

# T-former: An Efficient Transformer for Image Inpainting

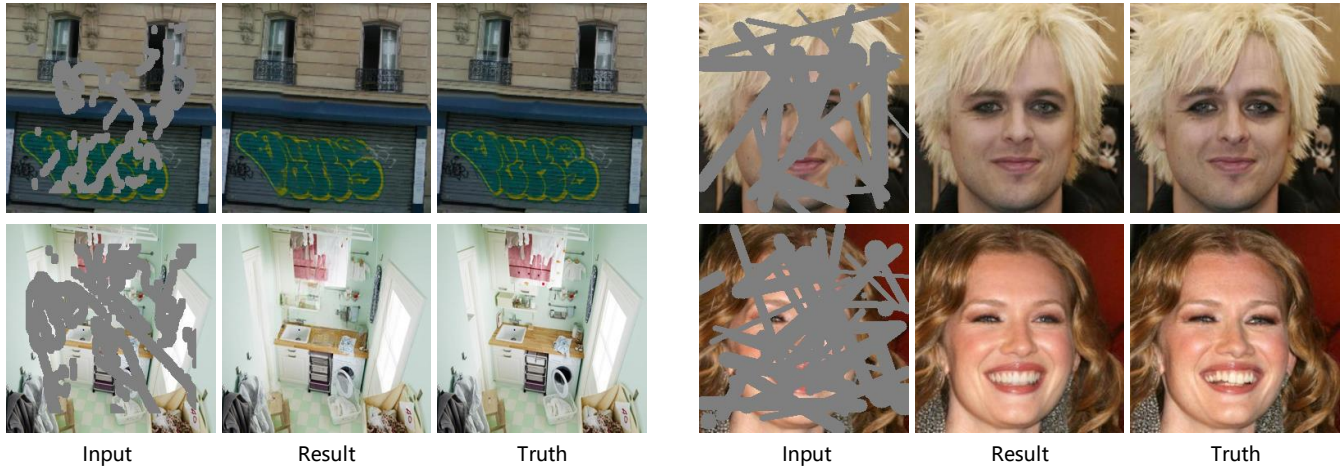
Ye Deng  
Xi'an Jiaotong University  
Xi'an, Shaanxi, China  
dengye@stu.xjtu.edu.cn

Siqi Hui  
Xi'an Jiaotong University  
Xi'an, Shaanxi, China  
huisiqi@stu.xjtu.edu.cn

Sanping Zhou\*  
Xi'an Jiaotong University  
Xi'an, Shaanxi, China  
spzhou@xjtu.edu.cn

Deyu Meng  
Xi'an Jiaotong University  
Xi'an, Shaanxi, China  
dymeng@mail.xjtu.edu.cn

Jinjun Wang  
Xi'an Jiaotong University  
Xi'an, Shaanxi, China  
jinjun@mail.xjtu.edu.cn



**Figure 1: Image inpainting outputs by our proposed  $T$ -former. In each group, the input image is shown on the left, with gray pixels representing the missing areas. (Best with color and zoomed-in view)**

## ABSTRACT

Benefiting from powerful convolutional neural networks (CNNs), learning-based image inpainting methods have made significant breakthroughs over the years. However, some nature of CNNs (e.g. local prior, spatially shared parameters) limit the performance in the face of broken images with diverse and complex forms. Recently, a class of attention-based network architectures, called transformer, has shown significant performance on natural language processing fields and high-level vision tasks. Compared with CNNs, attention operators are better at long-range modeling and have dynamic weights, but their computational complexity is quadratic

in spatial resolution, and thus less suitable for applications involving higher resolution images, such as image inpainting. In this paper, we design a novel attention linearly related to the resolution according to Taylor expansion. And based on this attention, a network called  $T$ -former is designed for image inpainting. Experiments on several benchmark datasets demonstrate that our proposed method achieves state-of-the-art accuracy while maintaining a relatively low number of parameters and computational complexity. The code can be found at [github.com/dengyecode/T-former\\_image\\_inpainting](https://github.com/dengyecode/T-former_image_inpainting)

## CCS CONCEPTS

• **Computing methodologies** → **Computer vision.**

## KEYWORDS

image inpainting, attention, neural networks, transformer

## ACM Reference Format:

Ye Deng, Siqi Hui, Sanping Zhou, Deyu Meng, and Jinjun Wang. 2022.  $T$ -former: An Efficient Transformer for Image Inpainting. In *Proceedings of the 30th ACM International Conference on Multimedia (MM '22)*, October 10–14, 2022, Lisboa, Portugal. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3503161.3548446>

\*The author is also with the Shunan Academy of Artificial Intelligence, Ningbo, Zhejiang 315000, China.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

MM '22, October 10–14, 2022, Lisboa, Portugal

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9203-7/22/10...\$15.00

<https://doi.org/10.1145/3503161.3548446>

## 1 INTRODUCTION

Image inpainting (or completion) [4] is the process of filling in corrupted or missing parts of an image, as Figure 1 shown. It is an important task in the field of computer vision and image processing and can benefit users in a wide range of applications, such as removing unwanted objects in image editing. The key challenge in image inpainting is to make the filled pixels blend in with the non-missing parts.

Prior to deep learning, non-learning inpainting algorithms can be roughly divided into two categories, diffusion-based approaches [1, 3, 4, 9] and exemplar-based approaches [2, 5, 28, 58]. Diffusion-based approaches smoothly propagate the information from observed boundaries to the interior of damaged areas. However, since diffusion-based approaches do not consider the global image structure, they are only effective in filling small holes and less effective in dealing with large scale breakage. To address the drawbacks of diffusion-based inpainting approaches, the exemplar-based approach searches for valid information from known regions of the entire image and copies or relocates this information to the missing locations. Although the exemplar-based algorithms perform well in the face of simple pattern breakage in larger areas, they do not perform well in filling images with complex patterns because they do not understand the semantic information of the image.

With the development of convolutional neural networks (CNNs), learning-based approaches have reached state-of-arts in the field of image inpainting. These inpainting models [24, 30, 42, 66] formulate the inpainting as a conditional image generation problem and customize a CNNs-based encoder-decoder as their corresponding conditional image generator. By training on sufficiently large scale datasets, CNNs show their strength in learning rich patterns and image semantics, and filling the target regions with such learned knowledge. In addition, the sparse connectivity and parameter sharing of CNNs in space make them computationally efficient. However, some basic characteristics of CNNs make them may have some limitations in handling the inpainting task. (a) the locality. CNNs are good at acquiring local relationships but not good at capturing long-range dependencies. Although the importance of locality for images has been demonstrated in various vision tasks for a long time, for image inpainting, focusing on non-local features (the whole image) is more likely to find the appropriate information for broken regions. (b) spatial-sharing and static parameters. The same convolution kernel operates on features across all spatial locations and the parameters of the kernel are static at the time of inference. This is somewhat inflexible in the face of inpainting tasks where images are mixed with broken and unbroken pixels and the damaged regions are variable.

Recently, (self-)attention [51], popular in the field of natural language processing, has been introduced in vision tasks [6, 13]. Compared to CNNs, attention operators whose weights dynamically adjust with the input are better able to capture long-range dependencies through explicit interaction with global features. And as a well-explored architecture in language tasks, the transformer model, based on the attention, is emerging in high-level vision tasks. Although the attention operator has advantages over CNNs in some aspects, its computational complexity grows quadratically with spatial resolution and is therefore not particularly suitable

for high-resolution images, a situation that occurs frequently in low-level vision tasks including image inpainting. Recently, some designs that can reduce the computational complexity of attention operators have been transferred to inpainting [12] or other low-level vision tasks [31, 56]. These methods either apply attention to a sequence of patches unfolded from the image [12], or divide the image into non-overlapping parts and compute the attention for each part independently [31, 56]. However, limiting the spatial extent of attention somewhat defeats the original purpose of capturing the long-range dependence between pixels.

