15 shows some examples of results in case of content-style mismatches that underline the importance of correctly choosing the couple content/style for a correct style transfer. The first row example uses the same style as in Figure 14 but with a different content image where the proportion of the face space is much larger than in the style image. As result, some red dots from the background are spread over the face. Note also that the lack of a beard in the content produces an undesirable effect on the face. The content image of the second row example lack of details and the style transfer creates undesirable noisy patterns in the sky and grass area. The third example is more extreme and results in loss of details from the content image and synthesis of phantom portrait silhouettes in the flat areas of the image.

As said in the conclusion of the main paper, allowing local control [GEB\*17] for UHR style transfer is an interesting direction for future developments.

### C.5. Comparison between SPST and AST Using Neurally-Guided Patch-Based Synthesis

Our paper focuses on neural style transfer methods. As said in the main paper, [TFF\*20] propose a hybrid method combining neural style transfer and patch-based transfer. Figure 16 shows a comparison between our algorithm and the one described in [TFF\*20]. It can be observed that compared to the UHR fast UST methods, the style details coherence at different scales is by far of better quality. Despite that, the content image is only used at the lowest scale for neural style transfer. This yields a final stylized image with high-resolution details unaligned with the original content image. For



**Figure 9:** UHR multiscale style transfer. Top row from left to right: style ( $3906 \times 4933$ ), content ( $3906 \times 5800$ ), transfer at scale 1 ( $976 \times 1450$ ), 2 ( $1953 \times 2900$ ), 3 ( $3906 \times 5800$ ). Bottom row: zoomed in detail for each image (from  $200^2$  to  $800^2$ ). While the transfer is globally stable from one scale to the other, each upscaling allows to add fine pictorial details which give an authentic painting aspect to the final output image. *We recommend a screen examination of the images after  $\times 8$  zoom in.*

instance, in the example of figure 16, the chairs present in the content image are no longer present in the stylized result of [TFF\*20] (third column) contrary to our result (second column).

### C.6. Full resolution images

To limit the main paper file size full resolution UHR images were not included. All UHR images have been downscaled by a factor  $\times 4$ , with highlighted details included with the true resolution.

In the following figures Figure 17, 18, 19, and 20, we include four experiments with the style and content images with a better resolution (only downgraded by a factor  $\times 2$ ) and our style transfer result in full resolution. Note that images have been compressed using jpeg quality 85 to limit the size of this supplementary material.

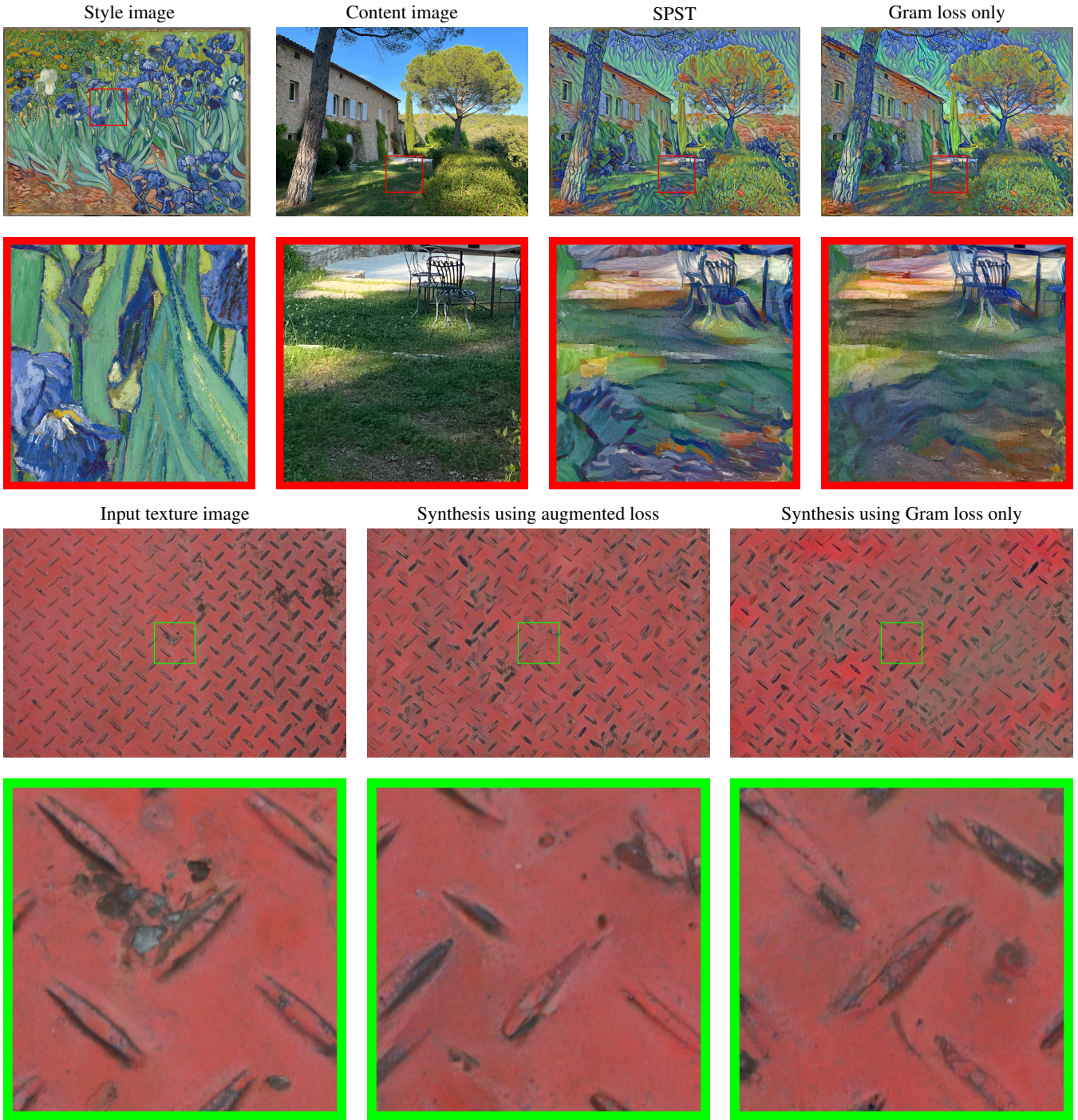
### D. Painting images dataset for the *identity test* experiment

