This ICCV paper is the Open Access version, provided by the Computer Vision Foundation. Except for this watermark, it is identical to the accepted version; the final published version of the proceedings is available on IEEE Xplore.

StructureFlow: Image Inpainting via Structure-aware Appearance Flow

Yurui Ren^{1,2} Xiaoming Yu^{1,2} Ruonan Zhang² Thomas H. Li^{3,1} Shan Liu⁴ Ge Li ⊠^{1,2} ¹School of Electronics and Computer Engineering, Peking University ²Peng Cheng Laboratory ³Advanced Institute of Information Technology, Peking University ⁴Tencent America

> yrren@pku.edu.cn tli@aiit.org.cn

xiaomingyu@pku.edu.cn shanl@tencent.com zhangrn@pcl.ac.cn geli@ece.pku.edu.cn

Abstract

Image inpainting techniques have shown significant improvements by using deep neural networks recently. However, most of them may either fail to reconstruct reasonable structures or restore fine-grained textures. In order to solve this problem, in this paper, we propose a two-stage model which splits the inpainting task into two parts: structure reconstruction and texture generation. In the first stage, edgepreserved smooth images are employed to train a structure reconstructor which completes the missing structures of the inputs. In the second stage, based on the reconstructed structures, a texture generator using appearance flow is designed to yield image details. Experiments on multiple publicly available datasets show the superior performance of the proposed network.

1. Introduction

Image inpainting refers to generating alternative structures and textures for missing regions of corrupted input images and obtaining visually realistic results. It has a wide range of applications. For example, users can remove unwanted objects or edit contents of images by using inpainting techniques. A major challenge of image inpainting tasks is to generate correct structures and realistic textures. Some early patch-based works attempt to fill missing holes with image patches from existing regions [1, 8]. By nearest-neighbor searching and copying relevant patches, these methods can synthesize vivid textures for background inpainting tasks. However, since these methods cannot capture high-level semantics, it is hard for them to generate realistic structures for images with non-repetitive patterns (*e.g.* faces).

With the advent of deep neural network techniques, some recent works [22, 12, 32, 33, 16] model the inpainting task as a conditional generation problem, which learns mapping functions between the input corrupted images and the ground truth images. These methods are able to learn mean-



Figure 1. (From left to right) Input corrupted images, reconstructed structure images, visualizations of the appearance flow fields, final output images. Our method first recovers global structures for missing regions, then generate textures by sampling features from existing regions according to the recovered structures. To visualize the appearance flow fields, we plot the sample points of some typical missing regions. The arrows show the direction of the appearance flow.

ingful semantics, so they can generate coherent structures for missing holes. However, since these methods do not effectively separate the structure and texture information, they often suffer from either over-smoothed boundaries or texture artifacts.

To solve this problem, some two-stage networks [33, 26, 21] are proposed. These methods recover missing structures in the first stage and generate the final results using the reconstructed information in the second stage. The method

proposed in [33] uses ground truth images as the labels of structure recovery. However, ground truth images contain high-frequency textures. These irrelevant details may mislead the structure reconstruction. Spg-net [26] predicts the semantic segmentation labels of the missing areas as structural information. However, regions with similar semantic labels may have different textures (e.g. the windows and walls of the same building), which creates difficulties for the final recovery. Using edge images as the structural guidance, EdgeConnect [21] achieves good results even for some highly structured scenes. However, the distribution of edge images differs greatly from the distribution of the target images. In other words, the edge extractor discards too much useful information, such as image color, making it difficult to generate vivid textures.

In this paper, we propose a novel two-stage network StructureFlow for image inpainting. Our network consists of a structure reconstructor and a texture generator. To recover meaningful structures, we employ edge-preserved smooth images to represent the global structures of image scenes. Edge-preserved smooth methods [30, 31] aim to remove high-frequency textures while retaining sharp edges and low-frequency structures. By using these images as the guidance of the structure reconstructor, the network is able to focus on recovering global structures without being disturbed by irrelevant texture information. After reconstructing the missing structures, the texture generator is used to synthesize high-frequency details. Since image neighborhoods with similar structures are highly correlated, the uncorrupted regions can be used to generate textures for missing regions. However, it is hard for convolutional neural networks to model long-term correlations [33]. In order to establish a clear relationship between different regions, we propose to use appearance flow [35] to sample features from regions with similar structures, as shown in Figure 1. Since appearance flow is easily stuck within bad local minima in the inpainting task [33], in this work, we made two modifications to ensure the convergence of the training process. First, Gaussian sampling is employed instead of Bilinear sampling to expand the receptive field of the sampling operation. Second, we introduce a new loss function, called sampling correctness loss, to determine if the correct regions are sampled.