Specifically, for the computational load caused by the dot product and softmax operator in the attention operator, we utilize Taylor's formula to approximate exponential function, and then reduce the computational complexity by swapping the computational order of matrix multiplication. In addition, to mitigate the performance loss caused by the error in the Taylor approximation, we introduced the gating mechanism [10] for the proposed attention operator. The previous work [61] showed that the gating mechanism on the convolution in the inpainting can be seen as controlling which features should flow forward. The gating mechanism we impose on the attention is equivalent to an adjustment of the "inaccurate" attention, allowing the subsequent layers in the network to focus on the information that will help the inpainting, thus producing a high quality complementation result.

In this paper, based on our designed linear attention, we propose an U-net [46] style network, called *T-former*, for image inpainting. Compared with the convolution-based encoder-decoder, in *T-former* we replace the convolution with the designed transformer module based on the proposed linear attention. Our proposed *T-former* combines the texture pattern learning capability of CNNs with the ability of the attention to capture long-range dependencies, and the complexity of this attention is linear rather than quadratically related to the resolution. Our proposed *T-former* is able to achieve performance comparable to other advanced models while maintaining a small complexity compared to those models.

## 2 RELATED WORK

### 2.1 Vision Transformer

The transformer model [51] is a neural network centered on the (self-)attention that plays an important role in natural language processing, and Carion *et al.* [6] were the first to introduce it into the field of vision for object detection. Dosovitskiy *et al.* [13] then designed a transformer structure more suitable for use in the visual field based on the characteristics of images. Touvron *et al.* [49] reduced the data requirements of visual transformer with the help of knowledge distillation. Wang *et al.* [55] then introduced the feature pyramid idea commonly used to build CNNs networks into transformer network construction, which improved the performance of transformer networks. Next, Vaswani *et al.* [50] reduce the computational demand of the model by limiting the range of attention so that the self-attention acts only on a local window. Subsequently, Liu *et al.* [37] extended the use and performance of the transformer model by more subtle design of window attention. These works demonstrated the potential of the transformer for high-level vision tasks, yet because its core self-attention excels in features such as long-range modeling, it also meets the needs

of low-level tasks such as inpainting. However, the computational complexity of the attention in the transformer grows quadratically with spatial resolution, making it inappropriate for direct use in low-vision tasks that require the generation of higher-resolution outputs. Therefore, a class of models chooses to process only low-resolution features of the image with transformer. VQGAN [15], an autoregressive transformer is utilized to learn the effective code-book. ImageBART [14] improves the quality of image synthesis by replacing the autoregressive model in VQGAN with the diffusion process model. MaskGIT [7], in contrast to VQGAN, abandons the autoregressive generation paradigm and introduces a mask, which determines the inference token by the probability value of the mask instead of synthesizing it sequentially as in autoregressive. ICT[53] is a two-stage inpainting where the first stage gets a coarse result by transformer and then feeds this result into a CNN to refine the details. BAT [62] improves on the first stage of ICT by introducing bidirectional and autoregressive transformer to improve the capability of the model. TFill [67] introduces a restricted CNN head on the transformer in ICT to mitigate the proximity influence. These approaches allowed the models to obtain more compact image encodings, but still did not change the limitation that they could not be applied to high resolution images. Subsequently, a different strategy to reduce the complexity was generally adopted. For example, Zamir *et al.* [63] propose replacing spatial attention with inter-channel attention. Or the replacement of inter-pixel attention with inter-patch attention as in [12, 64]. There is also the use of the window attention as in [31, 56] to reduce computational complexity by directly limiting the spatial range of action of attention in a similar way to [37]. Our *T*-former, which does not avoid the problem of attention between full-space pixels, learns long-range dependencies without imposing excessive complexity.

## 2.2 Image Inpainting

Prior to deep learning, non-learning methods could only fill pixels based on the content of the missing regions around [1, 3, 4, 9] or all observed regions [2, 5, 28, 58] because they could not understand the semantics of the image. These methods tend to be more effective for small missing holes or simple background filling, and have limited effect in the face of images with complex patterns. In order to enable the model to output semantic results, Pathak *et al.* [42] introduced the generative adversarial network (GAN) [17] framework to train a conditional image generation model with the help of convolutional neural networks (CNNs). Then, in response to the shared, static parameters of the convolution, some researchers have modified the convolution so that it can manually [34] or automatically [57, 61] adjust the features according to the image breakage. Next, since it is not easy for the model to recover complex patterns directly, some researchers have chosen to guide the model to complete the image with the help of additional extra image information (e.g., edges [40], structure [18, 29, 35, 45], semantics [32, 33]). To improve this, the researchers designed a class of attention operators called contextual attention [36, 54, 59, 60, 65]. Specifically, with the help of the attention module, they explicitly search the entire image for appropriate content to fill the missing regions. Nonetheless, the high burden of performing attention limits its large-scale deployment in the network, so the model is limited in the extent to which

it can improve its long-range modeling capabilities as well as its complementary quality. In contrast, our proposed linear attention in *T*-former is not only able to model long-range dependencies between features, but also reduces the complexity compared to the vanilla attention. This enables us to deploy more attention operators in the proposed *T*-former and achieve state-of-the-art in image inpainting.

## 3 APPROACH

The goal of image inpainting is to fill the target area of the input image  $I_m \in \mathbb{R}^{C \times H \times W}$  with the appropriate pixels so that the image looks intact. To achieve this goal, we designed an U-net [46] style network, based on our proposed linear attention module. In this section, we present our approach from bottom to top. We first describe our proposed linear attention module, and then introduce the architecture of our inpainting network.

### 3.1 Linear Attention

*Vanilla Attention.* We first explain why the attention operator of the vanilla transformer [51] model is not applicable to images with higher resolution. Considering a feature map  $X \in \mathbb{R}^{C \times H \times W}$ , assuming  $N = H \cdot W$ , the attention operator first feeds the feature  $X$  through three different transformations and reshapes them into the two-dimensional matrix to obtain the corresponding three embeddings: the query  $Q = [q_1, q_2, \dots, q_N]^T \in \mathbb{R}^{N \times C}$ , the key  $K = [k_1, k_2, \dots, k_N]^T \in \mathbb{R}^{N \times C}$ , and the value  $V = [v_1, v_2, \dots, v_N]^T \in \mathbb{R}^{N \times C}$ . Then the corresponding attention result  $O \in \mathbb{R}^{N \times C}$  can be obtained by:

$$\begin{aligned} O &= [o_1, o_2, \dots, o_3]^T \\ &= \mathcal{A}(X) \\ &= \text{Softmax}\left(\frac{QK^T}{\sqrt{C}}\right)V \end{aligned} \quad (1)$$

where  $\mathcal{A}(\cdot)$  the attention function which has quadratic space and time complexity with respect to  $N$ . And each  $o_i \in \mathbb{R}^C$  can be obtained as:

$$o_i = \sum_{j=1}^N \frac{\exp\left(\frac{q_i k_j^T}{\sqrt{C}}\right)}{\sum_{l=1}^N \exp\left(\frac{q_i k_l^T}{\sqrt{C}}\right)} v_j \quad (2)$$

The above equation is the dot-product attention with softmax normalization. We can find that the complexity of computing each row  $o_i$  in  $O$  is  $\mathcal{O}(NC)$ . Therefore, the computational complexity of  $O$  is obtained as  $\mathcal{O}(N^2C) = \mathcal{O}((HW)^2C)$ , which is quadratic with respect to the image resolution  $HW$ .

