

DGFont++: Robust Deformable Generative Networks for Unsupervised Font Generation

Xinyuan Chen, Yangchen Xie, Li Sun, Yue Lu,

Abstract—Automatic font generation without human experts is a practical and significant problem, especially for some languages that consist of a large number of characters. Existing methods for font generation are often in supervised learning. They require a large number of paired data, which are labor-intensive and expensive to collect. In contrast, common unsupervised image-to-image translation methods are not applicable to font generation, as they often define style as the set of textures and colors. In this work, we propose a robust deformable generative network for unsupervised font generation (abbreviated as DGFont++). We introduce a feature deformation skip connection (FDSC) to learn local patterns and geometric transformations between fonts. The FDSC predicts pairs of displacement maps and employs the predicted maps to apply deformable convolution to the low-level content feature maps. The outputs of FDSC are fed into a mixer to generate final results. Moreover, we introduce contrastive self-supervised learning to learn a robust style representation for fonts by understanding the similarity and dissimilarities of fonts. To distinguish different styles, we train our model with a multi-task discriminator, which ensures that each style can be discriminated independently. In addition to adversarial loss, another two reconstruction losses are adopted to constrain the domain-invariant characteristics between generated images and content images. Taking advantage of FDSC and the adopted loss functions, our model is able to maintain spatial information and generates high-quality character images in an unsupervised manner. Experiments demonstrate that our model is able to generate character images of higher quality than state-of-the-art methods.

Index Terms—Font Generation, Unsupervised Image-to-image Translation, Image-to-image Translation, and Image Generation.

I. INTRODUCTION

Every day, people consume a massive amount of text for information transfer and storage. As the representation of texts, the font is closely related to our daily life. Font generation is critical in many applications, *e.g.*, font library creation, personalized handwriting, historical handwriting imitation, and data augmentation for optical character recognition and handwriting identification. Traditional font library creating methods heavily rely on expert designers by drawing each glyph individually, which is especially expensive and labor-intensive for logographic languages such as Chinese (more than 60,000 characters), Japanese (more than 50,000 characters), and Korean (11,172 characters).

Recently, the development of convolutional neural networks enables automatic font generation without human experts. There have been some attempts to explore font generation and achieve promising results. [1]–[3] utilize deep neural networks to generate entire sets of letters for certain alphabet languages. Two notable projects, “Rewrite” [4] and “zi2zi” [5], generate logographic language characters by learning a mapping from one style to another with thousands of paired characters. After that, EMD [6] and SA-VAE [7] design neural networks to separate the content and style representation, which can extend to generate the character of new styles or contents. However, these methods are in supervised learning and required a large amount of paired training samples.

Some other methods exploit auxiliary annotations (*e.g.*, strokes, radicals) to facilitate high-quality font generation. For example, [8] utilizes labels for each stroke to generate glyphs by writing trajectories synthesis. [9] employ the radical decomposition (*e.g.*, radicals or sub-glyphs) of characters to achieve font generation for certain logographic language. DM-Font [10] and its improved version LF-Font [11] propose disentanglement strategies to disentangle complex glyph structures, which help capture local details in rich text design. MXFont [12] extracts localized features for few-shot font generation by exploiting sub-glyph and components of characters. However, these methods rely on prior knowledge and can only apply to specific writing systems. Some labels such as the stroke skeleton can be estimated by algorithms, but the estimation error would decrease the generated quality. Also, these methods still require thousands of paired data and annotated labels for training. Recently, there are some attempts [13], [14] for unsupervised font generation. [14] introduces a novel module that transfers the features across sequential DenseNet blocks [15]. [13] proposes a fast skeleton extraction method to obtain the skeleton of characters, and then utilize the extracted skeleton to facilitate font generation.

In the field of image-to-image translation, a series of unsupervised generative models have been proposed by combining adversarial training [16], [17] with carefully designed constraints [18]–[21]. However, classical unsupervised image-to-image translation methods cannot be directly applied to unsupervised font generation tasks. In image-to-image translation, the style of images is usually defined as the set of textures and colors. In contrast, different fonts have their own local patterns such as stroke thickness, tips of brushes, and joined-up writing patterns. Existing methods for image-to-image translation usually extract the style feature of the target class images and employ adaptive instance normalization (AdaIN)

X. Chen is with Shanghai Artificial Intelligence Laboratory, Shanghai, China (e-mail:xychen9191@gmail.com).

Y. Xie, L. Sun, and Y. Lu are with Shanghai Key Laboratory of Multidimensional Information Processing, East China Normal University, Shanghai, 200241, China (e-mail:yexie0702@126.com;sunli@ee.ecnu.edu.cn;ylu@cs.ecnu.edu.cn)

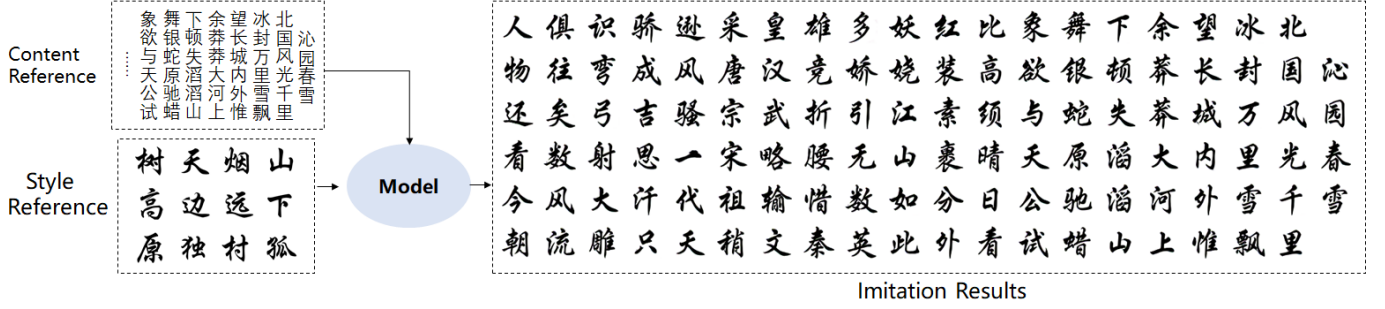


Fig. 1: **Unsupervised font generation results.** Given content images and a few style reference images, our model aims to generate imitations. In this example, the reference images are of a calligraphic font, and the imitation result is generated from our model.

[22], [23] to combine the content and the style features. The AdaIN-based methods transfer style by aligning feature statistics, which tends to transform texture and color, which is not suitable to transform local style patterns (e.g., geometric deformation) for the font.

Compelled by the above observations, we propose a robust deformable generative model for unsupervised font generation (DGFont++). The proposed method is designed to deform and transform the character of one font to another by leveraging the provided images of the target font. Figure 1 shows our scenario for font generation. Given a few style reference images (e.g., calligraphy from an artist), our model is able to generate imitations by transforming content reference character images. Our proposed DGFont++ separates style and content respectively and then mixes two representations to generate target characters. We introduce a feature deformation skip connection (FDSC) to predict pairs of displacement maps and employ the predicted maps to apply deformable convolution to the low-level feature maps from the content encoder. The outputs of FDSC are then fed into a mixer to generate the final results. Also, to learn a robust representation for fonts, we use contrastive learning for our style encoder to understand the similarity and dissimilarities of fonts. To this end, we define several data augmentation operations to construct positive counterparts for fonts. The model is then imposed to learn a feature space where characters and their positive counterparts are at a similar point. To distinguish different styles, we train our model with a multi-task discriminator, which ensures that each style can be discriminated independently. In addition to adversarial loss, another two reconstruction losses are adopted to constrain the domain-invariant characteristics between generated images and content images.

The feature deformation skip connection (FDSC) module is used to transform the low-level feature of content images, which preserves the pattern of character (e.g., strokes and radicals). Different from the image-to-image translation problem that defines style as a set of textures and colors, the style of font is basically defined as geometric transformation, stroke thickness, tips of brushes, and joined-up writing pattern. Two fonts with the same content, usually have correspondence for each stroke. Taking advantage of the spatial relationship of fonts, the feature deformation skip connection (FDSC) is used

to conduct spatial deformation, which effectively ensures the generated image has complete structures. Moreover, to further improve the quality of generated images, we integrate local spatial attention into an FDSC module, which predicts the local relationship between the encoder and mixer features. The local spatial attention model predicts similarity scores for features of each position regarding its neighboring positions.

This work is an extensive version of our conference paper [24]. Compared to the previous version, this paper includes the following additional contributions. (1) A more robust style feature extraction approach is proposed. We introduce data augmentation operations for fonts to construct positive counterparts of characters and introduce contrastive loss to help the model to learn better representation. (2) We integrate local spatial attention into an FDSC module, *i.e.*, FDSC-attn. Experiments prove the superiority of FDSC-attn compared to Vanilla FDSC. (3) Additional experiments are conducted to ablate and analyze the function of our DGFont++. (4) We conduct more comprehensive experiments of different image sizes and comparisons to state-of-the-art methods. Extensive experiments demonstrate that our model outperforms state-of-the-art font generation methods. Besides, results show that our model is able to extend to generate unseen character.

II. RELATED WORK

A. Font Generation

Font generation aims to automatically generate characters in a specific font and create a font library. Recent studies have employed image translation methods for font generation. “Zi2zi” [5] and “Rewrite” [4] implement font generation on the basis of GAN [25] with thousands of character pairs for strong supervision. After that, a series of models are proposed to improve the generated quality based on zi2zi [5]. PEGAN [26] sets up a multi-scale image pyramid to pass information through refinement connections. HAN [27] improves zi2zi by designing a hierarchical loss and skip connection. AEGG [28] adds an additional network to refine the training process. DC-Font [29] introduces a style classifier to get a better style representation. However, all the above methods are in supervised learning and require a large number of paired data. In unsupervised font generation, [13], [14]

achieve unsupervised font generation by learning a mapping between two fonts directly. However, they ignore the geometric deformation of the font, and their results are not satisfying.