Figure 21 shows the 79 UHR painting images used for the *identity test* where we evaluate whether a method was able to reproduce a painting when the content and style image were identical.

### E. Additional UHR texture synthesis results

Figures 22, 23 and 24 present additional UHR texture synthesis results.





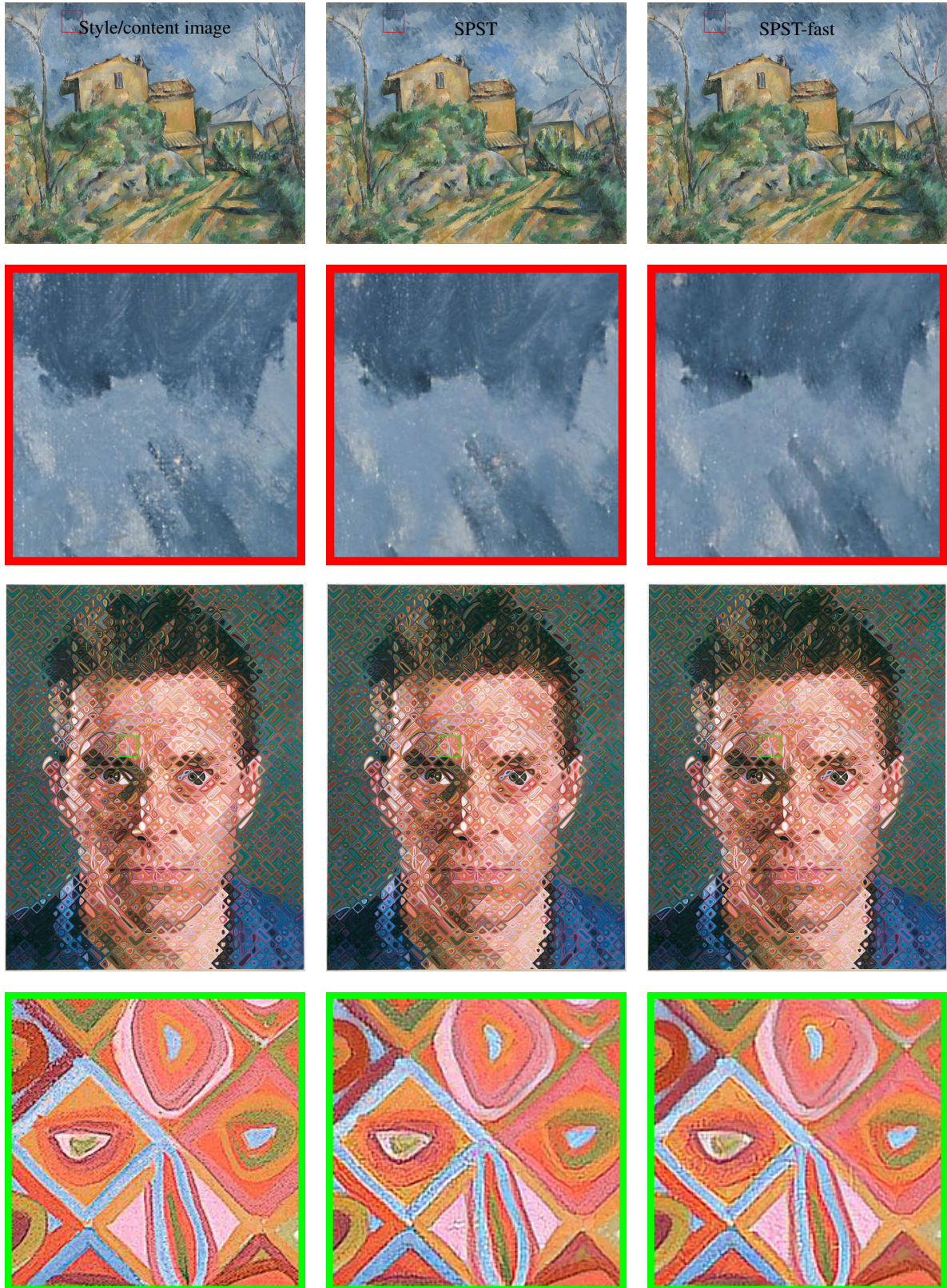
**Figure 10:** Interest of the augmented style loss for style transfer and texture synthesis: Top part: Style transfer with style of size  $6048 \times 7914$ , content and transfer results of size  $6048 \times 8064$ . Bottom part: Texture synthesis results for an image of size  $2848 \times 4272$ . Using only the original Gram loss produces gray color areas. Note that for the style transfer example this results in an image with darker colors and less brushstrokes.