*Linearization of Attention.* We can notice that the computational complexity of Eq. (2) mainly comes from the softmax term, therefore most linearizations of attention focus mainly on modifications to softmax. Revisiting Eq. (2), previous methods [27, 43, 44] compute the attention by using different kernel functions  $K(q, k)$  instead of  $\exp(qk^T)$ , by:

$$\begin{aligned} o_i &= \sum_j \frac{K(q_i, k_j)}{\sum_l K(q_i, k_l)} v_j = \sum_j \frac{f(q_i) f(k_j)^T v_j}{\sum_l f(q_i) f(k_l)^T} \\ &= \frac{f(q_i) \sum_j (f(k_j)^T v_j)}{\sum_j f(q_i) f(k_j)^T} \end{aligned} \quad (3)$$

Note that the property of kernel function  $K(q, k) = f(q)f(k)^\top$  is used here, and  $f(\cdot)$  is a projection. These methods obtain linear attention ( $O(HWC^2)$ ) by changing the order of computation of matrix multiplication from  $qk^\top v$  to  $q(k^\top v)$ .

Inspired by the above linear attention approaches, in this paper we take another perspective to linearize the attention by approximating the exponential function through Taylor expansion. Specifically, we note that Taylor's formula of the exponential function constituting the softmax operator is:

$$\exp(x) = e^x \approx 1 + x \quad (4)$$

Putting Eq. (4) into Eq. (2), we can get (the channel  $C$  is ignored for simplicity):

$$\begin{aligned} o_i &= \sum_{j=1}^N \frac{\exp(q_i k_j^\top)}{\sum_{l=1}^N \exp(q_i k_l^\top)} v_j \\ &= \sum_{j=1}^N \frac{1 + q_i k_j^\top}{\sum_{l=1}^N (1 + q_i k_l^\top)} v_j \\ &= \sum_{j=1}^N \frac{v_j + q_i k_j^\top v_j}{n + q_i \sum_{l=1}^N k_l^\top} \\ &= \sum_{j=1}^N \frac{v_j + q_i (k_j^\top v_j)}{n + q_i \sum_{l=1}^N k_l^\top} \end{aligned} \quad (5)$$

It is worth noting that the last line in Eq. 5 is obtained by the properties of vector multiplication.

*Analysis.* From the above, a linear complexity version of attention can be obtained from the properties of matrix multiplication:

$$V + (QK^\top)V = V + Q(K^\top V) \quad (6)$$

In Eq. (6), instead of calculating the attention matrix  $A = QK^\top \in \mathbb{R}^{N \times N}$  first,  $K^\top V \in \mathbb{R}^{C \times C}$  is computed first and then multiplying  $Q \in \mathbb{R}^{N \times C}$ . With the help of this trick, the computational complexity of the attention operation is  $O(NC^2) = O(HWC^2)$ . It is noted that in the task of image inpainting, the feature (channel) dimension  $C$  is always much smaller than the spatial resolution  $H \times W$ , so we reduce the computational complexity of the model by a large amount. Also similar to the vanilla transformer [51], we also use a multi-headed [51] version of attention to enhance the performance of our proposed linear attention operator. Furthermore, the term  $V+$  seems to be seen as a residual term with respect to  $QK^\top V$ , and from the ablation experiments (as seen in Table 3) we find that it improves the performance of our inpainting model.

### 3.2 Gated Mechanism for Linear Attention

The gating mechanism from recurrent neural networks (GRU [8], LSTM [22]) initially proved its effectiveness on language models [10]. And the gating mechanism is also widely used in the feed-forward networks (FFN) of the state-of-arts transformer networks [23, 47, 63]. A gating mechanism, (or gated linear unit) can be thought of as a neural network layer whose output  $O$  is the product of the components of two linear transformations of the input  $X$ , as:

$$O = I \odot G \quad I = \phi_i(W_u X) \quad G = \phi_g(W_g X) \quad (7)$$

where  $W_i, W_g$  are the learnable parameters,  $\phi_i, \phi_g$  are the corresponding activation functions (which can be absent), and  $\odot$  denotes the Hadamard product. The simple and effective gating mechanism significantly enhances the performance of the network making us want to generalize it to the proposed linear attention operator. Specifically, for an input feature  $X$  and the linear attention operator  $\mathcal{A}(\cdot)$ , then the out  $O$  of the attention with a gating mechanism can be written as:

$$O = A \odot G \quad A = \mathcal{A}(X) \quad G = \phi_g(W_g X) \quad (8)$$

The gating mechanism applied on the convolution [61] plays an important role in the field of image inpainting and can be seen as distinguishing invalid features caused by broken pixels in the input image. Since our proposed linear attention is an "inaccurate" attention, we complement our linear attention with a gating mechanism that allows subsequent layers in the network to focus on features that contribute to the inpainting.

### 3.3 Network Architecture

Our  $T$ -former is an U-net [46] style network based on the proposed transformer block, containing both encoder-decoder parts, as shown in Figure 2. The design of this transformer block we refer to the encoder block in vanilla transformer [51] and contains two sub-layers. The first is the proposed linear attention with gating mechanism (LAG), and the second is a simple feed-forward network (FFN). In addition, we adopt a residual connection [19] adhering to each of the sub-layers. Besides these transformer modules, we also use some convolutional layers to cope with scale changes (like upsampling) of features (inputs).

*Encoder Pipeline.* Given a masked images  $I_m \in \mathbb{R}^{3 \times H \times W}$ , our encoder part first feeds it into a  $7 \times 7$  convolution layer and then get the corresponding feature map  $E_0 \in \mathbb{R}^{C \times H \times W}$ . Here  $H$  and  $W$  represent the dimension of the spatial resolution and  $C$  denotes the dimension of the channel. Next these features are fed into 4-level encoder stages. Each level stage contains a stack of the designed transformer blocks, and we use a convolution with kernel size  $3 \times 3$  and stride 2 to downsample the features between every two stages. For the given feature map  $E_0 \in \mathbb{R}^{C \times H \times W}$ , the  $i$ -level encoder stage of the transformer block produces the feature map  $E_i \in \mathbb{R}^{2^{i-1}C \times \frac{H}{2^{i-1}} \times \frac{W}{2^{i-1}}}$ . By the way the feature map output by the final (4-level) stage of the encoder is  $E_4 \in \mathbb{R}^{8C \times \frac{H}{8} \times \frac{W}{8}}$ .