Lots of methods employ auxiliary annotations (*e.g.*, stroke and radical decomposition) to further improve the generation quality. SA-VAE [7] disentangles the style and content as two irrelevant domains with encoding Chinese characters into high-frequency character structure configurations and radicals. CalliGAN [30] further decomposes characters into components and offers low-level structure information including the order of strokes to guide the generation process. RD-GAN [9] proposes a radical extraction module to extract rough radicals which can improve the performance of the discriminator and achieves the few-shot Chinese font generation. DM-Font [10] and its improved version LF-Font [11] propose disentanglement strategies to disentangle complex glyph structures, which help capture local details in rich text design. MXFont [12] extracts localized features for few-shot font generation by exploiting sub-glyph and components of characters. Some other attempts have been made in Chinese character generation by adopting skeleton/stroke extraction algorithm [8], [13]. However, they need extra annotations or algorithms to guide font generation; while the estimation error would decrease the generation performance. In this work, our model, DGFont++, aims to generate high-quality character images in an unsupervised way without other annotations.

B. Image-to-Image Translation

The purpose of image-to-image translation is to learn a mapping from an image in the source domain to the target domain. Image-to-image translation has been applied in many fields such as artistic style transfer [31], [32], semantic segmentation [33], [34], image animation [35]–[37], object transfiguration [38], and video frames generation [39]–[41] *et al.* Pix2pix [42] is the first model proposed for image-to-image translation based on conditional GAN [43]. To achieve unsupervised image-to-image translation, a lot of works [18], [44]–[46] have been proposed, where Cycle-GAN [18] introduces a cycle consistency between source and target domain to discover the relationship of samples between two domain. However, the above-mentioned methods can only translate from one domain to another specific domain. To tackle this problem, recent works [21], [23], [47], [48] are proposed to simultaneously generate multiple style outputs given the same input. Gated-GAN [47] proposes a gated transformer to transfer multiple styles in a single model. FUNIT [23] encodes content image and class image respectively, and combines them with AdaIN [22]. TUNIT [21] further introduces a guiding network as an unsupervised domain classifier to automatically produce a domain label of a given image. DUNIT [48] extracts separate representations for the global image and for the instances to preserve the detailed content of object instances. To learn the mapping across geometry variations, [49] introduces a discriminator with dilated convolutions as well as a multi-scale perceptual loss that can represent errors in the underlying shape of objects. [50] disentangles image space into a Cartesian product of the appearance and the geometry latent spaces.

Our task is related to unsupervised image translation [21], [23], [51], which shares the same aim to translate images from one domain to another based on reference images of the target domain. However, these classical unsupervised image-to-image translation methods usually focus on images of wild animals or style transfer whose style is defined as poses or a set of textures and colors. There exist some works for improving shape deformation in the unsupervised image-to-image translation, such as GAN-imorph [49], but they still cannot generate high-quality font images. In contrast, characteristics of font lie in local patterns such as stroke thickness, tips of brushes, geometric deformation, and joined-up writing patterns. The unique characteristics of the font motivate the design of robust deformable networks for font generation.

C. Deformable Convolution and Attention Mechanism

CNNs have inherent limitations in modeling geometric transformations due to the fixed kernel configuration. To enhance the transformation modeling capability of CNNs, [52] proposes the deformable convolutional layer. It augments the spatial sampling locations in the modules with additional offsets. The deformable convolution has been applied to address several high-level vision tasks, such as object detection [52]–[54] video object detection [55] sampling, semantic segmentation [54], and human pose estimation [56]. Recently, some methods attempt to apply deformable convolution in image generation tasks. TDAN [57] addresses video super-resolution tasks by using deformable convolution to align two continuous frames and output a high-resolution frame. [58], [59] synthesized novel view images by deformable convolution given the view condition vectors. In our proposed DGFont, offsets are estimated by a learned latent style code.

Furthermore, attention modules and deformable convolution can be viewed from a unified perspective [60]. The proper combination of deformable convolution and attention module achieves higher accuracy than using one of them in the tasks of object detection and semantic segmentation. [59] proposes Soft and Hard Conditional Deformation Modules which employ deformable convolution and attention in a special skip connection way in supervised image synthesis tasks. They show that attention is more effective when applied to lower-resolution features, which is also demonstrated by [61], [62]. Our proposed model shows that the combination of deformable convolution and spatial attention works well in unsupervised font generation tasks.

D. Unsupervised Representation Learning

In this study, we utilize unsupervised learning to learn a robust style feature representation. Unsupervised representation learning aims to extract meaningful representations for downstream tasks without human supervision. Recently, a family of contrastive learning methods based on maximizing mutual information has been proposed [63]–[69]. These methods make use of noise contrastive estimation [70], learning an embedding where associated samples are brought together, in contrast to other samples in the dataset. The design choices of the contrastive loss, such as the number of negatives and how

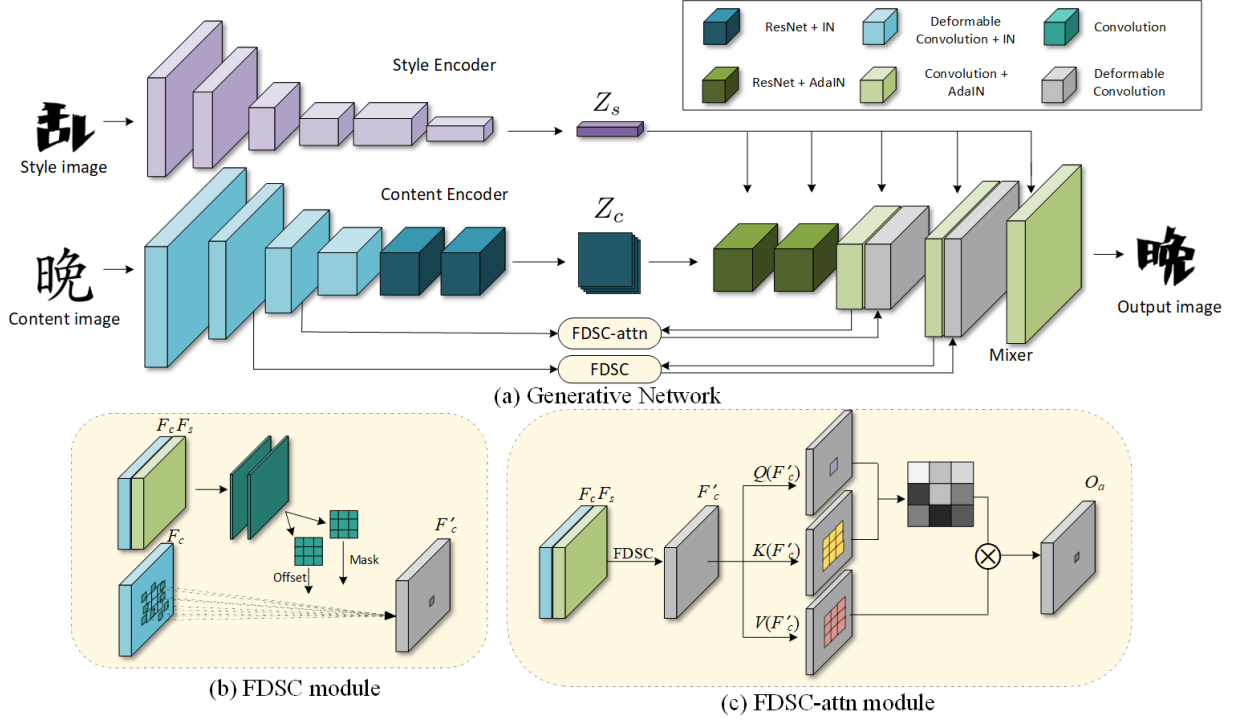


Fig. 2: **Overview of the proposed method.** a) Overview of our generative network. The style/content encoder maps style/content image to style/content representation Z_s/Z_c . FDSC and FDSC-attn apply transformation convolution to the low-level feature from the content encoder and inject the results into the mixer. The mixer generates the output image. b) A detailed illustration of the FDSC module. c) A detailed illustration of the FDSC-attn module.

to sample them, and data augmentation all play a critical role and need to be carefully studied. By maintaining the dictionary as a queue of negative data samples and employing the data augmentation images as positive samples, Moco [65] achieves outstanding performance in various downstream tasks under a reasonable mini-batch size on the ImageNet dataset. Due to the gap between natural images and font images, the data augmentation used in Moco cannot be directly applied to our framework. We adopt contrastive learning in our framework and design several data augmentations for font images.

III. METHODS

A. Overview

Given a content image I_c and a style image I_s , our model aims to generate the character of the content image with the font of the style image. As illustrated in Figure 2, the proposed generative network consists of a style encoder, a content encoder, a mixer, and two feature deformation skip connection (FDSC) modules. The architecture of the style encoder and discriminator is simplified in Figure 2. The detailed architecture is in Section IV-A1. **The style encoder** is designed to learn the style representation from input images. The style encoder takes a style image as the input and maps it to a style latent vector Z_s . **The content encoder** is introduced to extract the structure feature of the content images. The content encoder maps the content image into a spatial feature map Z_c . The content encoder module is made of three deformable convolution layers followed by two residual

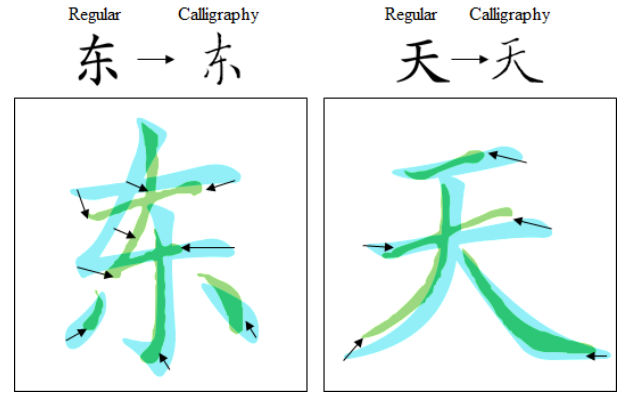


Fig. 3: **The geometric deformation of two fonts for a character.** There is a correspondence for each stroke between two fonts for the same character.

blocks. The introduced deformable convolution layer enables the content encoder to produce style-invariant features for images with the same content. **The mixer** aims to output characters by mixing the content feature representations Z_c and style feature representations Z_s . AdaIN [22] is adapted to inject the style feature into the mixer. Besides, the **feature deformation skip connection (FDSC)** modules transfer the deformed low-level feature from the content encoder to the mixer. Details are described in Sec III-B.