Both subjective and objective experiments compared with several state-of-the-art methods show that our method can achieve competitive results. Furthermore, we perform ablation studies to verify our hypothesis and modifications. The main contributions of our paper can be summarized as:

- We propose a structure reconstructor to generate edgepreserved smooth images as the global structure information.
- We introduce appearance flow to establish long-term

corrections between missing regions and existing regions for vivid texture generation.

- To ease the optimization of appearance flow, we propose to use Gaussian sampling instead of Bilinear sampling and introduce a novel sampling correctness loss.
- Experiments on multiple public datasets show that our method is able to achieve competitive results.

2. Related Work

2.1. Image Inpainting

Existing image inpainting works can be roughly divided into two categories: methods using diffusion-based or patch-based techniques and methods using deep neural networks. Diffusion-based methods [2, 6] synthesize textures by propagating the neighborhood region appearance to the target holes. However, these methods can only deal with small holes in background inpainting tasks. They may fail to generate meaningful structures. Unlike the diffusionbased methods using only neighborhood pixels of missing holes, patch-based methods can take advantage of remote information to recover the lost areas. Patch-based methods [1, 8, 3] fill target regions by searching and copying similar image patches from the uncorrupted regions of the source images. These methods can generate photo-realistic textures for relatively large missing holes. In order to find suitable image patches, bidirectional similarity [24] is proposed to capture more visual information and introduce less visual artifacts when calculating the patch similarity. To reduce the computational cost, PatchMatch [1] designs a fast nearest neighbor searching algorithm using natural coherence in the imagery as prior information. However, these patch-based methods assume that the non-hole regions have similar semantic contents with the missing regions, which may not be true in some tasks such as face image inpainting. Therefore, they may work well in some images with repetitive structures but cannot generate reasonable results for images with unique structures.

Recently, many deep learning based methods have been proposed to model the inpainting task as a conditional generation problem. A significant advantage of these methods is that they are able to extract meaningful semantics from the corrupted images and generate new content for images. Context Encoder [22] is one of the early attempts. It uses an encoder-decoder architecture to first extract features and then to reconstruct the outputs. However, this network struggles to maintain global consistency and often generate results with visual artifacts. Iizuka *et al.* [12] solve this problem by using both local and global discriminators which are responsible for generating realistic alternative contents for missing holes and maintaining the coherency of competed images respectively. Yu *et al.* [33] find



Figure 2. Overview of our StructureFlow. Our model first generates global structures (*i.e.* edge-preserved smooth images) using structure reconstructor. Then texture generator is used to yield high-frequency details and output the final results. We add the appearance flow to our texture generator to sample features from existing regions.

that convolutional neural networks are ineffective in building long-term correlations. To solve this problem, they propose contextual attention to borrow features from remote regions. Liu et al. [16] believe the substituting pixels in the masked holes of the inputs introduce artifacts to the final results. Therefore, they propose partial convolutions to force the network to use valid pixels (uncorrupted pixels) only. Gated convolution [32] further generalizes this idea by extending the feature selecting mechanism to be learnable for each location across all layers. EdgeConnect proposed in paper [21] has a similar motivation to our paper: generating reasonable structures by using additional prior information. EdgeConnect first recovers edge maps and then fills the missing regions in fine details. However, due to the limited representation ability of edge maps, this method may generate wrong details in the boundaries of objects.

2.2. Optical Flow and Appearance Flow

Optical flow is used to describe the motion of objects, surfaces, and edges between consecutive video frames. It has been widely used in video frame synthesis [37, 29], action recognition [25, 28], etc. Optical flow estimation is an important task in computer vision. Many methods [11, 27] have been proposed to accurately estimate optical flow between consecutive frames. Recently, some methods [5, 13]solve this problem by training deep neural networks. However, these techniques require sufficient ground truth optical flow fields which are extremely difficult to obtain. Therefore, some synthetic optical flow datasets [5] are created for training. Some other methods [18, 19] solve this problem by training the network in an unsupervised manner. However, many existing unsupervised optical flow estimation methods struggle to capture large motions. Some papers [18, 23] manage to use multi-scale approaches to improve the results. We believe it is due to the limited receptive field of Bilinear sampling. In this paper, we use Gaussian sampling as an improvement.