*Decoder Pipeline.* The decoder takes the final feature map of encoder  $E_4$  as input and progressively restores the high resolution representations. The decoder part consists of 3-level (arranged from largest to smallest) stages, each of which is stacked by several transformer blocks. In each stage of the decoder, the features are first passed through an upsampling layer consisting of nearest neighbor interpolation and  $3 \times 3$  convolution. Given a feature map, the  $i$ -level decoder stage of the upsampling produces the feature map  $D_i \in \mathbb{R}^{2^{i-1}C \times \frac{H}{2^{i-1}} \times \frac{W}{2^{i-1}}}$ . In addition, to help the decoder part, the encoder feature  $E_i$  are concatenated to the decoder feature  $D_i$  via a skip connection. And the concatenation operation is followed by a  $1 \times 1$  convolution layer to decrease the channels (by half). These fused features are then fed into the corresponding transformer block to obtain the dimensionally invariant features  $\tilde{D}_i$ . Finally, after last

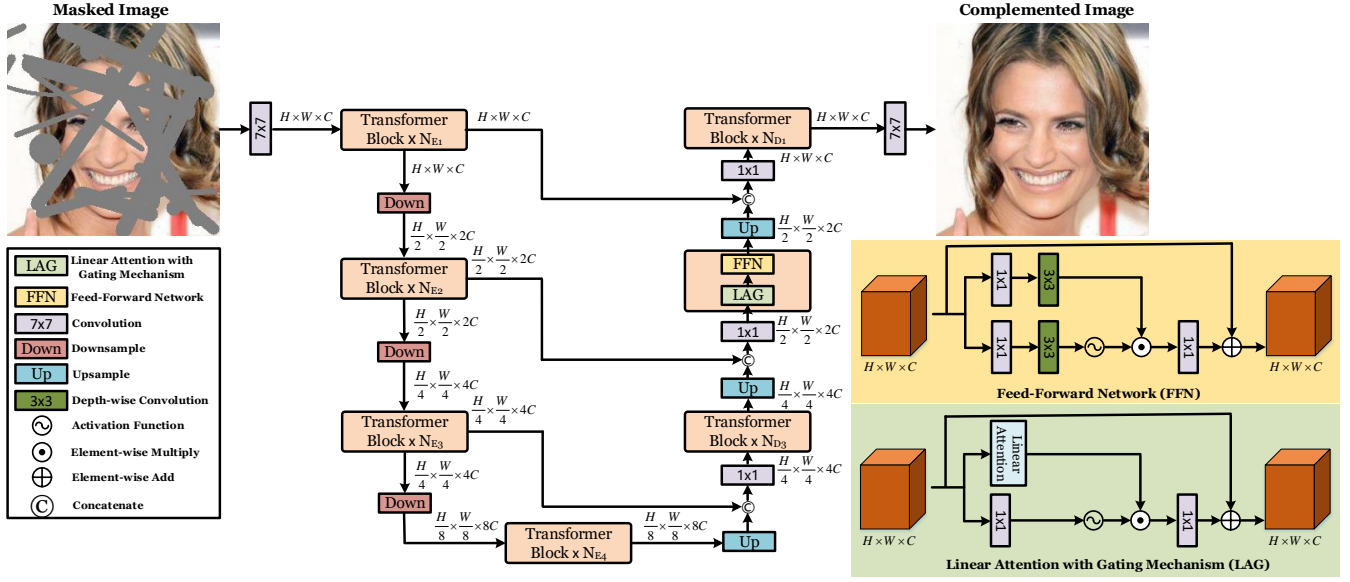


Figure 2: Overview of our proposed *T*-former. Our model accepts masked images as input and outputs complemented images. Our *T*-former which is an U-net style network composed of transformer blocks that we designed. The transformer block we designed contains two sublayers: (1) Linear attention with gating mechanism (LAG) that performs our proposed linear attention for full-space feature interaction, supplemented with a gating mechanism; (2) Feed-forward network (FFN) that transforms the features learned by the attention operator to send useful representations for subsequent layers.

(1-level) stage of the decoder we add a  $7 \times 7$  convolution layer to convert the feature map  $\bar{D}_1 \in \mathbb{R}^{C \times H \times W}$  into the complemented image  $I_{out} \in \mathbb{R}^{3 \times H \times W}$ .

**Transformer Block.** As shown in the Figure 2, the transformer block we used in *T*-former contains two sub-layers. The first is the proposed linear attention with the gating mechanism (LAG), and the second is a simple feed-forward network (FFN). In the LAG layer, the gating value we obtain by feeding the input  $X$  into a  $1 \times 1$  convolution layer with a GELU [20] activation function, i.e.  $\phi_g(\cdot)$  and  $W_g$  of Eq. (8).

For the design of the FFN we refer to the recent transformers [23, 63], which uses a gate-linear layer [10] with residual connections [19] instead of a residual block [19] composed of two convolutions in series. Specifically, to reduce the complexity, for the input  $X \in \mathbb{R}^{C \times H \times W}$ , whose parameter  $W_i, W_j$  in Eq. (7) we replace the standard convolution with a combination of a  $1 \times 1$  convolution and a  $3 \times 3$  depth-wise convolution.

### 3.4 Loss Function

The loss function  $L$  used to train our *T*-former can be written as:

$$L = \lambda_r L_{re} + \lambda_p L_{perc} + \lambda_s L_{style} + \lambda_a L_{adv} \quad (9)$$

where  $L_{re}$  represents the reconstruction Loss,  $L_{perc}$  denotes the perceptual loss [25],  $L_{style}$  denotes the style loss [16] and  $L_{adv}$  is the adversarial loss [17]. And we set  $\lambda_r = 1$ ,  $\lambda_p = 1$ ,  $\lambda_s = 250$ , and  $\lambda_a = 0.1$ . We will describe each loss function in detail below

**Reconstruction Loss.** The reconstruction loss  $L_{re}$  refers to the  $L_1$ -distance between the output  $I_{out}$  and the ground truth  $I_g$ , which

can be defined as:

$$L_{re} = \|I_{out} - I_g\|_1 \quad (10)$$

**Perceptual Loss.** The perceptual loss  $L_{perc}$  is formulated with:

$$L_{perc} = \mathbb{E} \left[ \sum_i \frac{1}{N_i} \|\phi_i(I_{out}) - \phi_i(I_g)\|_1 \right] \quad (11)$$

where  $\phi_i$  is the activation function of the  $i$ -th layer of the VGG-19 [48] pre-trained on ImageNet [11].

**Style Loss.** If the size of feature maps is  $C_j \times H_j \times W_j$ , then the style loss  $L_{style}$  is calculated by:

$$L_{style} = \mathbb{E}_j \left[ \left\| G_j^\Phi(I_{out}) - G_j^\Phi(I_g) \right\|_1 \right] \quad (12)$$

Where  $G_j^\Phi$  denotes a  $C_j \times C_j$  Gram matrix constructed by the corresponding activation maps  $\phi_j$ .

**Adversarial Loss.** The adversarial loss  $L_{adv}$  is formulated with:

$$L_{adv} = \mathbb{E}_{I_g} [\log D(I_g)] + \mathbb{E}_{I_{out}} [\log [1 - D(I_{out})]] \quad (13)$$

where  $D$  represents a patch GAN discriminator [69] with the spectral normalization [39].

