# CF-Font: Content Fusion for Few-shot Font Generation

Chi Wang[1,2*], Min Zhou[2], Tiezheng Ge[2], Yuning Jiang[2], Hujun Bao[1], Weiwei Xu[1†]

[1] State Key Lab of CAD&CG, Zhejiang University    [2] Alibaba Group

`wangchi1995@zju.edu.cn`, `{yunqi.zm, tiezheng.gtz, mengzhu.jyn}@alibaba-inc.com`
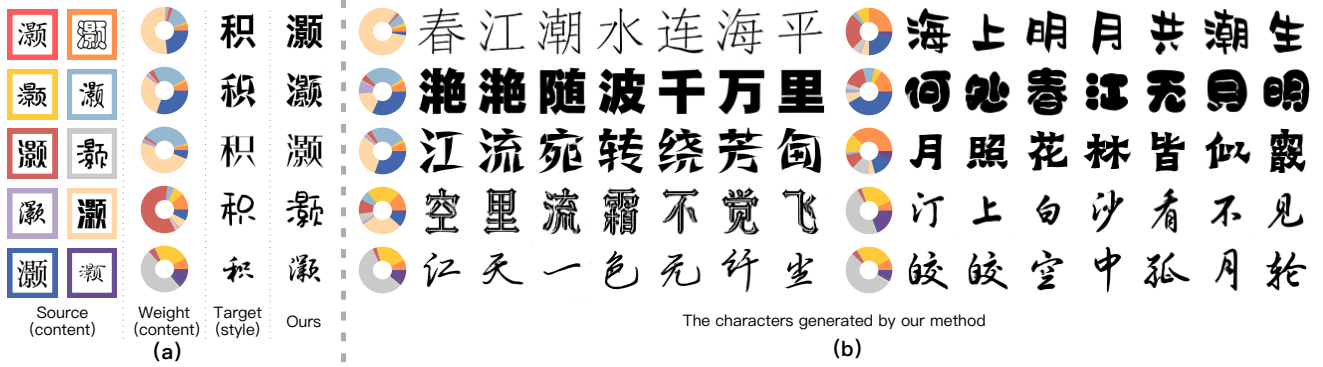`{bao, xww}@cad.zju.edu.cn`

Figure 1. Characters generated by our method. (a) Source: source character images selected from ten basis fonts for content feature fusion. Weights: different colors and their covered areas on the doughnut chart represent the weights used to blend content features adaptively. Ten colors correspond to source images in colored boxes. Target: few-shot target reference character images. One of those is performed as an example. Ours: images generated by our method with fused content features and style features. (b) Generated character images of the first ten lines from a famous Chinese poem, each line with an extracted style, *e.g.* thin, thick, swollen, cuneiform, inscription, or cursive style.

## Abstract

*Content and style disentanglement is an effective way to achieve few-shot font generation. It allows to transfer the style of the font image in a source domain to the style defined with a few reference images in a target domain. However, the content feature extracted using a representative font might not be optimal. In light of this, we propose a content fusion module (CFM) to project the content feature into a linear space defined by the content features of basis fonts, which can take the variation of content features caused by different fonts into consideration. Our method also allows to optimize the style representation vector of reference images through a lightweight iterative style-vector refinement (ISR) strategy. Moreover, we treat the 1D projection of a character image as a probability distribution and leverage the distance between two distributions as the reconstruction loss (namely projected character loss, PCL). Compared to L2 or L1 reconstruction loss, the distribution distance pays more attention to the global shape of characters. We*

*have evaluated our method on a dataset of 300 fonts with 6.5k characters each. Experimental results verify that our method outperforms existing state-of-the-art few-shot font generation methods by a large margin. The source code can be found at https://github.com/wangchi95/CF-Font.*

## 1. Introduction

Few-shot font generation aims to produce characters of a new font by transforming font images from a source domain to a target domain according to just a few reference images. It can greatly reduce the labor of expert designers to create a new style of fonts, especially for logographic languages that contain multiple characters, such as Chinese (over 60K characters), Japanese (over 50K characters), and Korean (over 11K characters), since only several reference images need to be manually designed. Therefore, font generation has wide applications in font completion for ancient books and monuments, personal font generation, etc.