B. Feature Deformation Skip Connection

As illustrated in Figure 3, there lies a geometric deformation of two fonts for a character and exists a correspondence for each stroke. Compelling this observation, we propose a feature deformation skip connection (FDSC) module to apply geometric deformation convolution to the content image in the feature space and directly transfer the deformation low-level feature to the mixer. Specifically, the module predicts offsets based on the guidance code to instruct the deformable convolution layer to perform a geometric transformation on the low-level feature. As demonstrated in Figure 2, the input of the FDSC module is a concatenation of two feature maps: a feature map F_c extracted from the content image and a style guidance map F_s . F_s is extracted from the mixer after injecting the style code F_s . The module estimates sampling parameters after applying convolution to the concatenation of F_s and F_c :

$$\Theta = f_{\theta}(F_s, F_c). \quad (1)$$

Here, f_{θ} refers to a deformable convolution layer, and $\Theta = \{\Delta p_k, \Delta m_k \mid k = 1, \dots, |\mathcal{R}|\}$ refers to the offsets and mask of the convolution kernel, where $\mathcal{R} = \{(-1, -1), (-1, 0), \dots, (0, 1), (1, 1)\}$ indicates a regular grid of a 3×3 kernel. Under the guidance of sampling parameter Θ , a geometrically deformed feature map F'_c is obtained from Θ and F_c based on deformable convolution $f_{DC}(\cdot)$:

$$F'_c = f_{DC}(F_c, \Theta). \quad (2)$$

Specifically, for each position p on the output F'_c , the deformable convolution $f_{DC}(\cdot)$ is applied as follow:

$$F'_c(p) = \sum_{k=1}^{\mathcal{R}} w(p_k) \cdot x(p + p_k + \Delta p_k) \cdot \Delta m_k, \quad (3)$$

where the $w(p_k)$ indicates the weight of the deformable convolution kernel at k -th location. The convolution is operated on the irregular positions $(p_k + \Delta p_k)$ where Δp_k may be fractional. Followed [52], the operation is implemented by using bilinear interpolation. At last, the output of the FDSC module is fed to the mixer and F'_c is then concatenated with F_s like a commonly used skip connection [71].

Deformable convolution introduces 2D offsets to the regular grid sampling locations in the standard convolution. It enables free-form deformation of the sampling grid. There are lots of areas of the same color in character images, such as background color and character color. By using deformable convolution, an area can be related to any other area with the same color. It is difficult to optimize the non-unique solution. To efficiently use our FDSC module, we impose a constraint on the offsets Δp . We introduce the constraint in detail in Subsection III-D, and demonstrate the visualization of the offsets Δp in Section IV-D2.

Our FDSC module aims to deform the spatial structure of the content image in the feature space. It is crucial to select which level of features to be transformed, as low-level features contain more structure and spatial information than high-level features. Experiments in Section IV-D demonstrate the analysis of the performance of the model with different numbers of the FDSC module.

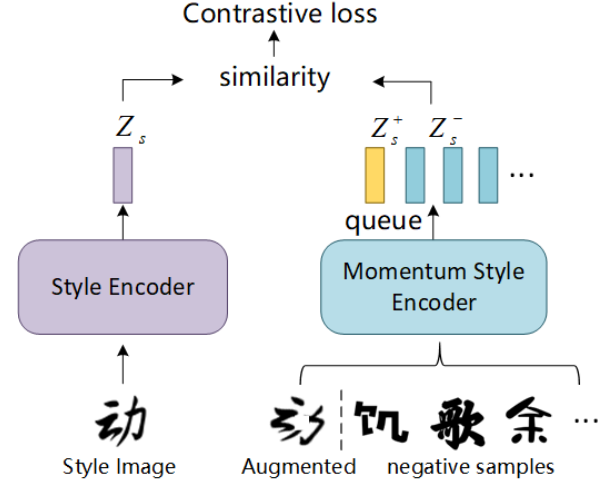


Fig. 4: Illustration of our contrastive representation learning. The model employs the style reference image and its augmented image as a positive pair. The other image samples of the dataset are negative samples.

FDSC-attn. In the inner skip connection, we integrate the local spatial attention model, aiming to explore local relationships by predicting weights at a position in regard to the features at its neighboring positions. The detail of FDSC-attn is shown in Figure 2(c). The deformable feature $F'_c \in \mathbb{R}^{h \times w \times c}$ produced by FDSC is first projected to the latent space as query, key, value $Q(F'_c), K(F'_c), V(F'_c) \in \mathbb{R}^{h \times w \times c}$ by using 1×1 convolution layers. For each location (i, j) within the spatial dimensions, we extract a patch with size s from K centered at (i, j) , denoted as $\mathbf{k} \in \mathbb{R}^{s \times s \times c}$. Then the weight $\mathbf{w} \in \mathbb{R}^{s \times s \times c}$ is obtained by reshaping the estimation of a feed-forward network (FFN):

$$\mathbf{w} = \text{reshape}(\text{FFN}(\text{concat}(\text{flatten}(\mathbf{k}), \mathbf{q}))). \quad (4)$$

The FFN is computed by giving the concatenation of the flattened patch query $\text{flatten}(\mathbf{k})$ and its corresponding query vector $\mathbf{q} \in \mathbb{R}^{1 \times 1 \times c}$ at (i, j) . The FFN consists of a fully-connected (FC) layer followed by leaky ReLU and a linear FC layer. Here w plays a conceptually equivalent role as the softmax attention map of the traditional key query aggregation [83, 98, 65]. The output vector $\mathbf{o} \in \mathbb{R}^{1 \times 1 \times c}$ at (i, j) is then calculated by element-wise multiplication of the predicted w and the value patch with the same size from V centered at (i, j) , denoted as $\mathbf{v} \in \mathbb{R}^{s \times s \times c}$: $\mathbf{o}(i, j) = \mathbf{w} \odot \mathbf{v}$. We loop over all the (i, j) to constitute an output O_a as the attention output. Different from deformable convolution, the additional attention module provides a heat map between the feature of a certain location and its context patch. It allows our model to look over the local attention information, then generates a 'soft' deformed feature according to its current position feature and context.

C. Unsupervised Feature Representation for Style Encoder

To learn a robust representation for a given font, we hope to learn consistent style features regardless of its content.

To this end, we define several data augmentation operations to construct positive counterparts for fonts. Then we adopt contrastive learning to force the positive pairs to be similar in the latent space. The contrastive representation learns the visual styles of images by maximizing the mutual information between an image and its augmented version in contrast to other negative images within the dataset.

A detailed illustration of the proposed contrastive representation learning is shown in Figure 4. Our style encoder aims to learn a robust map of the style image I_s to style code Z_s . Similar to MOCO [65], a momentum style encoder is used to map the positive counterpart I_s^+ as positive codes Z_s^+ . We construct a queue of data samples that are stored from the previously sampled images and mapped them as negative codes Z_s^{i-} . Then we use a contrastive loss to unitize the similarity of the positive pair (Z_s, Z_s^+) and the dissimilarity of negative pairs (Z_s, Z_s^{i-}) .

$$\mathcal{L}_{sty} = -\log \frac{\exp(Z_s \cdot Z_s^+ / \tau)}{\exp(Z_s \cdot Z_s^+ / \tau) + \sum_{i=1}^N \exp(Z_s \cdot Z_s^{i-} / \tau)}, \quad (5)$$

where τ is a temperature hyper-parameter [72]. The sum is over one positive and N negative sample. This loss is the log loss of an $(N+1)$ -way softmax-based classifier aiming to classify Z_s as Z_s^+ . The momentum style encoder is updated followed by [65]. The contrastive representation facilitates the style encoder to learn robust style features, which alleviates the impact of using different content images of a certain style. Experiments in Section IV-C2 demonstrate the influence by varying the number of the reference style images.

D. Loss Function

Our model aims to achieve automatic font generation via an unsupervised method. In addition to the contrastive loss in Eq. 5, we adopt four losses: 1) adversarial loss is used to produce realistic images; 2) content consistent loss is introduced to encourage the content of the generated image to be consistent with the content image; 3) image reconstruction loss is used to maintain the domain-invariant features; 4) deformation offset normalization is designed to prevent excessive offsets of the FDSC module. We introduce the formula of each loss and the full objective in this section.

Adversarial loss: When character images are generated from the generative network, a **multi-task discriminator** is adopted to conduct discrimination for each style simultaneously. For each style, the output of the discriminator is a binary classification whether the input image is a real image or a generated image. As there are several different styles of fonts in the training set, the discriminator outputs a binary vector whose length is the number of styles. In all, our model aims to generate plausible images by solving a mini-max optimization problem. The generative network G tries to fool discriminator D by generating fake images. The adversarial loss penalty is the wrong judgment when real/generated images are input to the discriminator.

$$\mathcal{L}_{adv} = \max_{D_s} \min_G \mathbb{E}_{I_s \in P_s, I_c \in P_c} [\log D_s(I_s) + \log(1 - D_s(G(I_s, I_c)))], \quad (6)$$

where $D_s(\cdot)$ denotes the logit from the corresponding style of the discriminator's output.