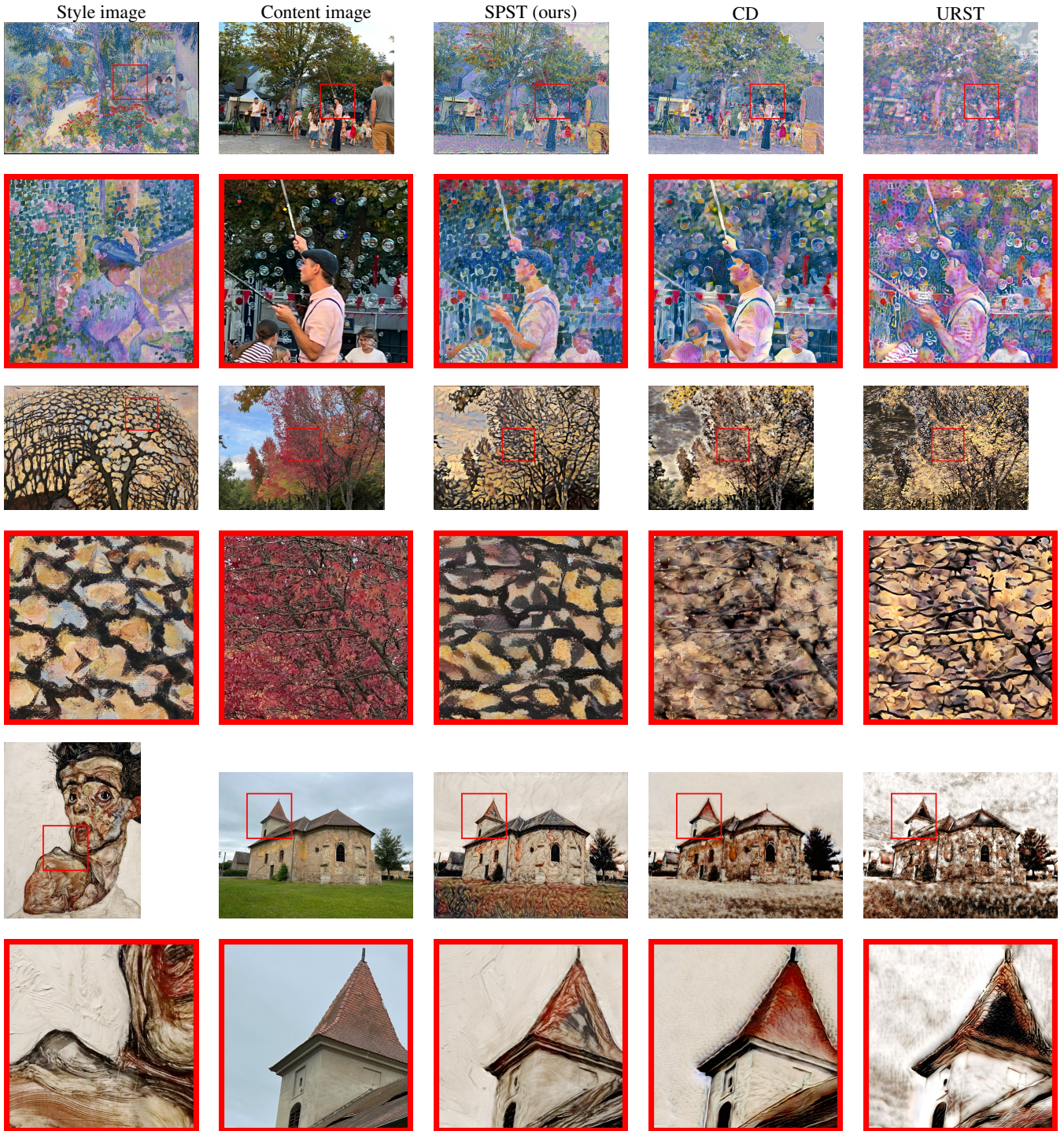
**Figure 11:** Comparison of SPST and SPST-fast style transfers (same images as the comparison figure of the main paper). First row: content ( $3168 \times 4752$ ), style ( $2606 \times 3176$ ), SPST and SPST-fast use three scales; third row: content ( $3024 \times 4032$ ), style ( $3024 \times 3787$ ), SPST and SPST-fast use three scales; fifth row: content ( $4480 \times 5973$ ), style ( $6000 \times 4747$ ), SPST and SPST-fast use four scales. Both algorithm produce visually similar style transfer. Some details of SPST are slightly better such as the eye contours in the last example.