Recently, with the rapid development of convolutional neural networks [22] and generative adversarial networks [9] (GAN), pioneers have made great progress in

---

generating gratifying logographic fonts. Zi2zi [38] introduces pix2pix [14] method to generate complex characters of logographic languages with high quality, but it cannot handle those fonts that do not appear in training (unseen fonts). For the few-shot font generation, many methods [3, 7, 31, 32, 34, 42, 47] verify that content and style disentanglement is effective to convert the style of a character in the source domain, denoted as *source character*, to the target style embodied with reference images of seen or unseen fonts. The neural networks in these methods usually have two branches to learn content and style features respectively, and the content features are usually obtained with the character image from a manually-chosen font, denoted as *source font*. However, since it's a difficult task to achieve a complete disentanglement between content and style features [17, 21], the choice of the font for content-feature encoding influences the font generation results substantially. For instance, *Song* and *Kai* are commonly selected as the source font [20, 28, 31, 42, 43, 47]. While such choices are effective in many cases, the generated images sometimes contain artifacts, such as incomplete and unwanted strokes.

The main contribution of this paper is a novel content feature fusion scheme to mitigate the influence of incomplete disentanglement by exploring the synchronization of content and style features, which significantly enhances the quality of few-shot font generation. Specifically, we design a content fusion module (CFM) to take the content features of different fonts into consideration during training and inference. It is realized by computing the content feature of a character of a target font through linearly blending content features of the corresponding characters in the automatically determined basis fonts, and the blending weights are determined through a carefully designed font-level distance measure. In this way, we can form a linear cluster for the content feature of a semantic character, and explore how to leverage the font-level similarity to seek for an optimized content feature in this cluster to improve the quality of generated characters.

In addition, we introduce an iterative style-vector refinement (ISR) strategy to find a better style feature vector for font-level style representation. For each font, we average the style vectors of reference images and treat it as a learnable parameter. Afterward, we fine-tune the style vector with a reconstruction loss, which further improves the quality of the generated fonts.

Most font-generation algorithms [3, 20, 31, 32, 38, 42] choose L1 loss as the character image reconstruction loss. However, L1 or L2 loss mainly supervises per-pixel accuracy and is easily disturbed by the local misalignment of details. Hence, we employ a distribution-based projected character loss (PCL) to measure the shape difference between characters. Specifically, by treating the 1D projection of 2D character images as a 1D probability distribution,

PCL computes the distribution distance to pay more attention to the global properties of character shapes, resulting in the large improvement of skeleton topology transfer results.

The CFM can be embedded into the few-shot font generation task to enhance the quality of generated results. Extensive experiments verify that our method, referred to as CF-Font, remarkably outperforms state-of-the-art methods on both seen and unseen fonts. Fig. 1 reveals that our method can generate high-quality fonts of various styles.

## 2. Related Works

### 2.1. Image-to-image Translation

Image-to-image translation is the task of converting a source image to the target domain of reference images. Early methods [6, 14, 33, 40, 48] utilize GAN [9] and yield vivid images. But they could only convert the source image to some specific domains (or categories), which is more limited in practical applications. Recently, some few-shot methods [1, 2, 5, 13, 18, 24] are proposed. These methods disentangle the content and style, and can convert the source image to arbitrary styles only if a few reference images are provided. Further, RG-UNIT [10] proposes an image retrieval strategy to help domain transfer, *i.e.* it finds images similar to the source in content but in the target domain, and extracts their content features as assistance. Though the retrieval strategy helps to generate more realistic images, it cannot be directly applied to font generation tasks. Because the retrieved image may still differ significantly from the target in content, as fonts are highly fine-grained. Thus, we build basis fonts and use fused content features to narrow the gap between the source and target domains.