Content consistent loss: adversarial loss is adopted to help the model to generate a realistic style while ignoring the correctness of the content. To prevent mode collapse and ensure that the features extracted from the same content can be content consistent after the content encoder f_c , we impose a content consistent loss here:

$$\mathcal{L}_{cnt} = \mathbb{E}_{I_s \in P_s, I_c \in P_c} \|Z_c - f_c(G(I_s, I_c))\|_1. \quad (7)$$

\mathcal{L}_{cnt} ensures that given a source content image I_c and corresponding generated images, their feature maps are consistent after content encoder f_c .

Image reconstruction loss: To ensure that the generator can reconstruct the source image I_c when given its origin style, we impose a reconstruction loss:

$$\mathcal{L}_{img} = \mathbb{E}_{I_c \in P_c} \|I_c - G(I_c, I_c)\|_1. \quad (8)$$

The objective helps preserve domain-invariant characteristics (e.g., content) of its input image I_c .

Deformation offset normalization: The deformable offsets enable free-form deformation of the sampling grid. As there are lots of areas of the same color between input images and generated images (such as background color and character color), it leads to a non-unique solution that is difficult to optimize. Meanwhile, the font generation focus on the stroke relationship between the content character image and the target character image, such as the thickness and tips of the stroke. However, given images with the same content but different style, the position of the same stroke in these two images are close. To efficiently use this deformable convolutional network, we impose a constrain on the offsets Δp :

$$\mathcal{L}_{offset} = \frac{1}{|\mathcal{R}|} \|\Delta p\|_1, \quad (9)$$

where Δp denotes offsets of the deformable convolution kernel, $|\mathcal{R}|$ denotes the number of the convolution kernel.

Overall objective loss: Combining all the above-mentioned losses, we have the overall loss function for training our proposed framework:

$$\mathcal{L} = \mathcal{L}_{adv} + \lambda_{img} \mathcal{L}_{img} + \lambda_{cnt} \mathcal{L}_{cnt} + \lambda_{sty} \mathcal{L}_{sty} + \lambda_{offset} \sum_i^N \mathcal{L}_{offset} / N, \quad (10)$$

where λ_{adv} , λ_{img} , λ_{cnt} , λ_{sty} , and λ_{offset} are hyper-parameters to control the weight of each loss, and N indicates the number of FDSC modules. In our model, the generative network aims to minimize the overall object loss, while the discriminator aims to maximize it.

IV. EXPERIMENTS

In this section, we evaluate our proposed model for the Chinese font generation task. We first introduce the implementation details of our experiments. Extensive experiments demonstrate the superiority of our model. We also provide ablation studies on the effects of FDSC, FDSC-attn, and objective function.

A. Implementation Detail

1) Network Architecture

The generative network consists of a style encoder, a content encoder, a mixer, and two FDSC modules. The style encoder is composed of 7 convolution layers and a fully connected (FC) layer. The convolution layer is of 3×3 kernel with nonlinear activation of ReLU. The FC layer maps the style feature maps into a style representation vector of 128 dimensions. The architecture of the content encoder and mixer is symmetrical. The detailed architectures of the style encoder, content encoder, and mixer are shown in Table I. The two FDSC modules consist of an FDSC and an FDSC-attn module. The FDSC is constructed by a deformable convolution layer with 3×3 kernel. FDSC-attn consists of a deformable convolution layer and a local spatial attention module with patch size $s = 13 \times 13$. The discriminator contains several residual blocks and three convolution layers. The detailed information on the discriminator is listed in Table II. FRN indicates filter response normalization [73].

2) Training Strategy

We initial the weights of convolutional layers with He initialization [74], in which all biases are set to zero and the weights of linear layers are sampled from $N(0, 0.01)$. We use Adam optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.99$ for style encoder, and RMSprop optimizer with $\alpha = 0.99$ for the content encoder and mixer. We train the whole framework with 200,000 iterations and the learning rate is set to 0.0001 with a weight decay of 0.001. We train the model with a hinge version adversarial loss [75] with R1 regularization [76] using $\gamma = 10$. We optimize our DGFont++ with $\lambda_{img} = 0.1$, $\lambda_{cnt} = 0.1$, $\lambda_{sty} = 0.1$, $\lambda_{offset} = 0.1$. The temperature hyper-parameter in Eq. 5 is set as $\tau = 0.07$ by default [65]. In the testing process, we use ten reference images to compute an average style code for the generation process.

3) Data augmentation

We use several data augmentation operations of spatial and geometric transformation to construct positive counterparts for fonts. Specifically, for a given style reference image, we conduct the following data augmentation operations successively: randomly scaling (from 0.7 to 1.2), rotation (degree from -20 to 20), horizontal and vertical translation (shift value range from 0 to 0.2 of image size) and random erasing with filled 255.

4) Compared Methods

We compare our model with the following methods for Chinese font generation:

- Cycle-GAN [18]: Cycle-GAN is an unsupervised image-to-image translation method, which consists of two generative networks which translate images from one domain to another using a cycle consistency loss.
- EMD [6]: EMD is a supervised font generation method that is optimized by L1 distance loss between ground-truth and generated images. It employs an encoder-decoder architecture and separates style/content representations.
- Zi2zi [5]: Zi2zi is a supervised font generation method based on pix2pix [42], it achieves font generation and

	Operation	Kernel	Resample	Padding	Features
Style encoder	Convolution	3	MaxPool	1	64
	Convolution	3	MaxPool	1	128
	Convolution	3	-	1	256
	Convolution	3	MaxPool	1	256
	Convolution	3	-	1	512
	Convolution	3	MaxPool	1	512
	Convolution	3	-	1	512
	Convolution	3	MaxPool	1	512
Content encoder	Avg pooling	-	-	-	512
	FC	-	-	-	128
	Deform. conv.	3	-	1	32
	Deform. conv.	4	stride-2	1	64
	Deform. conv.	4	stride-2	1	128
	Convolution	4	stride-2	1	256
Mixer	Res block	3	-	1	256
	Res block	3	-	1	256
	Res block	3	-	1	256
	Convolution	5	Upsample	2	128
	Convolution	5	Upsample	2	64
	Convolution	5	Upsample	2	32
	Convolution	3	-	1	3
	Res block	3	-	1	256
	Res block	3	-	1	256
	Convolution	5	Upsample	2	128

TABLE I: Architecture of generative network

Operation	Kernel	Resample	Features	Normalization
Convolution	3	-	64	-
Res block	3	-	64	FRN
Res block	3	AvgPool	128	FRN
Res block	3	-	128	FRN
Res block	3	AvgPool	256	FRN
Res block	3	-	256	FRN
Res block	3	AvgPool	512	FRN
Res block	3	-	512	FRN
Res block	3	AvgPool	1024	FRN
LeakyReLU	-	-	-	-
Convolution	4	-	1024	-
LeakyReLU	-	-	-	-
Convolution	1	AvgPool	221	-

AvgPool: Average Pooling; LeakyReLU: Slope 0.2.

TABLE II: Architecture of discriminative network

uses Gaussian Noise as category embedding to achieve multi-style transfer.

- GANimorph [49]: GANimorph is an unsupervised image-to-image translation method for representing shape deformation by introducing a discriminator with dilated convolutions to get a more context-aware generator.
- FUNIT [23]: FUNIT is a few-shot unsupervised image-to-image translation model which separates content and style of natural animal images and combines them with adaptive instance normalization (AdaIN) layer.
- MXFONT [12]: MXFont is the state-of-the-art few-shot font generation method with multiple localized experts by utilizing decomposable glyphs.

5) Dataset and Evaluation Metrics

We collect a dataset containing 231 fonts (styles) including printing and handwriting fonts, each of which has 1143 commonly used Chinese characters (content). The dataset is randomly partitioned into a training set and a testing set. The training set contains 221 fonts, and each font contains 800 characters. The testing set consists of two parts. One testing part is the remaining 343 characters of the 221 fonts. Another part is the remaining 10 fonts for testing the generalization for

Methods	one-to-many	training	L1 loss↓	RMSE↓	SSIM↑	LPIPS↓	FID↓
Seen fonts							
EMD [6] - 80	✓	paired	0.0538	0.1955	0.7676	0.1036	89.65
Zi2zi [5] - 80	×	paired	0.0521	0.1802	0.7789	0.1065	142.23
Cycle-GAN [18] - 80	×	unpaired	0.0863	0.2555	0.6392	0.1825	175.24
GANimorph [49] - 80	×	unpaired	0.0563	0.1759	0.7808	0.1403	72.89
FUNIT [23] - 80	✓	unpaired	0.0807	0.2510	0.6669	0.1216	53.77
DGFont - 80	✓	unpaired	0.0562	0.1994	0.7580	0.0814	46.15
DGFont++ - 80	✓	unpaired	0.0580	0.2032	0.7502	0.0875	43.56
FUNIT [23] - 128	✓	unpaired	0.0894	0.2766	0.7193	0.2113	36.40
MXFont [12] - 128	✓	paired	0.0704	0.2353	0.7622	0.1505	46.69
DGFont - 128	✓	unpaired	0.0584	0.2166	0.7871	0.1204	41.69
DGFont++ - 128	✓	unpaired	0.0567	0.2131	0.7910	0.1151	35.83
Unseen fonts							
EMD [6] - 80	✓	paired	0.0430	0.1755	0.7849	0.1255	82.53
FUNIT [23] - 80	✓	unpaired	0.0588	0.2089	0.7417	0.1125	59.98
DGFont - 80	✓	unpaired	0.0414	0.1709	0.7982	0.0867	50.29
FUNIT [23] - 128	✓	unpaired	0.0703	0.2426	0.7655	0.1734	35.76
MXFont [12] - 128	✓	paired	0.0561	0.2057	0.8121	0.1269	39.61
DGFont - 128	✓	unpaired	0.0517	0.2032	0.8146	0.1250	41.59
DGFont++ - 128	✓	unpaired	0.0532	0.2059	0.8125	0.1243	34.54

TABLE III: **Quantitative evaluation on the whole dataset.** We evaluate the methods on seen and unseen font sets. The bold number indicates the best. 80 and 128 in the first column indicate the training and testing image size.

unseen fonts. The image size is 128×128 , which is consistent with most mainstream methods [12], [23]. In order to compare with the latest font generation methods (*i.e.*, MXFont [12]), the collected characters can be decomposed into sub-characters. Differently, the dataset in our conference version is of 80×80 resolution and some of the characters in the conference version cannot be decomposed. For a fair comparison, all the experiments on 128×128 images use the newly collected dataset, and experiments on 80×80 images use the dataset of the conference version.