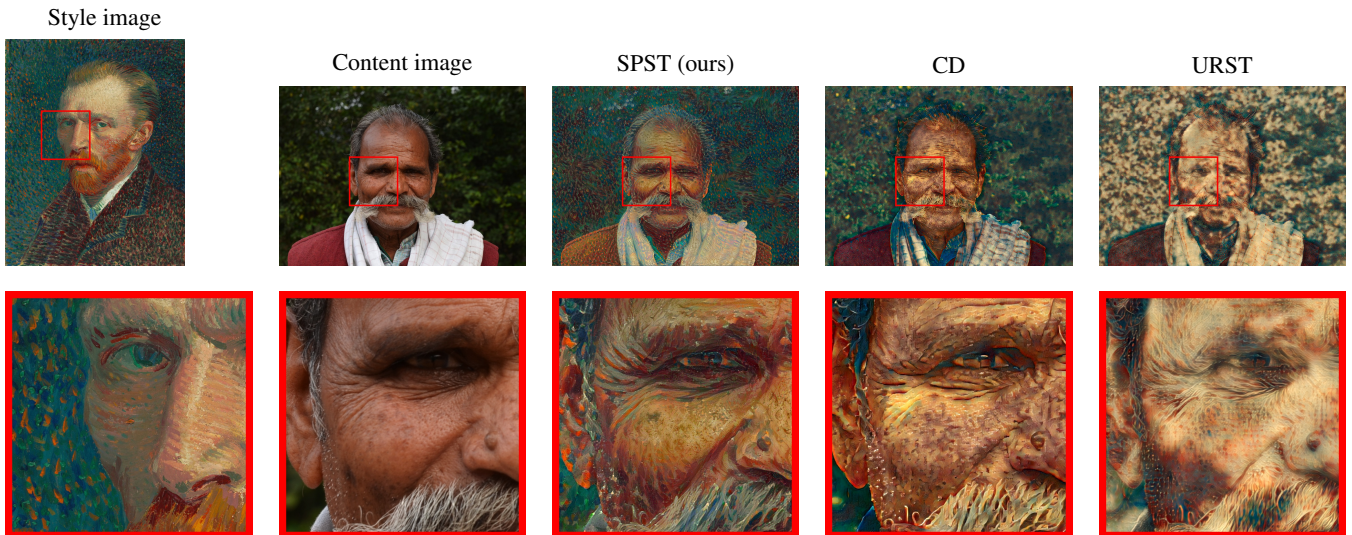
**Figure 12:** Identity test: Comparison between SPST and SPST-fast results. First row: The style image has resolution  $3375 \times 4201$ ; Third row: The style image has resolution  $3095 \times 4000$  (UHR images have been downsampled by  $\times 4$  factor to save memory). Second and fourth row: Close-up view with true resolution. SPST-fast outputs are visually close to SPST results but have slightly less details.