Appearance flow proposed by [35] is used to generate target scenes (objects) from source scenes (objects) using a flow-based method. It calculates the correlations between sources and targets to predict the 2-D coordinate vectors (*i.e.* appearance flow fields). This idea can be used in image inpainting tasks. To generate realistic alternative contents for missing holes, one can reasonably "flow" pixels (features) from source regions to missing regions. In this paper, we improve the appearance flow in [35] to make it suitable for image inpainting tasks.

3. Our Approach

The framework of our StructureFlow inpainting network is shown in Figure 2. Our model consists of two parts: the structure reconstructor G_s and the texture generator G_t . The structure reconstructor G_s is used to predict missing structures, thereby generating the global structure image $\hat{\mathbf{S}}$. The texture generator G_t draws details according to the reconstructed structures $\hat{\mathbf{S}}$ and outputs the final results $\hat{\mathbf{I}}$.

3.1. Structure Reconstructor

A major challenge of image inpainting tasks is to generate meaningful structures for missing regions. Therefore, we first design a structure reconstructor G_s to recover global structures of the input images. The edge-preserved smooth methods [30, 31] aim to remove high-frequency textures while retaining the sharp edges and low-frequency structures. Their results can well represent global structures. Let I_{gt} be the ground-truth image and S_{gt} be the edge-preserved smooth result of I_{gt} . The processing of our structure reconstructor G_s can be written as

$$\hat{\mathbf{S}} = G_s(\mathbf{I}_{in}, \mathbf{S}_{in}, \mathbf{M}) \tag{1}$$

where **M** is the mask of the input image I_{in} . It is a binarized matrix where 1 represents the missing region and 0 represents the background. $S_{in} = S_{gt} \circ (1 - M)$ is the structures of I_{in} . Here, \circ denotes element-wise product. \hat{S} is the predicted structures.

The reconstruction loss of G_s is defined as the ℓ_1 distance between the predicted structures $\hat{\mathbf{S}}$ and the ground-truth structures \mathbf{S}_{gt} .

$$\mathcal{L}_{\ell_1}^s = \|\hat{\mathbf{S}} - \mathbf{S}_{gt}\|_1 \tag{2}$$

Meanwhile, to mimic the distributions of the target structures S_{gt} , we apply generative adversarial framework [7] to our structure reconstructor. The adversarial loss of G_s can be written as

$$\mathcal{L}_{adv}^{s} = \mathbb{E}[\log(1 - D_{s}(G_{s}(\mathbf{I}_{in}, \mathbf{S}_{in}, \mathbf{M})))] + \mathbb{E}[\log D_{s}(\mathbf{S}_{qt})]$$
(3)

where D_s is the discriminator of the structure reconstructor. We jointly train the generator G_s and discriminator D_s using the following optimization.

$$\min_{G_s} \max_{D_s} \mathcal{L}^s(G_s, D_s) = \lambda^s_{\ell_1} \mathcal{L}^s_{\ell_1} + \lambda^s_{adv} \mathcal{L}^s_{adv}$$
(4)

where $\lambda_{\ell_1}^s$ and λ_{adv}^s are regularization parameters. We set $\lambda_{\ell_1}^s = 4$ and $\lambda_{adv}^s = 1$ in all experiments.

3.2. Texture Generator

After obtaining the reconstructed structure image S, our texture generator G_t is employed to yield vivid textures. The processing of the texture generator G_t can be written as

$$\hat{\mathbf{I}} = G_t(\mathbf{I}_{in}, \hat{\mathbf{S}}, \mathbf{M}) \tag{5}$$

where $\tilde{\mathbf{I}}$ denotes the final output result. We use ℓ_1 loss to calculate the reconstruction error.

$$\mathcal{L}_{\ell_1}^t = \|\hat{\mathbf{I}} - \mathbf{I}_{gt}\|_1 \tag{6}$$

To generate realistic results, we employ adversarial loss in our texture generator.

$$\mathcal{L}_{adv}^{t} = \mathbb{E}[\log(1 - D_t(G_t(\mathbf{I}_{in}, \mathbf{S}, \mathbf{M})))] \\ + \mathbb{E}[\log D_t(\mathbf{I}_{gt})]$$
(7)

Since image regions with similar structures are highly related, it is possible to extract these correlations using the reconstructed structures \hat{S} for texture generation to improve the performance. However, convolutional neural networks are not effective for capturing long-term dependency [33]. In order to establish a clear relationship between different regions, we introduce the appearance flow to our G_t . As shown in Figure 2, the appearance flow is used to warp the extracted features of the inputs. Thus, features containing vivid texture information can "flow" to the corrupted regions.