### 2.2. Few-shot Font Generation

Few-shot font generation aims to generate a new font library in the required style with only a few reference images. Early methods [4, 15, 29, 35, 38] for font generation train a cross-domain translation network to model mapping from the source to the target domain. These structures limit the model to generate unseen fonts. To address this issue, SA-VAE [34] and EMD [47] disentangle the representations of style and content, and can generate images of all style-content combinations. RD-GAN [8], SCFont [16], CalliGAN [41], and LF-Font [31] follow this way and employ component annotations to boost the style representation in local regions. To be less dependent on explicit component annotations, MX-Font [32] utilizes multiple experts and bipartite matching, and XMP-Font [25] employs a cross-modality encoder, which is conditioned jointly on character images and stroke labels. CG-GAN [20] supervises a font generator to decouple content and style on component level through a component-aware module. But these three methods still require the labels of component categories. Fs-Font
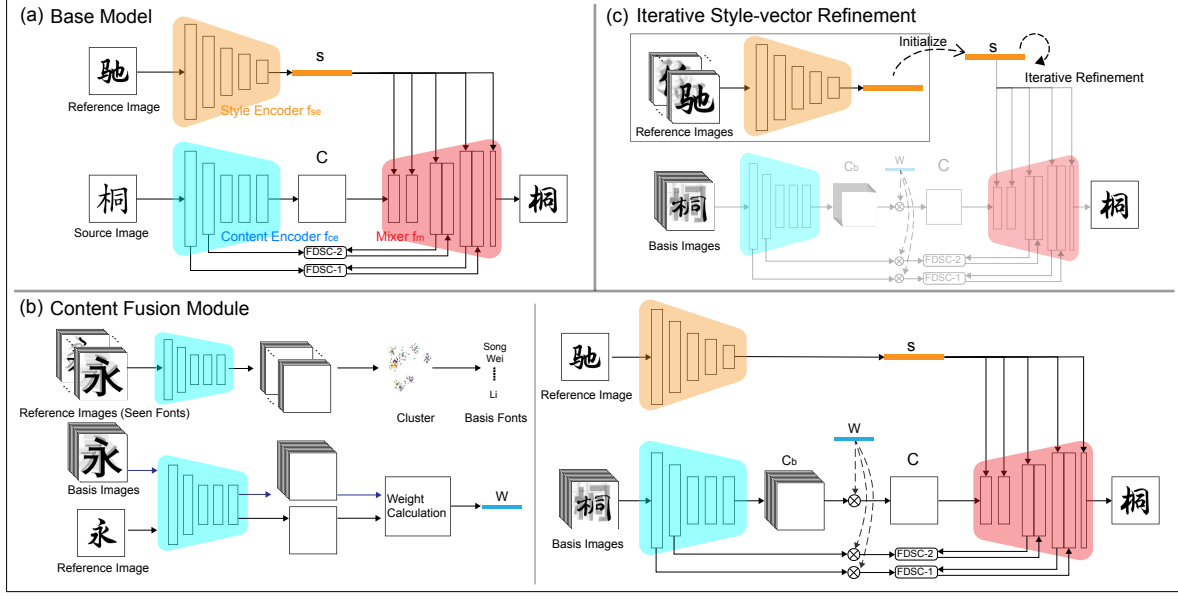
Figure 2. The framework of our model. (a) We first train the DGN [42] and use PCL to enhance the supervision of character skeletons. (b) After the model converges, content features of all training fonts are clustered and basis fonts are selected according to cluster centers. The original content encoder is replaced by CFM, and original content features are changed to fused features of basis fonts. Then we continue to train the model so that it adapts to fused content features. (c) In inference, we utilize ISR to polish the style of a font. The extracted mean style vector is treated as the only trainable variable to be fine-tuned for a few iterations.

is proposed to learn fine-grained local styles from reference images, and the spatial correspondence between the content and reference images. [36] However, it needs to select reference characters carefully to achieve high-quality generated results. DG-Font [42] introduces a feature deformation skip connection module and achieves excellent performance without any extra labels. However, it is difficult for these few-shot methods to generate new fonts if the source and target domains are very different, especially when the target font is unseen. Starting from this perspective, we propose the CFM to reduce the difficulty of domain transfer, and the PCL to enhance skeleton supervision.