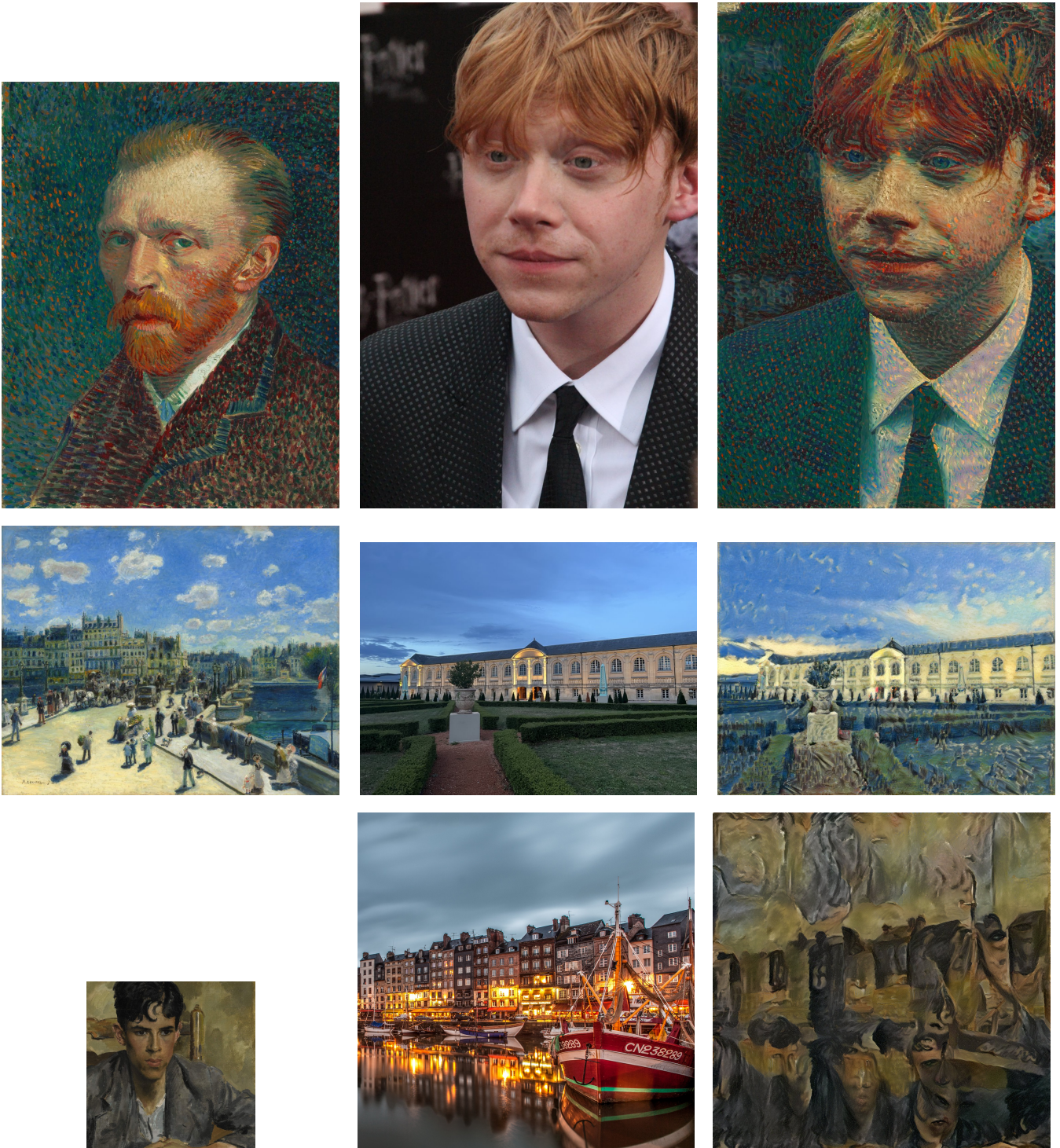


**Figure 13:** Comparison of UHR style transfers. For each example, top row, left to right: style, content, our result (SPST), CD [WLW\*20], URST [CWX\*22]. Bottom row: zoom in of the corresponding top row. First row: content ( $3024 \times 4032$ ), style ( $3024 \times 4477$ ). Third row: content ( $3024 \times 4032$ ), style ( $3024 \times 4738$ ). Fifth row: content ( $3024 \times 4032$ ), style ( $3655 \times 2836$ ). We used three scales for all our results. Observe the loss of details and the unrealistic looks of the outputs produced by both fast methods.



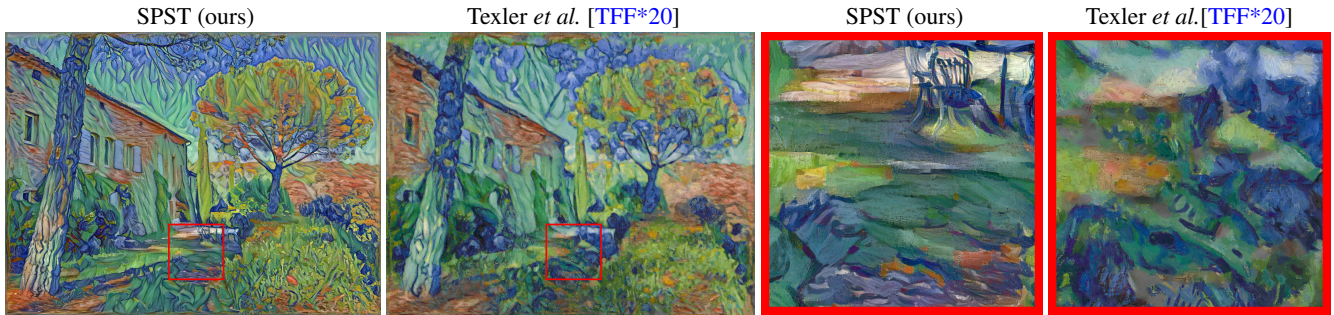


**Figure 14:** Comparison of UHR style transfer on portraits. For each example, top row, left to right: style, content, our result (SPST), CD [WLW\*20], URST [CWX\*22]. Bottom row: zoom in of the corresponding top row. Content ( $3264 \times 4512$ ), style ( $4126 \times 3264$ ), SPST uses three scales. Observe the loss of details and the unrealistic looks of the outputs produced by both fast methods, notably in the background and skin texture.



**Figure 15:** Examples of UHR style transfer with content-style mismatch. For each example, from left to right: style, content and our result. First row: content ( $5184 \times 3456$ ), style ( $4368 \times 3456$ ). Second row: content ( $3024 \times 4032$ ), style ( $3024 \times 3787$ ). Third row: content ( $4096 \times 4096$ ), style ( $2048 \times 2048$ ). We used three scales for the first two results and four scale for the last example.





**Figure 16:** Comparison to AST Using Neurally-Guided Patch-Based Synthesis: From left to right: Our result, result of [TFF\*20] (after a Gatys style transfer of size  $384 \times 512$ ), and close up detail. By design the result of [TFF\*20] is not faithful to the HR details of the content image. Here the chair simply disappear in the output result while it is reproduced using fine brushstrokes in the SPST output