## 4 EXPERIMENTS

We evaluated our proposed *T*-former on three datasets, including Paris street view (Paris) [42], CelebA-HQ [26] and Places2 [68]. For CelebA-HQ, we use the first 2000 images for test and the rest for training. For Paris and Places2, we follow the training, testing, and validation splits themselves. During the experiments, all images in datasets were resized to  $256 \times 256$ . Furthermore, during the

**Table 1: Numerical comparisons on the several datasets. The ↓ indicates lower is better, while ↑ indicates higher is better**

DataSet		Paris Street View				Celeba-HQ				Places2			
Mask Ratio		10-20%	20-30%	30-40%	40-50%	10-20%	20-30%	30-40%	40-50%	10-20%	20-30%	30-40%	40-50%
FID↓	GC	20.68	39.48	58.66	82.51	2.54	4.49	6.54	9.83	18.91	30.97	45.26	61.16
	RFR	20.33	28.93	39.84	49.96	3.17	4.01	4.89	6.11	17.88	22.94	30.68	38.69
	CTN	18.08	24.04	36.31	48.46	1.77	3.33	5.24	7.69	15.70	26.41	40.05	55.41
	DTS	16.66	31.94	47.30	65.44	2.08	3.86	6.06	8.58	15.72	27.88	42.44	57.78
	Ours	<b>12.15</b>	<b>22.63</b>	<b>34.47</b>	<b>46.60</b>	<b>1.40</b>	<b>2.55</b>	<b>3.88</b>	<b>5.42</b>	<b>10.85</b>	<b>17.96</b>	<b>26.56</b>	<b>34.52</b>
PSNR↑	GC	32.28	29.12	26.93	24.80	32.25	29.10	26.71	24.78	28.55	25.22	22.97	21.24
	RFR	30.18	27.76	25.99	24.25	30.93	28.94	27.11	25.47	27.26	24.83	22.75	21.11
	CTN	31.22	28.62	26.62	24.91	32.84	29.75	27.35	25.41	27.83	24.91	22.83	21.18
	DTS	32.69	29.28	26.89	24.97	32.91	29.51	27.02	25.13	28.91	25.36	22.94	21.21
	Ours	<b>32.79</b>	<b>29.72</b>	<b>27.47</b>	<b>25.47</b>	<b>33.36</b>	<b>30.15</b>	<b>27.67</b>	<b>25.67</b>	<b>29.06</b>	<b>25.69</b>	<b>23.36</b>	<b>21.52</b>
SSIM↑	GC	0.960	0.925	0.872	0.800	0.979	0.959	0.931	0.896	0.944	0.891	0.824	0.742
	RFR	0.943	0.908	0.861	0.799	0.970	0.958	0.939	0.913	0.929	0.891	0.830	0.756
	CTN	0.955	0.921	0.872	0.812	0.981	0.964	0.940	0.909	0.942	0.892	0.827	0.746
	DTS	0.963	0.929	0.875	0.812	0.981	0.962	0.937	0.905	0.952	0.901	0.834	0.755
	Ours	<b>0.964</b>	<b>0.933</b>	<b>0.887</b>	<b>0.825</b>	<b>0.983</b>	<b>0.967</b>	<b>0.945</b>	<b>0.915</b>	<b>0.953</b>	<b>0.907</b>	<b>0.846</b>	<b>0.770</b>

**Table 2: Complexity measure of different models. Including multiply-accumulate operation count (MAC) and number of parameters (Params). Compared to other baseline models, our *T*-former has a smaller number of parameters and computational complexity**

Model	GC	RFR	CTN	DTS	Ours
MAC	103.1G	206.1G	133.4G	75.9G	51.3G
Params	16.0M	30.6M	21.3M	52.1M	14.8M

experiments in image inpainting we have to specify the location of the broken areas. Therefore, we use the mask dataset from the PC [34] to simulate the location of the corruption. The *T*-former we propose was based on a Pytorch [41] implementation and was trained on one RTX3090 (24 GB) with a batch size of 6. From input to output, the number of transformer blocks of different levels is 1, 2, 3, 4, 3, 2, 1 in order. We used the AdamW [38] optimizer to train the model with a learning rate of  $10^{-4}$  and then fine-tune the model with a learning rate of  $10^{-5}$ . Specifically, on the CelebA-HQ and Paris street view we trained 450,000 iterations and then fine-tuned 200,000 iterations. As for the Places2 data set, we trained about 1 million iterations and then fine-tuned 500,000 iterations.

*Baselines.* To demonstrate the effectiveness of our *T*-former, We compare with the following baselines for their state-of-the-art performance:

- GC [61]: a CNNs-based inpainting model that exploits the gating mechanism and the contextual attention [60] to get high-quality complementary images..
- RFR [30]: a recurrent inpainting method with a special contextual attention that recurrently recovers the missing and progressively strengthens the result.

- CTN [12]: a transformer-style model for image inpainting relies on a patch-based version of the attention operator to model long-range dependencies between features.
- DTS [18]: a dual U-net inpainting model based CNNs, which recovers corrupted images by simultaneous modeling structure-constrained texture synthesis and texture-guided structure reconstruction

*Quantitative Comparison.* Following previous inpainting works [12, 18], we chosen FID (Fréchet Inception Distance) [21], PSNR (peak signal-to-noise ratio), SSIM (structural similarity index) to assess our model. And according to the masks with different masking percentage provided by the dataset [34], the performance of different models under different damage degrees (mask ratio) is tested in Table 1. In addition, we show the number of parameters and the computational complexity (multiply-accumulate operations, MAC) for each model in Table 2. SSIM and PSNR are widely used in image quality assessment for restoration tasks, quantifying the pixel and structural similarities between pairs of images. In addition we adopt the FID, a generally used numeric metric in the task of image generation, to evaluate the image distribution between the inpainting results and the original images. As shown in Table 1 and Table 2, benefiting from the long-distance dependency capture capability and the dynamic nature of the parameters brought by the proposed linear attention, our *T*-former, in the face of different scenarios (datasets) and encountering different breakage situations, can give relatively good complemented images with a low number of parameters and computational complexity.

*Qualitative Comparisons.* Figures 3, 4, and 5 show some comparison results between our and the baseline models on the three data sets Paris [42], CelebA-HQ [26] and Places2 [68] respectively. From these results, we can find that GC [61] is able to complete the basic semantic filling, but the filled position in the image is



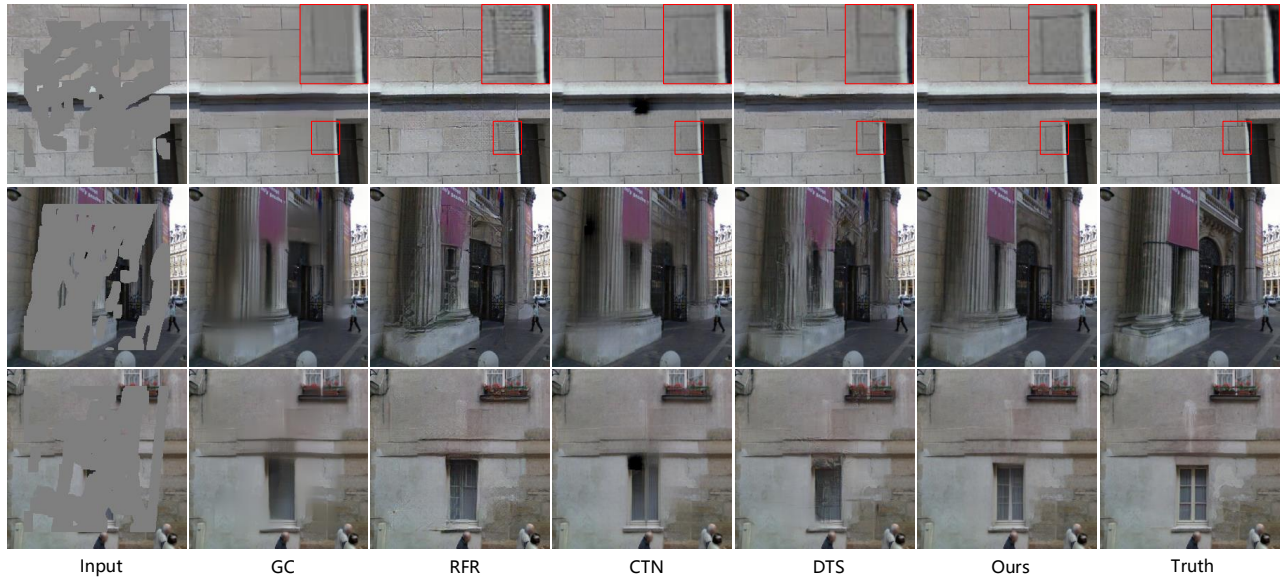


Figure 3: Qualitative results on the Paris with GC [61], RFR [30], CTN [12], DTS [18] and our *T*-former. (Best viewed with zoom-in)