However, training the appearance flow in an unsupervised manner is a difficult task [18, 23]. The networks may struggle to capture large motions and stuck in a bad local minima. To tackle this problem, we first propose to use Gaussian sampling instead of Bilinear sampling to expand the receptive field. Then, we propose a sampling correctness loss to constraint the possible convergence results.

The sampling process calculates the gradients according to the input pixels (features). If the receptive field of the sampling operation is limited, only a few pixels can participate in the operation. Since the adjacent pixels (features) are often highly correlated, a large receptive field is required to obtain correct and stable gradients. Therefore, Bilinear sampling with a very limited receptive field may not be suitable for tasks requiring establishing long-term correlations. To expand the receptive field, we use Gaussian sampling instead of Bilinear sampling in the appearance flow operation. The process of Gaussian sampling operation with kernel size n can be written as

$$\mathbf{F}_{o} = \sum_{i=1}^{n} \sum_{j=1}^{n} \frac{a_{i,j}}{\sum_{i=1}^{n} \sum_{j=1}^{n} a_{i,j}} \mathbf{F}_{i,j}$$
(8)

where $\mathbf{F}_{i,j}$ is the features around the sample center and \mathbf{F}_o is the output feature. The weights $a_{i,j}$ is calculated as

$$a_{i,j} = exp(-\frac{\Delta h^2 + \Delta v^2}{2\sigma^2}) \tag{9}$$

where Δh and Δv is the horizontal and vertical distance between the sampling center and feature $\mathbf{F}_{i,j}$ respectively. Parameter σ is used to denote the variance of the Gaussian sampling kernel.

The proposed sampling correctness loss is used to constraint the appearance flow fields. It determines whether the current sampled regions are "good" choices. We use the pre-trained VGG19 to calculate this loss. Specifically, we first calculate the VGG features of the input corrupted image I_{in} and the ground truth image I_{gt} . Let V^{in} and V^{gt} be the features generated by a specific layer of VGG19. Symbol M denotes a coordinate set containing the coordinates of missing areas, N is the number of elements in set M. Then, our sampling correctness loss calculate the relative cosine similarity between the ground truth features and the sampled features

$$\mathcal{L}_{c}^{t} = \frac{1}{N} \sum_{(x,y)\in M} exp(-\frac{\mu(\mathbf{V}_{x,y}^{gt}, \mathbf{V}_{x+\Delta x, y+\Delta y}^{in})}{\mu_{x,y}^{max}}) \quad (10)$$

where $\mathbf{V}_{x+\Delta x,y+\Delta y}^{in}$ is the sampled feature calculated by our Gaussian sampling and $\mu(*)$ denotes the cosine similarity. $\mu_{x,y}^{max}$ is a normalization term. For each feature $\mathbf{V}_{x,y}^{gt}$

	PSNR			SSIM			FID		
Mask	0-20%	20-40%	40-60%	0-20%	20-40%	40-60%	0-20%	20-40%	40-60%
CA	27.150	20.001	16.911	0.9269	0.7613	0.5718	4.8586	18.4190	37.9432
PConv	31.030	23.673	19.743	0.9070	0.7310	0.5325	-	-	-
EdgeConnect	29.972	23.321	19.641	0.9603	0.8600	0.6916	3.0097	7.2635	19.0003
Ours	32.029	25.218	21.090	0.9738	0.9026	0.7561	2.9420	7.0354	22.3803

Table 1. The evaluation results of CA [33], PConv [16], EdgeConnect [21], and our model over dataset Places2 [34]. Since the code and models of PConv are not available, we report the results presented in their paper.

where $(x, y) \in M$, we find the most similar feature from \mathbf{V}_{in} and calculate their cosine similarity as $\mu_{x,y}^{max}$.