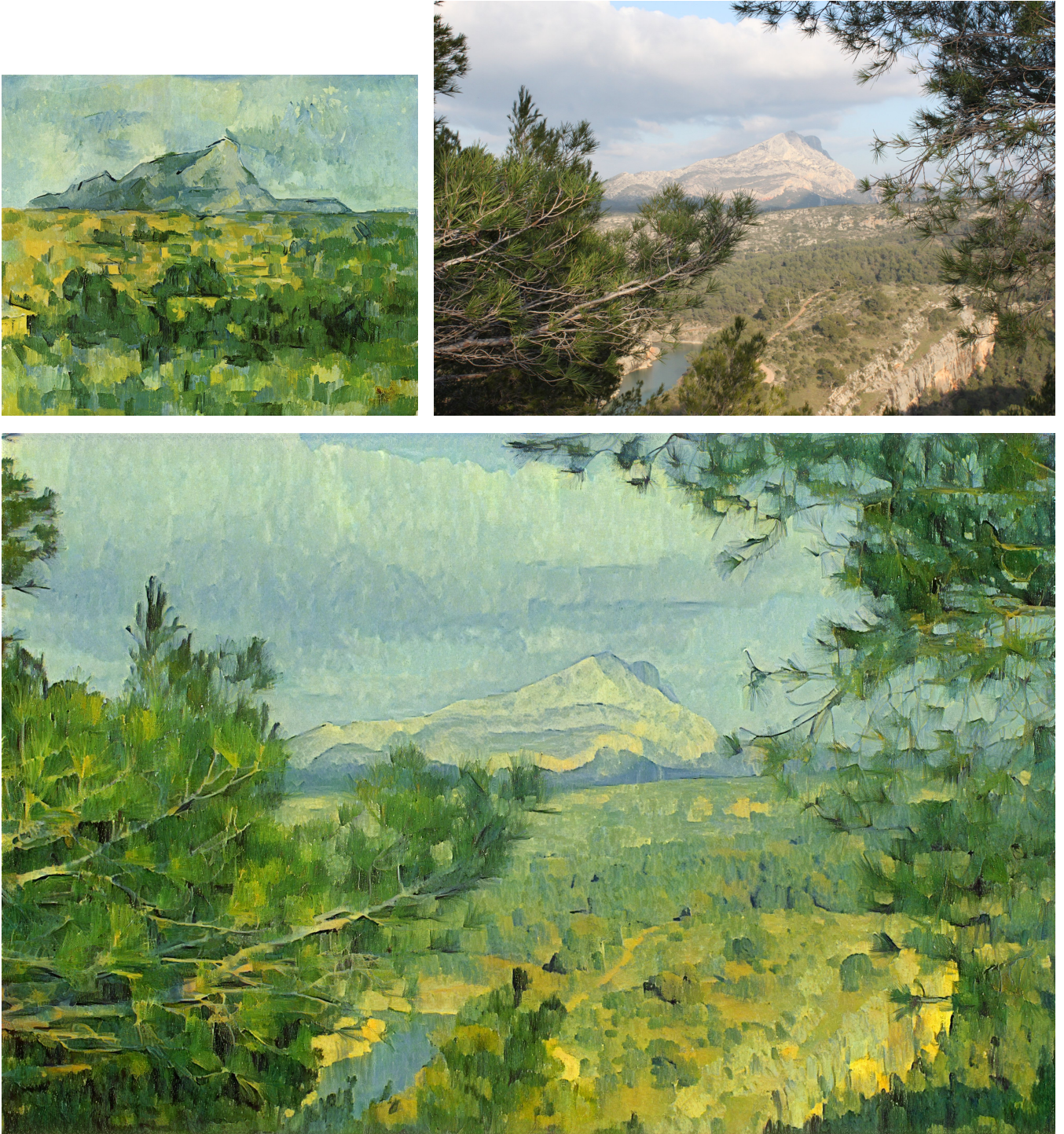
**Figure 17:** UHR style transfer in full resolution (from Figure 1 of the main paper). Top row: style image ( $4226 \times 5319$ ), content image ( $6048 \times 8064$ ). Bottom: result ( $6048 \times 8064$ ). Observe that fine details such as the canvas texture in the sky is well transferred.





**Figure 18:** UHR style transfer in full resolution (from Figure 2 of the main paper). Top row: style image ( $6048 \times 7914$ ), content image ( $6048 \times 8064$ ). Bottom: result ( $6048 \times 8064$ ). Observe how very fine details such as the chairs look as if painted.





**Figure 19:** UHR style transfer in full resolution (from Figure 4 of the main paper). Top row: style image ( $2606 \times 3176$ ), content image ( $3168 \times 4752$ ). Bottom: result ( $3168 \times 4752$ ).





**Figure 20:** UHR style transfer in full resolution (from Figure 4 of the main paper). Top row: style image ( $3024 \times 3787$ ), content image ( $3024 \times 4032$ ). Bottom: result ( $3024 \times 4032$ ).