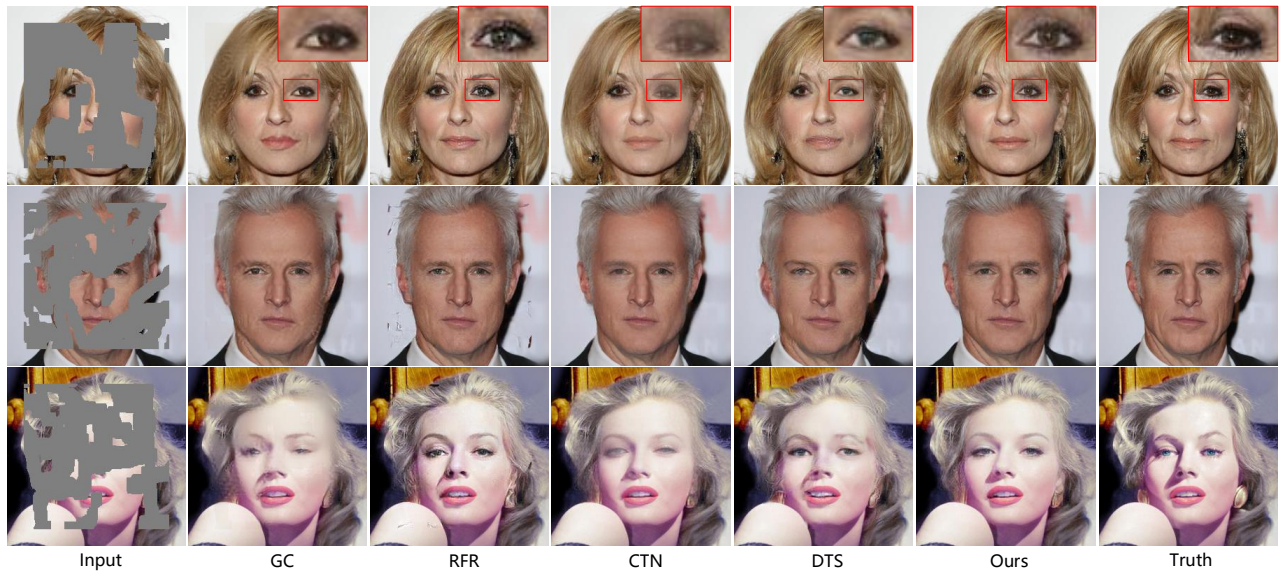


Figure 4: Qualitative results on the CelebA-HQ with GC [61], RFR [30], CTN [12], DTS [18] and our *T*-former. (Best viewed with zoom-in)

prone to blurring, especially when filling images with complex patterns, such as the 2nd row in Figure 3 and 2nd row in Figure 5. The detailed textures of the images complemented by RFR [30] look quite good, but the results are prone to obvious artifacts and are prone to semantic inconsistencies. As in the 1st and 2nd rows of Figure 4, both images generated by RFR show the problem of inconsistent eye color. CTN [12] also performs quite well, but its results are occasionally blurred (Figure 5, line 2) and also prone to

black artifacts as shown in (Figure 5, line 1). DTS [18] performs quite well with simple content images, but when it comes to images with complex patterns, the fill content appears to be significantly disorganized, as shown in the 1st row of Figure 5. Compared to these baselines, in most cases our complemented images look more reasonable and realistic.

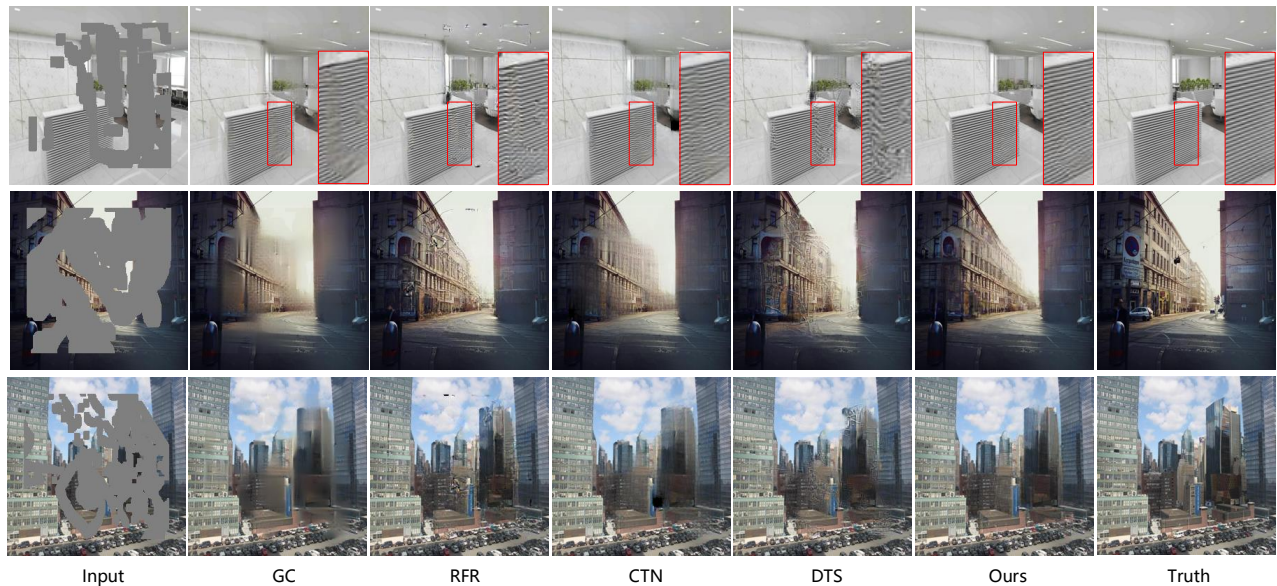


Figure 5: Qualitative results on the Places2 with GC [61], RFR [30], CTN [12], DTS [18] and our  $T$ -former. (Best viewed with zoom-in)

Table 3: Ablation study on the Paris. The  $\downarrow$  indicates lower is better, while  $\uparrow$  indicates higher is better

Mask Ratio		10-20%	20-30%	30-40%	40-50%
FID $\downarrow$	w/o $\mathcal{V}$	12.20	23.93	35.14	47.99
	w/o $\mathcal{G}$	12.74	22.89	36.08	47.48
	w/o $\mathcal{G}+\mathcal{V}$	12.94	24.19	37.08	48.47
	Ours	<b>12.15</b>	<b>22.63</b>	<b>34.47</b>	<b>46.60</b>
PSNR $\uparrow$	w/o $\mathcal{V}$	32.74	29.69	27.35	25.36
	w/o $\mathcal{G}$	32.72	29.68	27.36	25.33
	w/o $\mathcal{G}+\mathcal{V}$	32.70	29.66	27.33	25.28
	Ours	<b>32.79</b>	<b>29.72</b>	<b>27.47</b>	<b>25.47</b>
SSIM $\uparrow$	w/o $\mathcal{V}$	0.962	0.932	0.884	0.823
	w/o $\mathcal{G}$	0.963	0.931	0.884	0.823
	w/o $\mathcal{G}+\mathcal{V}$	0.961	0.930	0.883	0.821
	Ours	<b>0.964</b>	<b>0.933</b>	<b>0.887</b>	<b>0.825</b>