## 3. Approach

Our method is illustrated in Fig. 2. The whole training pipeline can be divided into two stages. Firstly, we train the neural network in DG-Font [42] as our base network, referred to as DGN. The network is used to learn basic, disentangled content and style features of character images in our dataset. Secondly, our content fusion module (CFM) is plugged into the model after the content encoder. Afterward, we replace the original content feature with the output of CFM, a linear content-feature interpolation of automatically-selected basis fonts. Then, we fix the content encoder and continue to train style encoder, feature deformation skip connection [42] (FSDC) and mixer together for a few epochs. The projected character loss (PCL) is used in training to supervise character skeletons. In addition, to further improve the generation quality, we utilize the iterative

style-vector refinement (ISR) strategy to polish the learned font-level style vector alone in inference. The motivation for ISR is seeking for a single and high-quality font-level style vector to generate images for all characters of the font. Specifically, for a font, we refine upon the average of the character style vector of all the given 16 characters in our few-shot setting.

### 3.1. Base Network

As illustrated in Fig. 2 (a), given a content image $I_c$ and a style image $I_s$, the DGN synthesizes an image with the character of the content image and the font of the style image. This generative network consists of four parts: a style encoder $f_{se}$ to extract style latent vector $s$, a content encoder $f_{ce}$ to obtain content feature map $C$, a mixer $f_m$ to mix style and content representations with AdaIN [13], and two FSDC modules. During training, a multi-task discriminator, fed with generated characters and their ground-truth images, is applied to conduct discrimination for each style simultaneously.

Four losses are adopted: 1) image reconstruction loss $\mathcal{L}_{img}$ for domain-invariant features maintaining; 2) content consistent loss $\mathcal{L}_{cnt}$ to guarantee consistency between generated and input content images; 3) adversarial loss $\mathcal{L}_{adv}$ in hinge version [23, 30, 45] for realistic image generation; 4) deformation offset normalization $\mathcal{L}_{offset}$ to avoid excessive offsets in FDSC. More details are in [42].
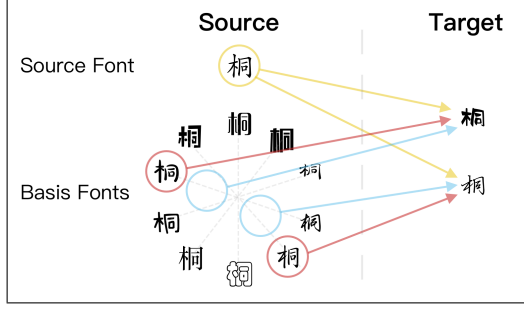
Figure 3. Visualization of content fusion. The yellow and red arrows are denoted for content features from the commonly used source font *Kai* [31] and the nearest font of the target respectively. The blue arrow represents the interpolation of content features of basis fonts to approximate the target.

## 3.2. Content Fusion Module

The content fusion module aims to adaptively extract content features by combining the content features of basis fonts. This network with CFM is constructed as in Fig. 2 (b). Firstly, to find representative fonts for content fusion, we cluster all training fonts through the concatenated content features of the given 16 few-shot characters and pick those nearest to the cluster centers as basis fonts. The basis fonts are fixed once selected. Then, for each target font, we calculate an L1-norm content fusion weight according to its similarity to basis fonts. As a result, the content features (input of the decoder) are replaced by the weighted sum of those of basis fonts. In addition, the network should be fine-tuned for a few epochs to adapt to fused content features (represented as the blue circles in Fig. 3).