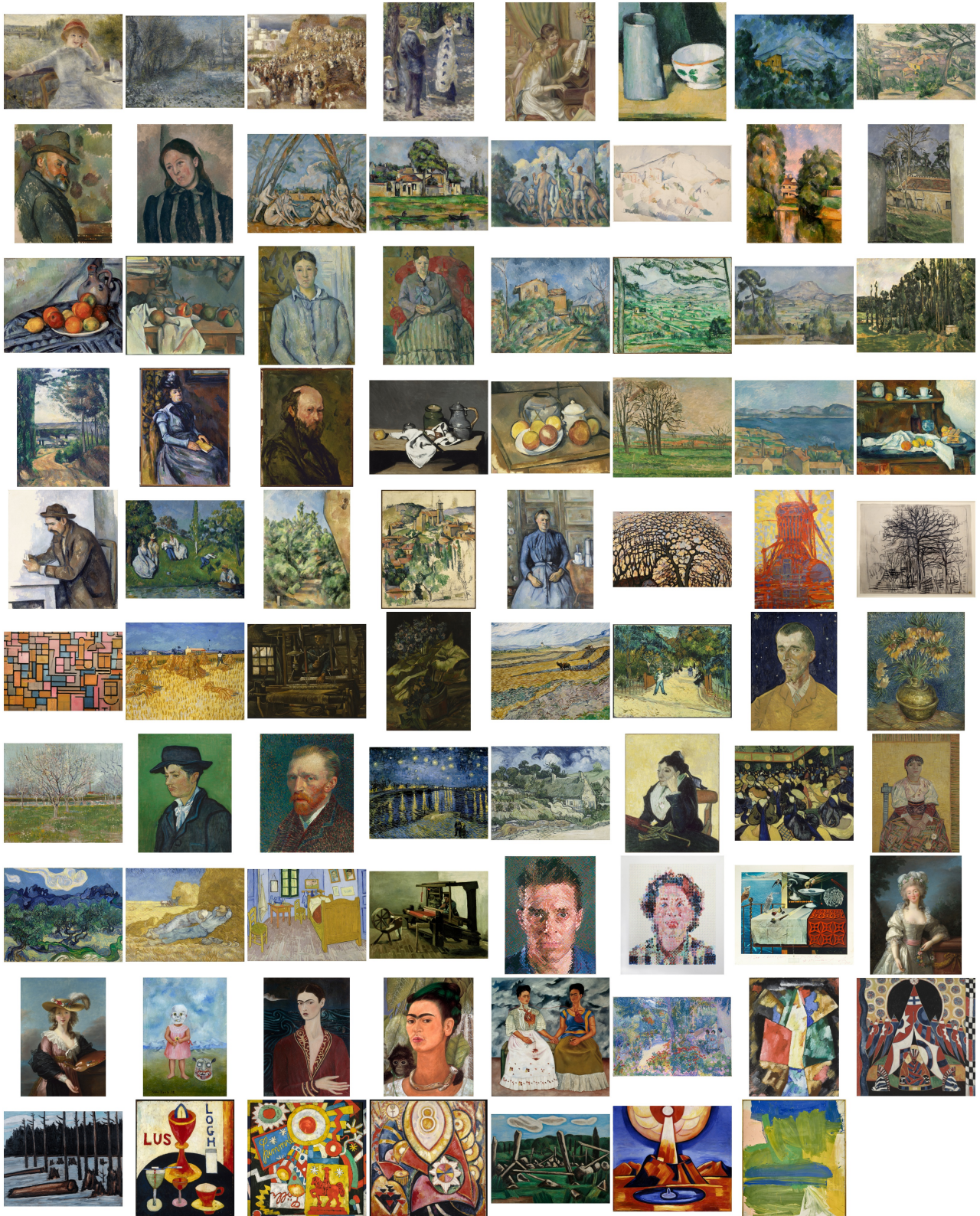


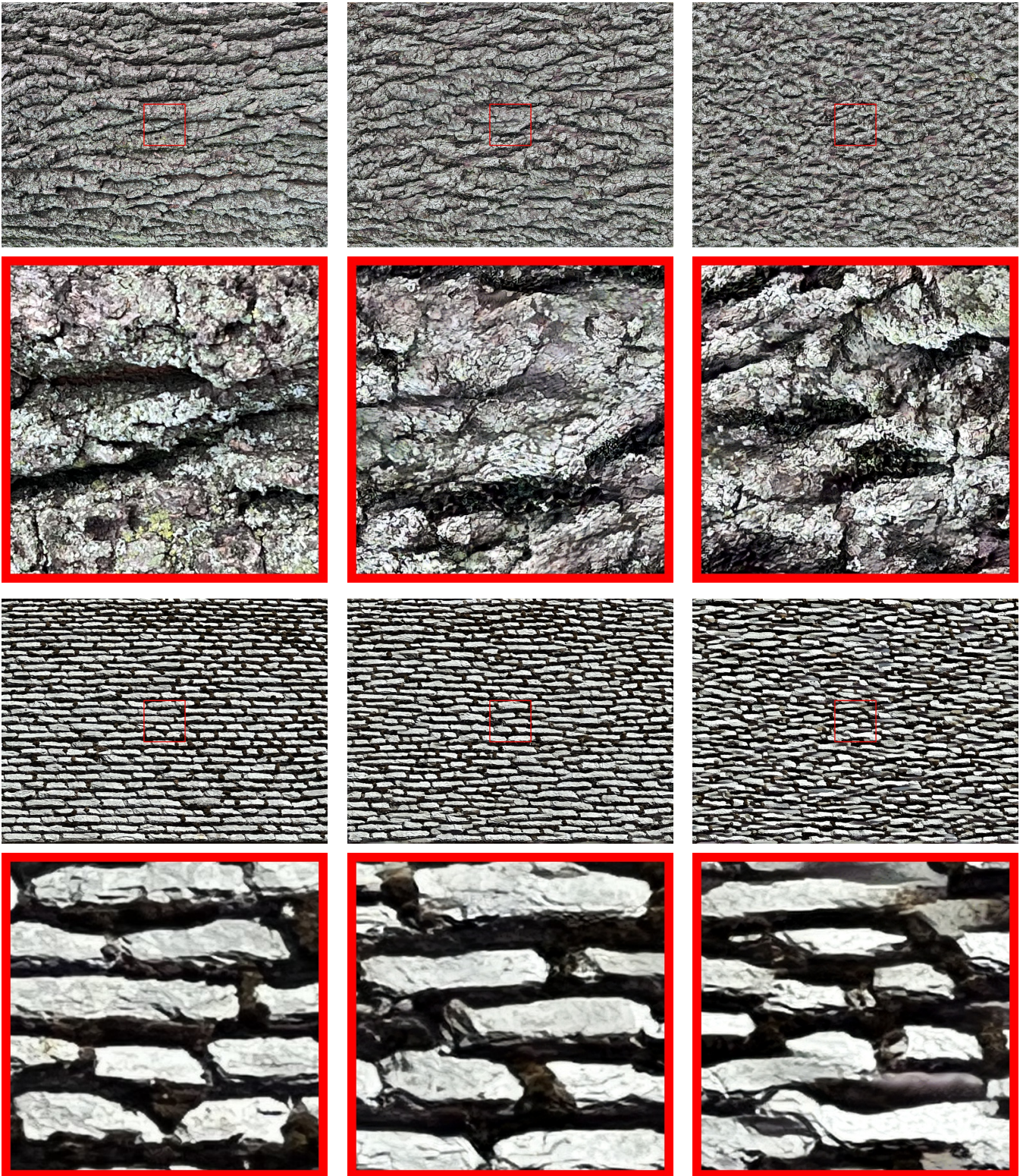
Figure 21: Overview of the 79 UHR painting images used for the identity test.