$$\mu_{x,y}^{max} = \max_{(x',y')\in\Omega} \mu(\mathbf{V}_{x,y}^{gt}, \mathbf{V}_{x',y'}^{in})$$
(11)

where Ω denotes a coordinate set containing all coordinates in \mathbf{V}_{in} . Our texture generator is trained using the following optimization

$$\min_{G_t} \max_{D_t} \mathcal{L}^t(G_t, D_t) = \lambda_{\ell_1}^t \mathcal{L}_{\ell_1}^t + \lambda_c^t \mathcal{L}_c^t + \lambda_{adv}^t \mathcal{L}_{adv}^t$$
(12)

where $\lambda_{\ell_1}^t$, λ_c^t and λ_{adv}^t are the hyperparameters. In our experiments, we set $\lambda_{\ell_1}^t = 5$, $\lambda_c^t = 0.25$ and $\lambda_{adv}^t = 1$.

4. Experiments

4.1. Implementation Details

Basically, autoencoder structures are employed to design our generators G_s and G_t . Several residual blocks [9] are added to further process the features. For the appearance flow, we concatenate the warped features with the features obtained by convolutional blocks. The architecture of our discriminators is similar to that of BicycleGAN [36]. We use two PatchGANs [14] with different scales to predict real *vs*. fake for overlapping image patches with different sizes. In order to solve the notorious problem of instability training of generative adversarial networks, spectral normalization [20] is used in our network.

We train our model on three public datasets including Places2 [34], Celeba [17], and Paris StreetView [4]. The most challenging dataset Places2 contains more than 10 million images comprising 400+ unique scene categories. Celeba and Paris StreetView contain highly structured face and building images respectively. We use the irregular mask dataset provided by [16]. The mask images are classified based on their hole sizes relative to the entire image (*e.g.* 0 - 20% *etc.*).

We employ edge-preserved smooth method RTV [31] to obtain the training labels of the structure reconstructor G_t . In RTV smooth method, parameter σ is used to control the spatial scale of smooth windows, thereby controlling the maximum size of texture elements. In section 4.3, we explore the impact of σ on the final results. We empirically find the best results obtained when we set $\sigma \approx 3$.

We train our model in stages. First, the structure reconstructor G_s and the texture generator G_t are trained separately using the edge-preserved image \mathbf{S}_{gt} . Then, we continue to fine-tune G_t using the reconstructed structures $\hat{\mathbf{S}}$. The network is trained using 256×256 images with batch size as 12. We use the Adam optimizer [15] with learning rate as 10^{-4} .

4.2. Comparisons

We subjectively and objectively compare our approach with several state-of-the-art methods including Contextual Attention (CA) [33], Partial Convolution (PConv) [16] and EdgeConnect [21].

Objective comparisons Image inpainting tasks lack specialized quantitative evaluation metrics. In order to compare the results as accurately as possible, we employ two types of metrics: distortion measurement metrics and perceptual quality measurement metrics. Structural similarity index (SSIM) and peak signal-to-noise ratio (PSNR) assume that the ideal recovered results are exactly the same as the target images. They are used to measure the distortions of the results. Fréchet Inception Distance (FID) [10] calculates the Wasserstein-2 distance between two distributions. Therefore, it can indicate the perceptual quality of the results. In this paper, we use the pre-trained Inception-V3 model to extract features of real and inpainted images when calculating FID scores. The final evaluation results over Places2 are reported in Table 1. We calculate the statistics over 10krandom images in the test set. It can be seen that our model achieves competitive results compared with other models.

Subjective comparisons We implement a human subjective study on the Amazon Mechanical Turk (MTurk). We ask volunteers to choose the more realistic image from image pairs of real and generated images. For each dataset, we randomly select 600 images and assign them random mask ratios from 0% - 60% for the evaluation. Each image is compared 5 times by different volunteers. The evaluation results are shown in Table 2. Our model achieves better results than the competitors in the highly-structured



Figure 3. The qualitative comparisons with existing models. (From top to bottom) Input corrupted images, results of CA [33], results of EdgeConnect [21], results of our StructureFlow, and Ground truth images.

	CA	EdgeConnect	Ours
Celeba	5.68%	26.28%	32.04%
Paris	17.36%	33.44%	33.68%
Places2	8.72%	26.36%	23.56%

Table 2. The evaluation results of user study. The volunteers are asked to select the more realistic image from image pairs of real and generated images. The fooling rate is provided in the table.

scenes, such as face dataset Celeba and street view dataset Paris. This indicates that our model can generate meaningful structures for missing regions. We also achieve competitive results in dataset Places2.