We employ five metrics L1 loss, Root Mean Square Error (RMSE), SSIM, LPIPS [77] and FID [78] for evaluation. L1 loss, RSME, and Structural Similarity (SSIM) is widely used in font generation task with known ground truths [6]. L1 loss calculates the L1-norm of the distance between generated image and ground-truth. RMSE utilizes the mean square error to give an overall evaluation. SSIM calculates the global mean and variance to assess the structural similarity. Meanwhile, Learned Perceptual Image Patch Similarity (LPIPS) is another metric to compute the distance between the generations and ground truths in the perceptual domain. Besides, Fréchet Inception Distance (FID) is employed to measure the realism of the generated images. The FID score is calculated between the distributions of the generated and real data for each style.

B. Comparison with State-of-art Methods

1) Quantitative comparison

The quantitative results are shown in Table III. In addition to comparing methods of generating 80×80 images as our previous version [24], we adapt our model to generate 128×128 and compare with recent work [12] designed for the image of 128×128 . For the one-to-one image translation models (*i.e.*, Cycle-GAN and GANimorph), we train models individually for each target style. In the experiment

of generating images of 80×80 size, results show that our methods (*i.e.*, DGFont and DGFont++) are comparable to the compared methods in pixel-level evaluation metrics, *e.g.*, L1 loss, RMSE, SSIM. It is noted that these metrics focus on the pixel-wise comparison between generated image and ground-truth while ignoring the feature similarity. Pixel-wise metrics tend to predict high scores for blur images that are not consistent with human perception [31]. In perceptual-level metrics (*i.e.*, FID [78] and LPIPS [77]), our methods outperform all the compared methods for both seen font and unseen font.

In the experiment of generating images of 128×128 sizes, we compare our method to FUNIT and MXFont. We observe MXFont achieves a better score than FUNIT in terms of L1 loss, RMSE, SSIM, and LPIPS, while FUNIT outperforms MXFont in terms of FID score. It is because MXFont uses pixel-wise loss to supervise the training of the model. FUNIT has a higher score in FID which measures the realism and quality of generated images. In contrast, our methods (*i.e.*, DGFont and DGFont++) outperform FUNIT and MXFont in all metrics and achieve state-of-the-art performance. Our DGFont++ has achieved better performance than DGFont. Detailed analysis for DGFont and DGFont++ is discussed in Section IV-C2. To explore the generality of our model, we also generate imitations for unseen styles whose results are shown in Table IV. It demonstrates that our model outperforms all the compared methods, including the few-shot font generation method, MXFont.

2) Qualitative comparison

Figure 5 demonstrates the qualitative comparison between our method and the state-of-the-art methods. To explore the capability of deforming and transforming source character patterns (*e.g.*, stroke, skeleton), we display the visual comparisons as two cases: easy cases and challenge cases. Characters

C-GAN:	怀 饭 化 政 形	性 用 那 面 社	实 到 浓 是 全	你 很 国 起 情
EMD:	怀 饭 化 政 形	性 用 那 面 社	实 到 浓 是 全	你 很 国 起 情
Zi2zi:	怀 饭 化 政 形	性 用 那 面 社	实 到 浓 是 全	你 很 国 起 情
GAN-imorph:	怀 饭 化 政 形	性 用 那 面 社	实 到 浓 是 全	你 很 国 起 情
FUNIT:	怀 饭 化 政 形	性 用 那 面 社	实 到 浓 是 全	你 很 国 起 情
DGFont:	怀 饭 化 政 形	性 用 那 面 社	实 到 浓 是 全	你 很 国 起 情
MXFont:	怀 饭 化 政 形	性 用 那 面 社	实 到 浓 是 全	你 很 国 起 情
DGFont++:	怀 饭 化 政 形	性 用 那 面 社	实 到 浓 是 全	你 很 国 起 情
Target:	怀 饭 化 政 形	性 用 那 面 社	实 到 浓 是 全	你 很 国 起 情

(a) Easy Cases (e.g., printing typeface).

C-GAN:	邵 性 家 过 琛	我 全 机 把 羊	第 或 数 好 能	在 有 去 点 上
EMD:	邵 性 家 过 琛	我 全 机 把 羊	第 或 数 好 能	在 有 去 点 上
Zi2zi:	邵 性 家 过 琛	我 全 机 把 羊	第 或 数 好 能	在 有 去 点 上
GAN-imorph:	邵 性 家 过 琛	我 全 机 把 羊	第 或 数 好 能	在 有 去 点 上
FUNIT:	邵 性 家 过 琛	我 全 机 把 羊	第 或 数 好 能	在 有 去 点 上
DGFont:	邵 性 家 过 琛	我 全 机 把 羊	第 或 数 好 能	在 有 去 点 上
MXFont:	邵 性 家 过 琛	我 全 机 把 羊	第 或 数 好 能	在 有 去 点 上
DGFont++:	邵 性 家 过 琛	我 全 机 把 羊	第 或 数 好 能	在 有 去 点 上
Target:	邵 性 家 过 琛	我 全 机 把 羊	第 或 数 好 能	在 有 去 点 上

(b) Challenging Cases (e.g., calligraphy, wordart).

Fig. 5: Comparisons between our model and the state-of-the-art methods. The red boxes highlight failures of structure preservation, the blue boxes highlight failures of style transfer, and the orange boxes highlight the blur and noisy outputs.

in easy cases are close to printing typefaces that do not have cursive writing, while in challenge cases characters are of WordArt and calligraphy fonts that are hollow or have joined-up writing. We observe that the results of Cycle-GAN can hardly recognize and tend to be noisy. In easy cases, characters generated by Zi2zi, EMD, and GANimorph can maintain a complete structure, but they are usually vague. In challenging cases, they can only generate parts of characters or sometimes unreasonable structures. FUNIT can generate characters with a clear background but the generated characters lose their structure to some degree. In the challenging cases (Figure 5(b)), FUNIT generates characters of incomplete structure

when the target font is the cursive writing style. In the case of generating hollow font, FUNIT fails to generate the character of reference style (see the last column in Figure 5(b)). MXFont is the state-of-the-art few-shot font generation method. In most cases, MXFont and our method generate images of comparable quality. Notably, MXFont requires paired datasets while our methods are in an unsupervised manner. Also, we observe that MXFont fails in generating hollow font. In contrast, our proposed methods, DGFont, and DGFont++ can not only generate characters with complete structure but also learn joined-up writing. Also, DGFont++ can produce more clear and more complete images than DGFont. More ablation study

Source	Baseline	(a)	(b)	(c)	(d)	Target
	面	面	面	面	面	面
	多	多	多	多	多	多
	到	到	到	到	到	到
L1 loss:	0.0632	0.0600	0.0595	0.0587	0.0582	
RMSE:	0.2199	0.2126	0.2120	0.2097	0.2080	
SSIM:	0.7304	0.7420	0.7427	0.7460	0.7469	
LPIPS:	0.1108	0.1048	0.1026	0.1022	0.1006	
FID:	64.86	56.58	50.23	48.87	46.39	

Fig. 6: **Effect of different components in DGFont.** We add different parts into our baseline successively. (a) Replace the first three convolution layers of content encoder with deformable convolution layers; (b) add one FDSC module (without normalization); (c) impose normalization on the FDSC module; (d) add another FDSC module (*i.e.*, DGFont).

for DGFont and DGFont++ is discussed in the next section.

C. Ablation Study

In this part, we conduct ablation studies for our model. We set our baseline model as a normal encoder-decoder model which consists of a style encoder, content encoder, and mixer, but no deformable convolution layer, FDSC module, and contrastive loss. In Section IV-C1, we conduct ablation studies for incremental modules and losses between DGFont and the baseline model. In Section IV-C2, we conduct ablation studies for incremental modules and losses between DGFont and DGFont++.

1) Baseline model vs. DGFont

We add deformable convolution, feature deformation skip connection, and deformable offset normalization successively on the baseline model and get the full model of DGFont. The experiments are conducted in challenging cases to explore the functionality of each component for DGFont. Qualitative and quantitative comparisons are shown in Figure 6.

1) Effectiveness of deformable convolution in the content encoder. Figure 6(a) shows the results by replacing the first three convolution layers of the content encoder with deformable convolution layers. We can see that the quantitative results improve obviously in terms of L1 loss, RMSE, and SSIM. This indicates that deformable convolution layers in the content encoder effectively help improve the performance of our model.

2) The influence of the FDSC module. In this part, we add an FDSC module (without offset normalization in Eq. 9) that connects the features after the first layer and penultimate layer. Results are shown in Figure 6(b). Comparing with Figure 6(a), we observe that the generated characters preserve more structure information and are able to reconstruct the complete structure of characters.