#### 4.1 Ablation Study

We analyze the effectiveness of our proposed module. And all the ablation experiments are conducted on the Paris street view. In the ablation experiments we explored two main components: (1) for the effect of the residual-like connection resulting from  $V+$  in Eq. (6) (or Eq. (5)), i.e.  $1$  in  $e \approx 1+x$  in the Taylor expansion, denoted by  $\mathcal{V}$ . When  $\mathcal{V}$  does not exist, our linear attention is somewhat similar to the current implementation of a series of linear attentions [27, 43, 44] in the field of natural language processing where softmax operators are replaced by kernel functions. And the  $\mathcal{V}$  is equivalent to adding a new residual term to this family of linear attention operators; (2) for the effect of the gating mechanism on the model

performance, denoted by  $\mathcal{G}$ . It can be noticed that both components have a positive impact on the inpainting task. A more interesting point is that it can be found that  $\mathcal{V}$  acts more significantly when the input image is more broken, while the effect of  $\mathcal{G}$  is independent of the degree of input image breakage. The paper [52] has showed that the residual connections can be seen as an ensemble of the model, and one more connection is equivalent to one more sub-network. We speculate that when the model encounters difficult scenes (more broken parts of the input image), more sub-networks (with  $\mathcal{V}$ ) are needed to assist the model get the proper content to fill the missing areas.

## 5 CONCLUSION

In this paper, we propose  $T$ -former, a U-net style network built by the proposed linear attention for image inpainting. To address the problem that CNNs-based inpainting networks have insufficient long-range modeling capability and the standard self-attention operator has high computational load, we propose a linear attention operator based on Taylor’s formula that captures the long-range dependence between features at a small computational cost. In addition, we utilize a gating mechanism to enhance the performance of the proposed linear attentional operator. Quantitative and qualitative results demonstrate that our proposed  $T$ -former outperforms state-of-the-art methods in terms of performance and also maintains a relatively small complexity.

## ACKNOWLEDGMENTS

This work is jointly supported by the National Key Research and Development Program of China under Grant No. 2017YFA0700800, the General Program of China Postdoctoral Science Foundation under Grant No. 2020M683490, and the Youth program of Shaanxi Natural Science Foundation under Grant No. 2021JQ-054.



## REFERENCES

- [1] C. Ballester, M. Bertalmio, V. Caselles, G. Sapiro, and J. Verdera. 2001. Filling-in by joint interpolation of vector fields and gray levels. *IEEE Transactions on Image Processing* 10, 8 (2001), 1200–1211. <https://doi.org/10.1109/83.935036>
- [2] Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. 2009. PatchMatch: A Randomized Correspondence Algorithm for Structural Image Editing. *ACM Transactions on Graphics (Proc. SIGGRAPH)* 28, 3 (Aug. 2009).
- [3] M. Bertalmio. 2006. Strong-continuation, contrast-invariant inpainting with a third-order optimal PDE. *IEEE Transactions on Image Processing* 15, 7 (2006), 1934–1938. <https://doi.org/10.1109/TIP.2006.877067>
- [4] Marcelo Bertalmio, Guillermo Sapiro, Vincent Caselles, and Coloma Ballester. 2000. Image Inpainting. In *Proceedings of the 27th Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH '00)*. ACM Press/Addison-Wesley Publishing Co., USA, 417–424. <https://doi.org/10.1145/344779.344972>
- [5] Pierre Buysens, Maxime Daisy, David Tschumperlé, and Olivier Lézoray. 2015. Exemplar-Based Inpainting: Technical Review and New Heuristics for Better Geometric Reconstructions. *IEEE Transactions on Image Processing* 24, 6 (2015), 1809–1824. <https://doi.org/10.1109/TIP.2015.2411437>
- [6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-End Object Detection with Transformers. In *European Conference on Computer Vision*, Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm (Eds.). Springer International Publishing, 213–229.
- [7] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T. Freeman. 2022. MaskGIT: Masked Generative Image Transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 11315–11325.
- [8] Kyunghyun Cho, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259* (2014).
- [9] Antonio Criminisi, Patrick Pérez, and Kentaro Toyama. 2004. Region filling and object removal by exemplar-based image inpainting. *IEEE Transactions on image processing* 13, 9 (2004), 1200–1212.
- [10] Yann N. Dauphin, Angela Fan, Michael Auli, and David Grangier. 2017. Language Modeling with Gated Convolutional Networks. In *Proceedings of the 34th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 70)*, Doina Precup and Yee Whye Teh (Eds.). PMLR, 933–941.
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*. 248–255.
- [12] Ye Deng, Siqi Hui, Sanping Zhou, Deyu Meng, and Jinjun Wang. 2021. Learning Contextual Transformer Network for Image Inpainting. In *Proceedings of the 29th ACM International Conference on Multimedia (Virtual Event, China) (MM '21)*. Association for Computing Machinery, New York, NY, USA, 2529–2538. <https://doi.org/10.1145/3474085.3475426>
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xi-aohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *International Conference on Learning Representations*.
- [14] Patrick Esser, Robin Rombach, Andreas Blattmann, and Bjorn Ommer. 2021. ImageBART: Bidirectional Context with Multinomial Diffusion for Autoregressive Image Synthesis. In *Advances in Neural Information Processing Systems*, Vol. 34. 3518–3532.
- [15] Patrick Esser, Robin Rombach, and Bjorn Ommer. 2021. Taming Transformers for High-Resolution Image Synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 12873–12883.
- [16] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. 2016. Image Style Transfer Using Convolutional Neural Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [17] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative Adversarial Nets. In *Advances in Neural Information Processing Systems*, Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Q. Weinberger (Eds.), Vol. 27. Curran Associates, Inc.
- [18] Xiefan Guo, Hongyu Yang, and Di Huang. 2021. Image Inpainting via Conditional Texture and Structure Dual Generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 14134–14143.
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [20] Dan Hendrycks and Kevin Gimpel. 2016. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415* (2016).
- [21] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. In *Advances in Neural Information Processing Systems*, Vol. 30.
- [22] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [23] Weizhe Hua, Zihang Dai, Hanxiao Liu, and Quoc V Le. 2022. Transformer Quality in Linear Time. *arXiv preprint arXiv:2202.10447* (2022).
- [24] Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa. 2017. Globally and Locally Consistent Image Completion. *ACM Transactions on Graphics (Proc. of SIGGRAPH 2017)* 36, 4, Article 107 (2017), 107:1–107:14 pages.
- [25] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. 2016. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*. Springer, 694–711.
- [26] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. 2018. Progressive Growing of GANs for Improved Quality, Stability, and Variation. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=Hk99zCeAb>
- [27] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. 2020. Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, 5156–5165.
- [28] Nikos Komodakis and Georgios Tziritas. 2007. Image Completion Using Efficient Belief Propagation Via Priority Scheduling and Dynamic Pruning. *IEEE Transactions on Image Processing* 16, 11 (2007), 2649–2661. <https://doi.org/10.1109/TIP.2007.906269>
- [29] Jingyuan Li, Fengxiang He, Lefei Zhang, Bo Du, and Dacheng Tao. 2019. Progressive Reconstruction of Visual Structure for Image Inpainting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- [30] Jingyuan Li, Ning Wang, Lefei Zhang, Bo Du, and Dacheng Tao. 2020. Recurrent Feature Reasoning for Image Inpainting. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [31] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. 2021. SwinIR: Image Restoration Using Swin Transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*. 1833–1844.
- [32] Liang Liao, Jing Xiao, Zheng Wang, Chia-Wen Lin, and Shin'ichi Satoh. 2020. Guidance and evaluation: Semantic-aware image inpainting for mixed scenes. In *European Conference on Computer Vision*. Springer, 683–700.
- [33] Liang Liao, Jing Xiao, Zheng Wang, Chia-Wen Lin, and Shin'ichi Satoh. 2021. Image Inpainting Guided by Coherence Priors of Semantics and Textures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 6539–6548.
- [34] Guilin Liu, Fitsum A. Reda, Kevin J. Shih, Ting-Chun Wang, Andrew Tao, and Bryan Catanzaro. 2018. Image Inpainting for Irregular Holes Using Partial Convolutions. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- [35] Hongyu Liu, Bin Jiang, Yibing Song, Wei Huang, and Chao Yang. 2020. Rethinking Image Inpainting via a Mutual Encoder-Decoder with Feature Equalizations. In *European Conference on Computer Vision*. Springer International Publishing, Cham, 725–741.
- [36] Hongyu Liu, Bin Jiang, Yi Xiao, and Chao Yang. 2019. Coherent Semantic Attention for Image Inpainting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- [37] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 10012–10022.
- [38] Ilya Loshchilov and Frank Hutter. 2019. Decoupled Weight Decay Regularization. In *International Conference on Learning Representations*.
- [39] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. 2018. Spectral Normalization for Generative Adversarial Networks. In *International Conference on Learning Representations*.
- [40] Kamyar Nazeri, Eric Ng, Tony Joseph, Faisal Qureshi, and Mehran Ebrahimi. 2019. EdgeConnect: Structure Guided Image Inpainting using Edge Prediction. In *The IEEE International Conference on Computer Vision (ICCV) Workshops*.
- [41] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc.
- [42] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A. Efros. 2016. Context Encoders: Feature Learning by Inpainting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [43] Hao Peng, Nikolaos Pappas, Dani Yogatama, Roy Schwartz, Noah Smith, and Lingpeng Kong. 2021. Random Feature Attention. In *International Conference on Learning Representations*.
- [44] Zhen Qin, Weixuan Sun, Hui Deng, Dongxu Li, Yunshen Wei, Baohong Lv, Junjie Yan, Lingpeng Kong, and Yiran Zhong. 2022. cosFormer: Rethinking Softmax In