Figure 3 shows some example results of different models. It can be seen that the results of CA suffer from artifacts, which means that this method may struggle to balance the generation of textures and structures. EdgeConnect is able to recover correct global structures. However, it may generate wrong details at the edges of objects. Our method can generate meaningful structures as well as vivid textures. We also provide the reconstructed structures of EdgeConnect and our model in Figure 4. We find that the edge maps loss too much useful information, such as image color when recovering the global structures. Therefore, EdgeConnect may fill incorrect details for some missing areas. Meanwhile, edges of different objects may be mixed together in edge maps, which makes it difficult to generate textures. In contrast, our edge-preserved smooth images can well represent the structures of images. Therefore, our model can well balance structure reconstruction and texture generation. Photo-realistic results are obtained even for some highly structured images with large hole ratios.

4.3. Ablation Studies

In this section, we analyze how each component of our StructureFlow contributes to the final performance from two perspectives: structures and appearance flow.

Structure Ablation In this paper, we assume that the structure information is important for image inpainting tasks. Therefore, we first reconstruct structures and use them as prior information to generate the final results. To verify this assumption, we remove our structure reconstructor and train an inpainting model with only the texture generator. The corrupted images along with its masks are directly inputted into the model. Please note that we also keep appearance flow in the network for fair comparisons. The results are shown in Table 3. It can be seen that our structure structure structure shown in Table 3.



Figure 4. The visual comparisons between EdgeConnect [21] and our StructureFlow. (From left to right) Input corrupted images, reconstructed structures of EdgeConnect, inpainted results of Edge-Connect, reconstructed structures of our StructureFlow, inpainted results of our StructureFlow.

ture reconstructor can bring stable performance gain to the model.

Then we turn our attention to the edge-preserved smooth images. We believe the edge-preserved smooth images are able to represent the structures since the smooth operations remove high-frequency textures. To verify this, we train a model using ground truth images I_{gt} as the labels of the structure reconstructor. The results can be found in Table 3. Compared with StructureFlow, we can find that using images containing high-frequency textures as structures leads to performance degradation.

However, it is difficult to accurately distinguish the textures and the structures of an image. What is the appropriate degree of smooth operation? We find there exists a trade-off between the structure reconstructor and the texture generator. If very few textures are removed, the structure reconstruction will be more difficult, since it needs to



Figure 5. The influence of the parameter σ in RTV edge-preserved smooth method on final results. The last three columns show results of models trained by smooth images generated with $\sigma = 0, 3, 9$, where the first row shows the results of structure reconstructor and the second row shows the generated results.

		PSNR	SSIM
Paris	w/o Structure	28.46	0.8879
	w/o Smooth	28.41	0.8848
	w/o Flow	28.77	0.8906
	StructureFlow	29.25	0.8979
Celeba	w/o Structure	29.42	0.9324
	w/o Smooth	29.61	0.9335
	w/o Flow	29.91	0.9368
	StructureFlow	30.31	0.9420

Table 3. The evaluation results of ablation studies. We provide the statistics of four models: the model trained without the structure reconstructor (*i.e.* w/o Structure), the model trained using ground truth images as the labels of the structure reconstructor (*i.e.* w/o Smooth), the model trained without the appearance flow operation (*i.e.* w/o Flow) and our full model (*i.e.* StructureFlow). The statistics are based on random masks with mask ratio 0%-60%.

	$\sigma = 0$	$\sigma = 1$	$\sigma = 3$	$\sigma = 6$	$\sigma = 9$
PSNR	28.41	28.81	29.25	29.14	28.98
SSIM	0.8848	0.8896	0.8979	0.8962	0.8990

Table 4. The evaluation results over dataset Paris of models trained using edge-preserved images generated with $\sigma = 0, 1, 3, 6, 9$. The statistics are based on random masks with mask ratio 0%-60%.

recover more information. However, the texture generation will be easier. Therefore, we need to balance the difficulties of these two tasks to achieve better results. We use σ in RTV [31] smooth method to control the maximum size of texture elements in \mathbf{S}_{gt} . Smoother results are obtained with larger σ value. We train our StructureFlow using smooth images generated from $\sigma = 0, 1, 3, 6, 9$. The evaluation results over dataset Paris are obtained when $\sigma = 3$. Both too small and too large σ values lead to model performance degradation. An example can be found in Figure 5. When



Figure 6. Ablation studies of Gaussian sampling and the sampling correctness loss. The appearance flow fields are visualized using the provided color coding map. Flow direction is encoded with color and magnitude with color intensity.