3) Effectiveness of deformable offset constraint. We investigate the impact of deformable offset normalization

by comparing FDSC modules without and with offset normalization. As shown in Figure 6(b) and (c), adding offset normalization helps the model generate images whose style becomes more similar to the target.

4) Effectiveness of two FDSC modules. Figure 6(d) shows the results of full DGFont with two FDSC modules. It is noted that the generated images get more details, less noise, and achieve better quantitative results.

2) DGFont vs. DGFont++

Based on DGFont, we then replace FDSC as FDSC-attn in the inner skip connection (indicated as DGFont+attn), and successively add contrastive learning (indicated as DGFont+attn+ \mathcal{L}_{sty}) to get the full model of DGFont++. Results of ablation study for DGFont++ (*i.e.*, DGFont+attn+ \mathcal{L}_{sty}) are shown in Table IV and Figure 7. Experiments are conducted on the image of size 128×128 . We analyze the results and effectiveness as follows.

1) Effectiveness of the FDSC-attn module. In this part, we investigate the effectiveness of FDSC-attn compared with FDSC. Quantitative results are shown in the second row of Table IV (indicated as “DGFont+attn”). DGFont+attn achieves a close score in pixel-wise metrics (*i.e.*, L1 loss, RMSE, SSIM) while outperforming DGFont in terms of perceptual metrics (*i.e.*, FID and LPIPS). The qualitative results for FDSC-attn are shown in Figure 7. Compared to DGFont, the generated images learn the target style well. For example, in the fourth and fifth columns, DGFont fails to transfer style whose results remain the style as the source style, while the model with FDSC-attn succeeds in generating images of the target style.

2) Effectiveness of contrastive loss. We investigate the impact of style contrastive loss. In Table IV, we observe that the model with contrastive loss produces results of the highest score in all metrics. The fourth row of Figure 7 (indicated as “DGFont+attn+ \mathcal{L}_{sty} ”) shows the qualitative results for contrastive loss. Results show that our DGFont++ is able to generate high-quality images without missing strokes or noisy points. It is because the introduced contrastive learning helps to learn robust and content-invariant style representation, which helps to maintain complete structure and strokes of generated glyph images. More analysis for contrastive representation learning is demonstrated in the next section.

3) Robustness Analysis. In this subsection, we investigate the robustness of our model by varying the size of the style reference images. We report the FID score of DGFont++ and DGFont by varying the size of the reference set from 1 (2^0) to 256 (2^8). Results are shown in Figure 8. For each case, we evaluate the performance by randomly sampling the style images ten times, and the results are demonstrated as box plots. We use the same y-axis interval in Figure 8(a) and Figure 8(b) for better comparison. Results show that DGFont++ achieves better FID scores, whose results also have a smaller fluctuation under ten times of inference with the same number of references. Meanwhile, Figure 8(c) shows the variance of FID scores for DGFont++ and DGFont, which demonstrates that DGFont++ has a lower variance than DGFont especially when given a small number of reference images. We observe that the variance of the two models declines with an increase

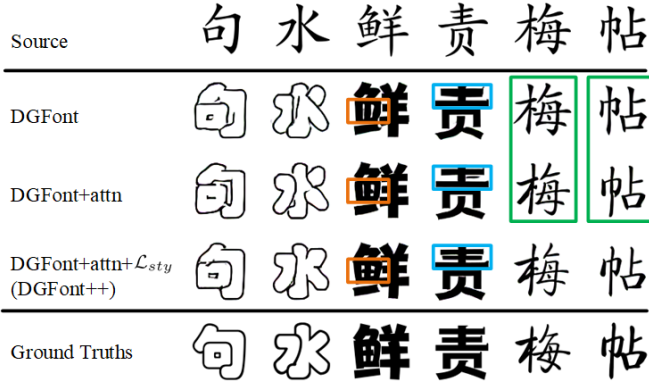


Fig. 7: Visual samples of ablation study for DGFont++.

Methods	L1 loss↓	RMSE↓	SSIM↑	LPIPS↓	FID↓
DGFont	0.0584	0.2166	0.7871	0.1204	41.69
DGFont+attn	0.0585	0.2170	0.7870	0.1190	39.67
DGFont+attn+ \mathcal{L}_{sty}	0.0567	0.2131	0.7910	0.1151	35.83

TABLE IV: Ablation study for DGFont++. “DGFont+attn+ \mathcal{L}_{sty} ” indicates the model of DGFont++.

in reference images number. It is because the contrastive representation model tends to learn robust and content-invariant style representation, which alleviates the influence of the variation of reference style images.

D. Analysis for FDSC Module

1) FDSC vs skip-connection

We compare our proposed FDSC module with commonly used skip-connection [71]. Skip-connection is often adopted to transfer feature maps with different resolutions directly from the encoder to the decoder, which is effective in semantic segmentation [79], [80] whose content of inputs and outputs share the same structure. However, font generation requires a geometric deformation between content inputs and the corresponding generated images in the structure. To compare the FDSC module with skip-connection, We replace two FDSC modules with skip-connection in our proposed DGFont network. The comparison results are shown in Table V. We can observe that models with FDSC modules outperform models with skip-connection, which proves the effectiveness of FDSC.

2) Visualization

In order to show the effectiveness of FDSC, we visualize the feature maps generated by the FDSC module. As shown in Figure 9, the feature maps F_c preserve the pattern of characters well, which helps generate a character with complete structure. On the other hand, we can observe that the FDSC module effectively transforms features extracted from the content encoder.

In addition, we visualize the learned offsets from the FDSC module using optical flow and character flow respectively. To visualize the offsets, the kernel of deformable convolution in the FDSC module is set to 1×1 . As demonstrated in Figure 10, we observe that the learned offsets mainly affect the character region. The offset value of the background tends to zero,

Method	L1 loss↓	RMSE↓	SSIM↑	LPIPS↓	FID↓
SC	0.0641	0.2212	0.7252	0.1114	46.88
FDSC	0.0582	0.2080	0.7469	0.1006	46.39

TABLE V: Comparison with skip-connection (SC) proposed by U-Net [71]. We replace two FDSC modules with skip-connections and then compare the new model with the full model of DGFont.

which proves the usefulness of the proposed offset loss Eq. 9. In character flow, we can see that most of the offset vectors point from the stroke in target characters to the corresponding source stroke. The results show that in the convolution process, the sampling locations of target characters tend to shift to corresponding locations in the source character by the learned offsets.

E. Style interpolation

To demonstrate that our style representations are semantically meaningful, we provide the style interpolation results in Figure 11. We first extract style features from two different fonts, and linearly interpolate the style factors Z_s to generate interpolated images. The interpolated factor is varied from 0 to 1 with an interval of 0.1. The results show that DGFont++ produces a smooth transition from one font to another and provides well-interpolated style features such as thickness, stroke, or joined-up writing while not hurting the character content.

F. User Study

To further compare the quality of images generated with other methods, we conduct an experiment on a human study by pairwise A/B tests. We compare our method with six compared methods. For each comparison, we randomly select 100 characters from the test set to make 100 paired images generated by our method and another compared method. The participants are 10 people who use Chinese characters every day. The participants are asked to select a more similar image compared to ground-truth from each pair within seven seconds. For each pair, we choose the image with more votes as the judgment result. Figure 12 shows the participants’ preference among the four tasks. We observe that more than 90 results of our methods outperform the results of Zi2zi, cycleGAN, EMD, which indicates that our method generates more realistic characters. Compared to FUNIT which gets impressive results in image-to-image translation tasks, our method still performs better on 85% results. MXFont is the most recent method for font generation. 76% results of our method are better than MXFont, which validates the superiority of our model over these state-of-the-art methods.

V. CONCLUSION

This paper proposes an effective unsupervised font generation model which is capable to generate realistic characters without paired images and can extend to unseen font well. We propose a Feature Deformation Skip Connection (FDSC) and

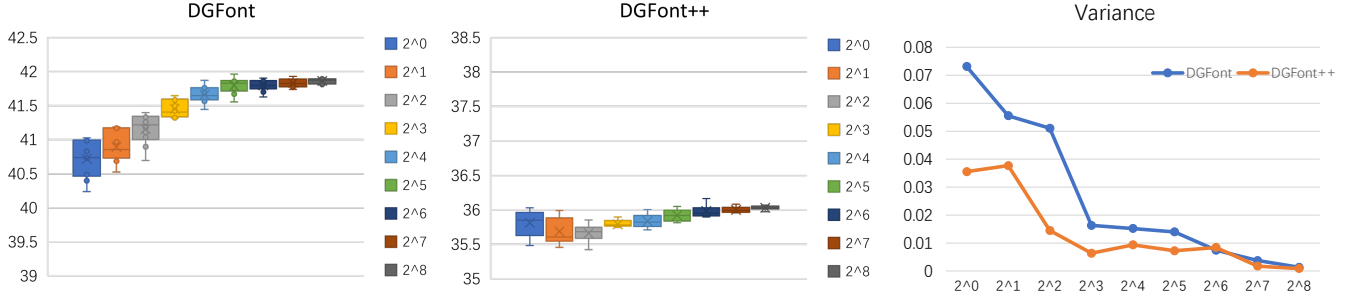


Fig. 8: Performance changes by varying the size of the reference set. We report how the performance is affected by the size of the reference style images. For each case, we calculate the FID score by random sampling the style images ten times. Figures (a) and (b) show the statistics of the FID score as box plots. Figure (c) shows the variance of FID for each model.