- Attention. In *International Conference on Learning Representations*.
- [45] Yurui Ren, Xiaoming Yu, Ruonan Zhang, Thomas H. Li, Shan Liu, and Ge Li. 2019. StructureFlow: Image Inpainting via Structure-Aware Appearance Flow. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- [46] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-Net: Convolutional Networks for Biomedical Image Segmentation. In *Medical Image Computing and Computer-Assisted Intervention*, Nassir Navab, Joachim Hornegger, William M. Wells, and Alejandro F. Frangi (Eds.). Springer International Publishing, Cham, 234–241.
- [47] Noam Shazeer. 2020. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202* (2020).
- [48] Karen Simonyan and Andrew Zisserman. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *International Conference on Learning Representations*.
- [49] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve Jegou. 2021. Training data-efficient image transformers & distillation through attention. In *Proceedings of the 38th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 10347–10357.
- [50] Ashish Vaswani, Prajit Ramachandran, Aravind Srinivas, Niki Parmar, Blake Hechtman, and Jonathon Shlens. 2021. Scaling Local Self-Attention for Parameter Efficient Visual Backbones. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 12894–12904.
- [51] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*, Vol. 30. Curran Associates, Inc.
- [52] Andreas Veit, Michael J Wilber, and Serge Belongie. 2016. Residual Networks Behave Like Ensembles of Relatively Shallow Networks. In *Advances in Neural Information Processing Systems*, D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (Eds.), Vol. 29. Curran Associates, Inc.
- [53] Ziyu Wan, Jingbo Zhang, Dongdong Chen, and Jing Liao. 2021. High-Fidelity Pluralistic Image Completion With Transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 4692–4701.
- [54] Ning Wang, Jingyuan Li, Lefei Zhang, and Bo Du. 2019. MUSICAL: Multi-Scale Image Contextual Attention Learning for Inpainting. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*. International Joint Conferences on Artificial Intelligence Organization, 3748–3754. <https://doi.org/10.24963/ijcai.2019/520>
- [55] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. 2021. Pyramid Vision Transformer: A Versatile Backbone for Dense Prediction Without Convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 568–578.
- [56] Zhendong Wang, Xiaodong Cun, Jianmin Bao, and Jianzhuang Liu. 2021. Uformer: A general u-shaped transformer for image restoration. *arXiv preprint arXiv:2106.03106* (2021).
- [57] Chaohao Xie, Shaohui Liu, Chao Li, Ming-Ming Cheng, Wangmeng Zuo, Xiao Liu, Shilei Wen, and Errui Ding. 2019. Image Inpainting With Learnable Bidirectional Attention Maps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- [58] Zongben Xu and Jian Sun. 2010. Image Inpainting by Patch Propagation Using Patch Sparsity. *IEEE Transactions on Image Processing* 19, 5 (2010), 1153–1165. <https://doi.org/10.1109/TIP.2010.2042098>
- [59] Zhaoyi Yan, Xiaoming Li, Mu Li, Wangmeng Zuo, and Shiguang Shan. 2018. Shift-Net: Image Inpainting via Deep Feature Rearrangement. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- [60] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S. Huang. 2018. Generative Image Inpainting With Contextual Attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [61] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S. Huang. 2019. Free-Form Image Inpainting With Gated Convolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- [62] Yingchen Yu, Fangneng Zhan, Rongliang Wu, Jianxiang Pan, Kaiwen Cui, Shijian Lu, Feiyang Ma, Xuansong Xie, and Chunyan Miao. 2021. Diverse image inpainting with bidirectional and autoregressive transformers. In *Proceedings of the 29th ACM International Conference on Multimedia*. 69–78.
- [63] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. 2021. Restormer: Efficient Transformer for High-Resolution Image Restoration. *arXiv preprint arXiv:2111.09881* (2021).
- [64] Yanhong Zeng, Jianlong Fu, and Hongyang Chao. 2020. Learning Joint Spatial-Temporal Transformations for Video Inpainting. In *European Conference on Computer Vision*, Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm (Eds.). Springer International Publishing.
- [65] Yanhong Zeng, Jianlong Fu, Hongyang Chao, and Baining Guo. 2019. Learning Pyramid-Context Encoder Network for High-Quality Image Inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [66] Chuanxia Zheng, Tat-Jen Cham, and Jianfei Cai. 2019. Pluralistic Image Completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [67] Chuanxia Zheng, Tat-Jen Cham, Jianfei Cai, and Dinh Phung. 2022. Bridging Global Context Interactions for High-Fidelity Image Completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 11512–11522.
- [68] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. 2018. Places: A 10 Million Image Database for Scene Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 40, 6 (2018), 1452–1464. <https://doi.org/10.1109/TPAMI.2017.2723009>
- [69] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. 2017. Unpaired Image-To-Image Translation Using Cycle-Consistent Adversarial Networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.