 $\sigma=0$, the structure reconstructor fail to generate reasonable structures, as it is disturbed by irrelevant texture information. The texture generator fails to yield realistic images when trained with $\sigma=9$ since some useful structural information is removed.

Flow Ablation In this ablation study, we first evaluate the performance gain bought by our appearance flow. Then, we illustrate the effectiveness of Gaussian sampling and the sampling correctness loss.

To verify the validity of our appearance flow, we train a model without using the appearance flow blocks in the texture generator. The evaluation results can be found in Table 3. It can be seen that our StructureFlow has better performance than the model trained without the appearance flow operation, which means that our appearance flow can help with the texture generation and improve model performance.

Next, we test our Gaussian sampling and the sampling correctness loss. Two models are trained for this ablation study: a model trained using Bilinear sampling in the warp operation of appearance flow and a model trained without using the sampling correctness loss. Figure 6 shows the appearance flow fields obtained by these models. It can be seen that the model trained without using the sampling correctness loss is unable to sample correct features for large missing regions. Bilinear sampling also fails to capture long-term correlations. Our StructureFlow obtains a reasonable flow field and generates realistic textures for missing regions.

4.4. User case

Our method can be used for some image editing applications. Figure 7 provides some usage examples. Users can remove the unwanted objects by interactively drawing masks in the inputs. Our model is able to generate realistic alternative contents for the missing regions. In addition, by



Figure 7. Examples of object removing and image editing using our StructureFlow. Our model is able to generate realistic alternative contents for missing regions.

directly editing the structure images, users can copy or add new objects and contents to images.

5. Conclusion

In this paper, we propose an effective structure-aware framework for recovering corrupted images with meaningful structures and vivid textures. Our method divides the inpainting task into two subtasks: structure reconstruction and texture generation. We demonstrate that edge-preserved smooth images can well represent the global structure information and play an important role in inpainting tasks. As for texture generation, we use appearance flow to sample features from relative regions. We verify that our flow operation can bring stable performance gain to the final results. Our method can obtain competitive results compared with several state-of-the-art methods. Our source code is available at: https://github.com/RenYurui/ StructureFlow.

Acknowledgements. This work was supported by National Engineering Laboratory for Video Technology-Shenzhen Division, Shenzhen Municipal Science and Technology Program (JCYJ20170818141146428), and Shenzhen Key Laboratory for Intelligent Multimedia and Virtual Reality (ZDSYS201703031405467). In addition, we thank the anonymous reviewers for their valuable comments.

References

- Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. Patchmatch: A randomized correspondence algorithm for structural image editing. In *ACM Transactions on Graphics (ToG)*, volume 28, page 24. ACM, 2009.
 1, 2
- [2] Marcelo Bertalmio, Guillermo Sapiro, Vincent Caselles, and Coloma Ballester. Image inpainting. In *Proceedings of* the 27th annual conference on Computer graphics and interactive techniques, pages 417–424. ACM Press/Addison-Wesley Publishing Co., 2000. 2
- [3] Soheil Darabi, Eli Shechtman, Connelly Barnes, Dan B Goldman, and Pradeep Sen. Image melding: Combining inconsistent images using patch-based synthesis. ACM Trans. Graph., 31(4):82–1, 2012. 2
- [4] Carl Doersch, Saurabh Singh, Abhinav Gupta, Josef Sivic, and Alexei Efros. What makes paris look like paris? ACM Transactions on Graphics, 31(4), 2012. 5
- [5] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. Flownet: Learning optical flow with convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2758–2766, 2015. 3
- [6] Alexei A Efros and William T Freeman. Image quilting for texture synthesis and transfer. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 341–346. ACM, 2001. 2
- [7] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 4
- [8] James Hays and Alexei A Efros. Scene completion using millions of photographs. ACM Transactions on Graphics (TOG), 26(3):4, 2007. 1, 2
- [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceed-ings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5
- [10] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In Advances in Neural Information Processing Systems, pages 6626–6637, 2017. 5
- [11] Berthold KP Horn and Brian G Schunck. Determining optical flow. Artificial intelligence, 17(1-3):185–203, 1981. 3
- [12] Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa. Globally and Locally Consistent Image Completion. ACM Transactions on Graphics (Proc. of SIGGRAPH 2017), 36(4):107:1–107:14, 2017. 1, 2
- [13] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. Flownet 2.0: Evolution of optical flow estimation with deep networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2462–2470, 2017. 3