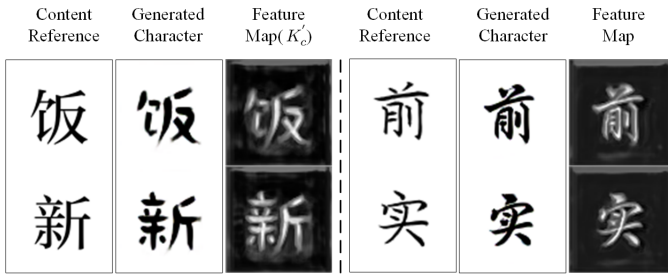


Fig. 9: **Feature visualization.** We visualize the features F_c' generated from the FDSC module. For each case, from left to right: content reference characters, the corresponding generated characters, and the visualization of feature maps. For feature map images, the whiter area represents the larger activation value.

FDSC-attn module to transfer the global and local deformable low-level spatial information to the mixer. Besides, we employ deformable convolution layers in the content encoder to learn style-invariant feature representations. Extensive experiments on font generation verify the effectiveness of our proposed model.

REFERENCES

- [1] P. Upchurch, N. Snively, and K. Bala, “From A to Z: supervised transfer of style and content using deep neural network generators,” *CoRR*, vol. abs/1603.02003, 2016.
- [2] S. Azadi, M. Fisher, V. G. Kim, Z. Wang, E. Shechtman, and T. Darrell, “Multi-content GAN for few-shot font style transfer,” in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*. IEEE Computer Society, 2018, pp. 7564–7573.
- [3] S. Fogel, H. Averbuch-Elor, S. Cohen, S. Mazor, and R. Litman, “ScrabbleGAN: Semi-supervised varying length handwritten text generation,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*. IEEE, 2020, pp. 4323–4332.
- [4] Rewrite, “<https://github.com/kaonashi-tyc/rewrite>.”
- [5] Zi2zi, “<https://github.com/kaonashi-tyc/zi2zi>.”
- [6] Y. Zhang, Y. Zhang, and W. Cai, “Separating style and content for generalized style transfer,” in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*. IEEE Computer Society, 2018, pp. 8447–8455.
- [7] D. Sun, T. Ren, C. Li, H. Su, and J. Zhu, “Learning to write stylized chinese characters by reading a handful of examples,” in *Proceedings*

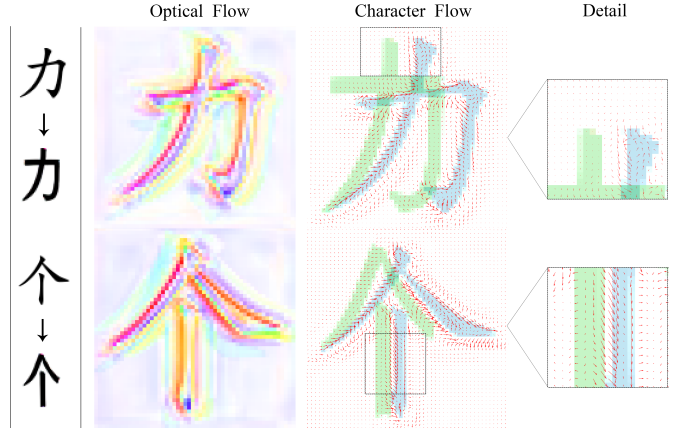


Fig. 10: **The visualization of learned offsets.** First column: source images to generated images. Second column: the deformation flows of the estimated offsets Δp . Third column: visualization of the estimated offsets Δp by quiver plot. Fourth column: zoomed-in details. The source and generated images are in blue and green respectively.



Fig. 11: Style interpolation. Generated glyphs in each row correspond to identical content representations. The leftmost and rightmost images are generated given two fonts. Images in the middle are convex combinations of two fonts.

of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden. ijcai.org, 2018, pp. 920–927.

- [8] Y. Jiang, Z. Lian, Y. Tang, and J. Xiao, “Scfont: Structure-guided chinese font generation via deep stacked networks,” in *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, Honolulu, Hawaii, USA, January 27 - February 1, 2019*. AAAI Press, 2019, pp.

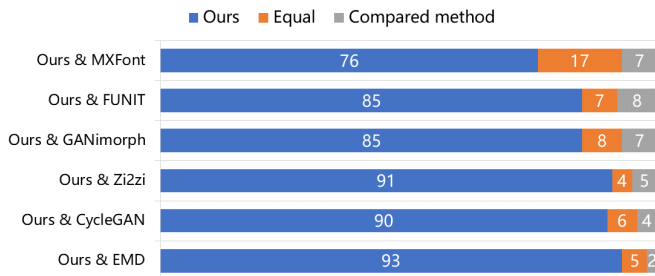


Fig. 12: **Results of user study.** For each compared method, we randomly sample 100 generated images for the participants' preferences test. The blue bar indicates the number of images that more participants prefer our results. The gray bar indicates the number of images that more participants prefer results from the compared methods. The orange bar indicates the number of images that get equal votes.

- 4015–4022.
- [9] Y. Huang, M. He, L. Jin, and Y. Wang, "RD-GAN: few/zero-shot chinese character style transfer via radical decomposition and rendering," in *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part VI*, ser. Lecture Notes in Computer Science, A. Vedaldi, H. Bischof, T. Brox, and J. Frahm, Eds., vol. 12351. Springer, 2020, pp. 156–172.
 - [10] J. Cha, S. Chun, G. Lee, B. Lee, S. Kim, and H. Lee, "Few-shot compositional font generation with dual memory," in *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XIX*, ser. Lecture Notes in Computer Science, A. Vedaldi, H. Bischof, T. Brox, and J. Frahm, Eds., vol. 12364. Springer, 2020, pp. 735–751.
 - [11] S. Park, S. Chun, J. Cha, B. Lee, and H. Shim, "Few-shot font generation with localized style representations and factorization," *CoRR*, vol. abs/2009.11042, 2020.
 - [12] —, "Multiple heads are better than one: Few-shot font generation with multiple localized experts," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 13 900–13 909.
 - [13] Y. Gao and J. Wu, "Gan-based unpaired chinese character image translation via skeleton transformation and stroke rendering," in *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, New York, NY, USA, February 7-12, 2020*. AAAI Press, 2020, pp. 646–653.
 - [14] B. Chang, Q. Zhang, S. Pan, and L. Meng, "Generating handwritten chinese characters using cyclegan," in *2018 IEEE Winter Conference on Applications of Computer Vision, WACV 2018, Lake Tahoe, NV, USA, March 12-15, 2018*. IEEE Computer Society, 2018, pp. 199–207.
 - [15] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*. IEEE Computer Society, 2017, pp. 2261–2269.
 - [16] T. Kim, M. Cha, H. Kim, J. K. Lee, and J. Kim, "Learning to discover cross-domain relations with generative adversarial networks," in *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, 2017, pp. 1857–1865.
 - [17] Z. Yi, H. Zhang, P. Tan, and M. Gong, "Dualgan: Unsupervised dual learning for image-to-image translation," in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
 - [18] J. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*. IEEE Computer Society, 2017, pp. 2242–2251.
 - [19] Y. Taigman, A. Polyak, and L. Wolf, "Unsupervised cross-domain image generation," *CoRR*, vol. abs/1611.02200, 2016.
 - [20] S. Benaim and L. Wolf, "One-sided unsupervised domain mapping," in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, 2017, pp. 752–762.
 - [21] K. Baek, Y. Choi, Y. Uh, J. Yoo, and H. Shim, "Rethinking the truly unsupervised image-to-image translation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 14 154–14 163.
 - [22] X. Huang and S. J. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*. IEEE Computer Society, 2017, pp. 1510–1519.
 - [23] M. Liu, X. Huang, A. Mallya, T. Karras, T. Aila, J. Lehtinen, and J. Kautz, "Few-shot unsupervised image-to-image translation," in *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*. IEEE, 2019, pp. 10 550–10 559.
 - [24] Y. Xie, X. Chen, L. Sun, and Y. Lu, "Dg-font: Deformable generative networks for unsupervised font generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 5130–5140.
 - [25] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, 2014, pp. 2672–2680.
 - [26] D. Sun, Q. Zhang, and J. Yang, "Pyramid embedded generative adversarial network for automated font generation," in *24th International Conference on Pattern Recognition, ICPR 2018, Beijing, China, August 20-24, 2018*. IEEE Computer Society, 2018, pp. 976–981.
 - [27] J. Chang, Y. Gu, Y. Zhang, and Y. Wang, "Chinese handwriting imitation with hierarchical generative adversarial network," in *British Machine Vision Conference 2018, BMVC 2018, Newcastle, UK, September 3-6, 2018*. BMVA Press, 2018, p. 290.
 - [28] P. Lyu, X. Bai, C. Yao, Z. Zhu, T. Huang, and W. Liu, "Auto-encoder guided GAN for chinese calligraphy synthesis," in *14th IAPR International Conference on Document Analysis and Recognition, ICDAR 2017, Kyoto, Japan, November 9-15, 2017*. IEEE, 2017, pp. 1095–1100.
 - [29] Y. Jiang, Z. Lian, Y. Tang, and J. Xiao, "Dcfont: An end-to-end deep chinese font generation system," in *SIGGRAPH Asia 2017 Technical Briefs, Bangkok, Thailand, November 27 - 30, 2017*. ACM, 2017, pp. 22:1–22:4.
 - [30] S.-J. Wu, C.-Y. Yang, and J. Hsu, "Calligan: Style and structure-aware chinese calligraphy character generator," *ArXiv*, vol. abs/2005.12500, 2020.
 - [31] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II*, ser. Lecture Notes in Computer Science, vol. 9906. Springer, 2016, pp. 694–711.
 - [32] H. Zhang and K. J. Dana, "Multi-style generative network for real-time transfer," in *Computer Vision - ECCV 2018 Workshops - Munich, Germany, September 8-14, 2018, Proceedings, Part IV*, ser. Lecture Notes in Computer Science, vol. 11132. Springer, 2018, pp. 349–365.
 - [33] S. Shukla, L. V. Gool, and R. Timofte, "Extremely weak supervised image-to-image translation for semantic segmentation," in *2019 IEEE/CVF International Conference on Computer Vision Workshops, ICCV Workshops 2019, Seoul, Korea (South), October 27-28, 2019*. IEEE, 2019, pp. 3368–3377.
 - [34] L. Musto and A. Zinelli, "Semantically adaptive image-to-image translation for domain adaptation of semantic segmentation," *BMVC*, 2020.
 - [35] C. Wang, C. Xu, and D. Tao, "Self-supervised pose adaptation for cross-domain image animation," *IEEE Transactions on Artificial Intelligence*, vol. 1, no. 1, pp. 34–46, 2020.
 - [36] E. Yang, C. Deng, W. Liu, X. Liu, D. Tao, and X. Gao, "Pairwise relationship guided deep hashing for cross-modal retrieval," in *AAAI*, 2017, pp. 1618–1625.
 - [37] Z. Chen, C. Wang, B. Yuan, and D. Tao, "Puppeteergan: Arbitrary portrait animation with semantic-aware appearance transformation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
 - [38] X. Chen, C. Xu, X. Yang, and D. Tao, "Attention-gan for object transfiguration in wild images," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 164–180.
 - [39] C. Chan, S. Ginosar, T. Zhou, and A. A. Efros, "Everybody dance now," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 5933–5942.
 - [40] X. Chen, C. Xu, X. Yang, and D. Tao, "Long-term video prediction via criticization and retrospection," *IEEE Transactions on Image Processing*, vol. 29, pp. 7090–7103, 2020.