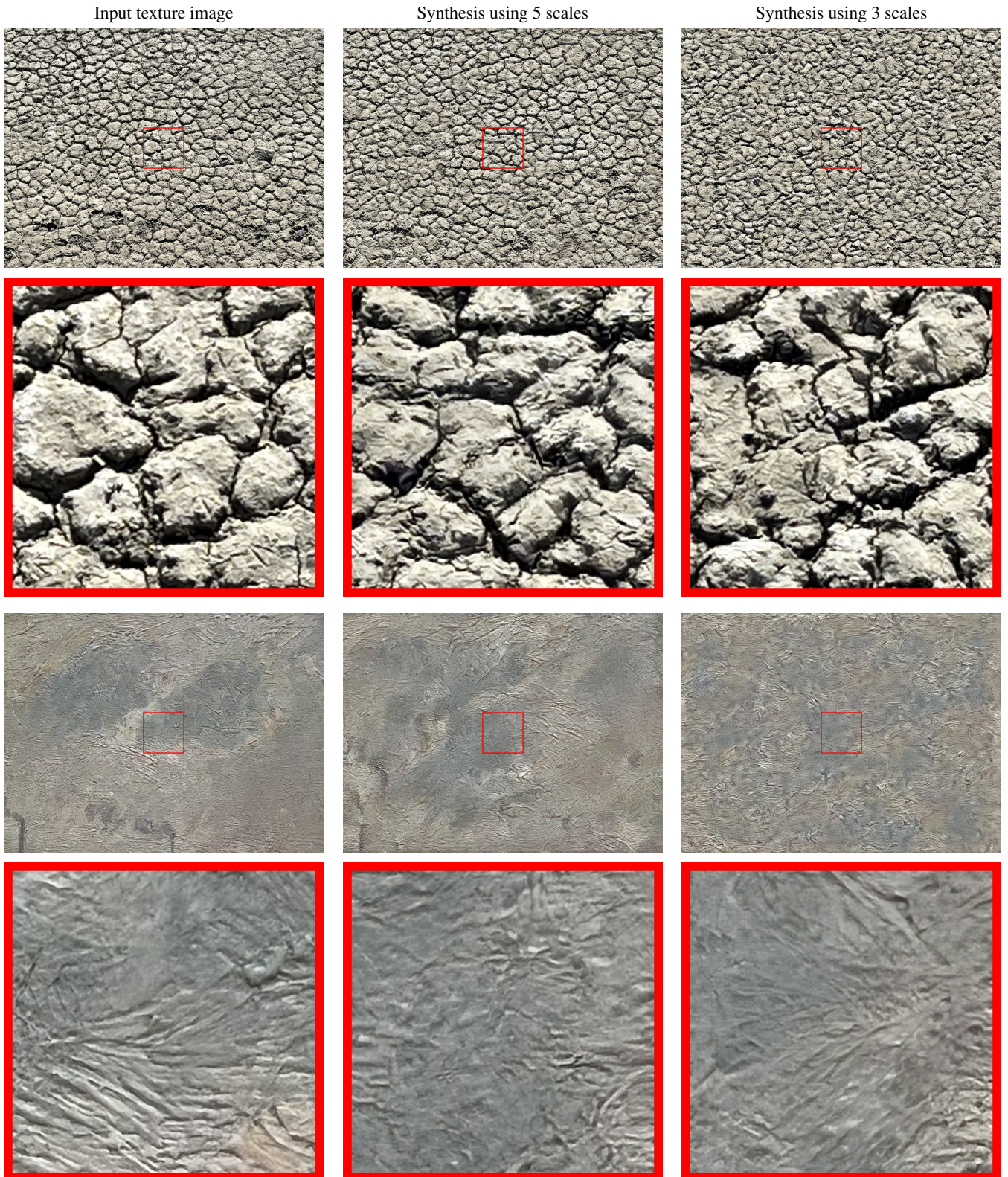
**Figure 22:** UHR texture synthesis: From left to right: Input texture image, synthesis using 5 scales, synthesis using 3 scales. Image have size  $(3024 \times 4032)$  (downscaled by a factor 4 for inclusion in the .pdf) and true resolution details have size  $(512 \times 512)$ .





**Figure 23:** UHR texture synthesis (same as Figure 22): From left to right: Input texture image, synthesis using 5 scales, synthesis using 3 scales. Image have size  $(3024 \times 4032)$  (downscaled by a factor 4 for inclusion in the .pdf) and true resolution details have size  $(512 \times 512)$ .





**Figure 24:** UHR texture synthesis (same as Figure 22): From left to right: Input texture image, synthesis using 5 scales, synthesis using 3 scales. Image have size  $(3024 \times 4032)$  (downscaled by a factor 4 for inclusion in the .pdf) and true resolution details have size  $(512 \times 512)$ .