- [14] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. 5
- [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014. 5
- [16] Guilin Liu, Fitsum A Reda, Kevin J Shih, Ting-Chun Wang, Andrew Tao, and Bryan Catanzaro. Image inpainting for irregular holes using partial convolutions. In *Proceedings* of the European Conference on Computer Vision (ECCV), pages 85–100, 2018. 1, 3, 5
- [17] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738, 2015. 5
- [18] Ziwei Liu, Raymond A Yeh, Xiaoou Tang, Yiming Liu, and Aseem Agarwala. Video frame synthesis using deep voxel flow. In *Proceedings of the IEEE International Conference* on Computer Vision, pages 4463–4471, 2017. 3, 4
- [19] Simon Meister, Junhwa Hur, and Stefan Roth. Unflow: Unsupervised learning of optical flow with a bidirectional census loss. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018. 3
- [20] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.
 5
- [21] Kamyar Nazeri, Eric Ng, Tony Joseph, Faisal Qureshi, and Mehran Ebrahimi. Edgeconnect: Generative image inpainting with adversarial edge learning. arXiv preprint arXiv:1901.00212, 2019. 1, 2, 3, 5, 6, 7
- [22] Deepak Pathak, Philipp Krähenbühl, Jeff Donahue, Trevor Darrell, and Alexei Efros. Context encoders: Feature learning by inpainting. 2016. 1, 2
- [23] Anurag Ranjan and Michael J Black. Optical flow estimation using a spatial pyramid network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4161–4170, 2017. 3, 4
- [24] D. Simakov, Y. Caspi, E. Shechtman, and M. Irani. Summarizing visual data using bidirectional similarity. In 2008 IEEE Conference on Computer Vision and Pattern Recognition, pages 1–8, June 2008. 2
- [25] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In Advances in neural information processing systems, pages 568– 576, 2014. 3
- [26] Yuhang Song, Chao Yang, Yeji Shen, Peng Wang, Qin Huang, and C-C Jay Kuo. Spg-net: Segmentation prediction and guidance network for image inpainting. *arXiv preprint arXiv:1805.03356*, 2018. 1, 2
- [27] Deqing Sun, Stefan Roth, JP Lewis, and Michael J Black. Learning optical flow. In *European Conference on Computer Vision*, pages 83–97. Springer, 2008. 3
- [28] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal segment networks: Towards good practices for deep action recognition.

In European conference on computer vision, pages 20–36. Springer, 2016. 3

- [29] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-tovideo synthesis. In Advances in Neural Information Processing Systems (NeurIPS), 2018. 3
- [30] Li Xu, Cewu Lu, Yi Xu, and Jiaya Jia. Image smoothing via 10 gradient minimization. ACM Transactions on Graphics (SIGGRAPH Asia), 2011. 2, 3
- [31] Li Xu, Qiong Yan, Yang Xia, and Jiaya Jia. Structure extraction from texture via relative total variation. ACM Transactions on Graphics (TOG), 31(6):139, 2012. 2, 3, 5, 7
- [32] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Free-form image inpainting with gated convolution. *arXiv preprint arXiv:1806.03589*, 2018. 1, 3
- [33] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Generative image inpainting with contextual attention. In *Proceedings of the IEEE Conference* on Computer Vision and Pattern Recognition, pages 5505– 5514, 2018. 1, 2, 4, 5, 6
- [34] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE transactions on pattern analysis* and machine intelligence, 40(6):1452–1464, 2018. 5
- [35] Tinghui Zhou, Shubham Tulsiani, Weilun Sun, Jitendra Malik, and Alexei A Efros. View synthesis by appearance flow. In *European conference on computer vision*, pages 286–301. Springer, 2016. 2, 3
- [36] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. In Advances in Neural Information Processing Systems, pages 465–476, 2017. 5
- [37] Xiaoou Tang Yiming Liu Ziwei Liu, Raymond Yeh and Aseem Agarwala. Video frame synthesis using deep voxel flow. In *Proceedings of International Conference on Computer Vision (ICCV)*, October 2017. 3