- [41] J. Dong, X. Li, C. Xu, X. Yang, G. Yang, X. Wang, and M. Wang, "Dual encoding for video retrieval by text," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [42] P. Isola, J. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*. IEEE Computer Society, 2017, pp. 5967–5976.
- [43] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *CoRR*, vol. abs/1411.1784, 2014.
- [44] M. Liu and O. Tuzel, "Coupled generative adversarial networks," in *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain, 2016*, pp. 469–477.
- [45] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan, "Unsupervised pixel-level domain adaptation with generative adversarial networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*. IEEE Computer Society, 2017, pp. 95–104.
- [46] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb, "Learning from simulated and unsupervised images through adversarial training," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*. IEEE Computer Society, 2017, pp. 2242–2251.
- [47] X. Chen, C. Xu, X. Yang, L. Song, and D. Tao, "Gated-gan: Adversarial gated networks for multi-collection style transfer," *IEEE Trans. Image Process.*, vol. 28, no. 2, pp. 546–560, 2019.
- [48] D. Bhattacharjee, S. Kim, G. Vizier, and M. Salzmann, "DUNIT: detection-based unsupervised image-to-image translation," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*. IEEE, 2020, pp. 4786–4795.
- [49] A. Gokaslan, V. Ramanujan, D. Ritchie, K. I. Kim, and J. Tompkin, "Improving shape deformation in unsupervised image-to-image translation," in *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XII*, ser. Lecture Notes in Computer Science, V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, Eds., vol. 11216. Springer, 2018, pp. 662–678.
- [50] W. Wu, K. Cao, C. Li, C. Qian, and C. C. Loy, "Transgaga: Geometry-aware unsupervised image-to-image translation," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*. Computer Vision Foundation / IEEE, 2019, pp. 8012–8021.
- [51] X. Huang, M. Liu, S. J. Belongie, and J. Kautz, "Multimodal unsupervised image-to-image translation," in *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part III*, ser. Lecture Notes in Computer Science, V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, Eds., vol. 11207. Springer, 2018, pp. 179–196.
- [52] J. Dai, H. Qi, Y. Xiong, Y. Li, G. Zhang, H. Hu, and Y. Wei, "Deformable convolutional networks," in *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*. IEEE Computer Society, 2017, pp. 764–773.
- [53] G. Bertasius, L. Torresani, and J. Shi, "Object detection in video with spatiotemporal sampling networks," in *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XII*, ser. Lecture Notes in Computer Science, vol. 11216. Springer, 2018, pp. 342–357.
- [54] X. Zhu, H. Hu, S. Lin, and J. Dai, "Deformable convnets V2: more deformable, better results," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*. Computer Vision Foundation / IEEE, 2019, pp. 9308–9316.
- [55] Z. Chen, W. Li, C. Fei, B. Liu, and N. Yu, "Spatial-temporal feature aggregation network for video object detection," in *2020 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2020, Barcelona, Spain, May 4-8, 2020*. IEEE, 2020, pp. 1858–1862.
- [56] X. Sun, B. Xiao, F. Wei, S. Liang, and Y. Wei, "Integral human pose regression," in *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part VI*, ser. Lecture Notes in Computer Science, vol. 11210. Springer, 2018, pp. 536–553.
- [57] Y. Tian, Y. Zhang, Y. Fu, and C. Xu, "TDAN: temporally-deformable alignment network for video super-resolution," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*. IEEE, 2020, pp. 3357–3366.
- [58] M. Yin, L. Sun, and Q. Li, "Novel view synthesis on unpaired data by conditional deformable variational auto-encoder," in *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XXVIII*, ser. Lecture Notes in Computer Science, vol. 12373. Springer, 2020, pp. 87–103.
- [59] —, "Id-unet: Iterative soft and hard deformation for view synthesis," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*. Computer Vision Foundation / IEEE, 2021, pp. 7220–7229.
- [60] X. Zhu, D. Cheng, Z. Zhang, S. Lin, and J. Dai, "An empirical study of spatial attention mechanisms in deep networks," in *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*. IEEE, 2019, pp. 6687–6696.
- [61] H. Zhang, I. J. Goodfellow, D. N. Metaxas, and A. Odena, "Self-attention generative adversarial networks," in *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. PMLR, 2019, pp. 7354–7363.
- [62] N. Yu, G. Liu, A. Dundar, A. Tao, B. Catanzaro, L. Davis, and M. Fritz, "Dual contrastive loss and attention for gans," *CoRR*, vol. abs/2103.16748, 2021.
- [63] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton, "A simple framework for contrastive learning of visual representations," in *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 1597–1607.
- [64] T. Chen, S. Kornblith, K. Swersky, M. Norouzi, and G. E. Hinton, "Big self-supervised models are strong semi-supervised learners," in *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., 2020.
- [65] K. He, H. Fan, Y. Wu, S. Xie, and R. B. Girshick, "Momentum contrast for unsupervised visual representation learning," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*. Computer Vision Foundation / IEEE, 2020, pp. 9726–9735.
- [66] X. Chen, H. Fan, R. B. Girshick, and K. He, "Improved baselines with momentum contrastive learning," *CoRR*, vol. abs/2003.04297, 2020.
- [67] O. J. Hénaff, "Data-efficient image recognition with contrastive predictive coding," in *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 4182–4192.
- [68] R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, and Y. Bengio, "Learning deep representations by mutual information estimation and maximization," in *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.
- [69] T. Park, A. A. Efros, R. Zhang, and J. Zhu, "Contrastive learning for unpaired image-to-image translation," in *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part IX*, ser. Lecture Notes in Computer Science, A. Vedaldi, H. Bischof, T. Brox, and J. Frahm, Eds., vol. 12354. Springer, 2020, pp. 319–345.
- [70] M. Gutmann and A. Hyvärinen, "Noise-contrastive estimation: A new estimation principle for unnormalized statistical models," in *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2010, Chia Laguna Resort, Sardinia, Italy, May 13-15, 2010*, ser. JMLR Proceedings, Y. W. Teh and D. M. Titterton, Eds., vol. 9. JMLR.org, 2010, pp. 297–304.
- [71] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention - MICCAI 2015 - 18th International Conference Munich, Germany, October 5 - 9, 2015, Proceedings, Part III*, ser. Lecture Notes in Computer Science, vol. 9351. Springer, 2015, pp. 234–241.
- [72] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, "Unsupervised feature learning via non-parametric instance discrimination," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3733–3742.
- [73] S. Singh and S. Krishnan, "Filter response normalization layer: Eliminating batch dependence in the training of deep neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 237–11 246.
- [74] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *2015 IEEE International Conference on Computer Vision, ICCV 2015*,

- Santiago, Chile, December 7-13, 2015. IEEE Computer Society, 2015, pp. 1026–1034.
- [75] D. Tran, R. Ranganath, and D. M. Blei, “Deep and hierarchical implicit models,” *CoRR*, vol. abs/1702.08896, 2017.
 - [76] L. M. Mescheder, A. Geiger, and S. Nowozin, “Which training methods for gans do actually converge?” in *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, ser. Proceedings of Machine Learning Research, J. G. Dy and A. Krause, Eds., vol. 80. PMLR, 2018, pp. 3478–3487.
 - [77] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, “The unreasonable effectiveness of deep features as a perceptual metric,” in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*. IEEE Computer Society, 2018, pp. 586–595.
 - [78] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, “Gans trained by a two time-scale update rule converge to a local nash equilibrium,” in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA, 2017*, pp. 6626–6637.
 - [79] S. Jégou, M. Drozdal, D. Vázquez, A. Romero, and Y. Bengio, “The one hundred layers tiramisu: Fully convolutional densenets for semantic segmentation,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2017, Honolulu, HI, USA, July 21-26, 2017*. IEEE Computer Society, 2017, pp. 1175–1183.
 - [80] J. Long, E. Shelhamer, and T. Darrell, “Fully convolutional networks for semantic segmentation,” in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*. IEEE Computer Society, 2015, pp. 3431–3440.