**Basis selection.** Suppose we need to choose $M$ basis fonts from $N$ training fonts. It can be realized by clustering the content features $\{\boldsymbol{C}_i\}_{i=1}^N$ and selecting fonts lying in the cluster centroids. In our practice, since the dimension of $\boldsymbol{C}_i$ is too large while $N$ is relatively small, we follow [37] to map $\boldsymbol{C}_i$ to the vector of the distances between it and features of all fonts $\boldsymbol{e}_i \in \mathbb{R}^N$. More specifically:

$$
\begin{aligned}
\boldsymbol{C}_i &= f_{ce}(\boldsymbol{I}_i),\,^1 \\
\boldsymbol{d}_i &= (d_{i1}, d_{i2}, ..., d_{iN}), \quad d_{ij} = \|\boldsymbol{C}_i - \boldsymbol{C}_j\|_1, \\
\boldsymbol{e}_i &= \sigma(\boldsymbol{d}_i), \\
\mathcal{B} &= \textbf{Cluster}(M, \{\boldsymbol{e}_1, \boldsymbol{e}_2, ..., \boldsymbol{e}_N\}),
\end{aligned}
\tag{1}
$$

where $\sigma(\cdot)$ is the softmax operation, $d_{ij}$ is the L1 distance between two fonts, **Cluster** is the classical K-Medoids

---

[1] $\boldsymbol{C}_i$ is actually the concatenated content features extracted from several characters of font $i$. For the sake of brevity, we omit the superscript for characters here.
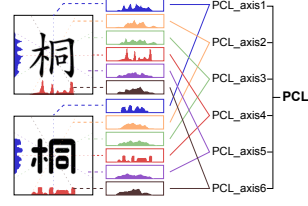


Figure 4. Illustration of PCL. We project the binary characters into multi-direction 1D spaces (distinguished by color) and calculate normalized histograms for each. It is obvious that for the different fonts with the same character, the projected distributions vary along with the skeletons and are less sensitive to textures or colors.

cluster algorithm [39], and set $\mathcal{B}$ is the indices of selected fonts.

**Weight calculation.** For the target font $t$ and its content feature $\boldsymbol{C}_t$, we measure its similarity to the basis fonts $\{\boldsymbol{C}_m\}_{m=1}^M$, namely $\boldsymbol{d}'_t \in \mathbb{R}^M$. Then the content fusion weight $\boldsymbol{w}_t \in \mathbb{R}^M$ is calculated as follow:

$$
\begin{aligned}
\boldsymbol{d}'_t &= (d_{t1}, d_{t2}, ..., d_{tM}), \quad d_{tm} = \|\boldsymbol{C}_t - \boldsymbol{C}_m\|_1, \\
\boldsymbol{w}_t &= \sigma(-\boldsymbol{d}'_t/\tau),
\end{aligned}
\tag{2}
$$

where $\tau$ is the temperature of the softmax operation.

**Content fusion.** Once the basis fonts and content fusion weights are obtained, the original content feature map $\boldsymbol{C}$ is replaced with the fused one $\boldsymbol{C}'_t$, where the content fusion weight of CFM is related to its target font $t$.

$$
\boldsymbol{C}'_t = \sum_{m \in \mathcal{B}} w_{tm} \cdot \boldsymbol{C}_m.
\tag{3}
$$

## 3.3. Projected Character Loss

To better supervise the skeleton, we design a projected character loss, which measures image difference with marginal distribution distances on multiple 1D projections, as shown in Fig. 4. Since the distribution is sensitive to the relative relationship, PCL pays more attention to the global shape of characters.

$$
\mathcal{L}_p(\boldsymbol{Y}, \hat{\boldsymbol{Y}}) = \frac{1}{P} \sum_{p=1}^{P} \mathcal{L}_{1d}(\phi_p(\boldsymbol{Y}), \phi_p(\hat{\boldsymbol{Y}})),
\tag{4}
$$

where $\boldsymbol{Y}$ and $\hat{\boldsymbol{Y}}$ represent the generated and ground-truth image respectively, $P$ is the number of projections, and $\phi_p(\cdot)$ denotes a projection function with the $p$-th direction.

There are lots of metrics to measure the alignment between 1D distributions, such as the KL-divergence and

Figure 5. L1 vs PCL. We retrieve the closest character of all training fonts to the top-left one by L1, PC-WDL, and PC-KL, respectively. The top ten results of each loss are listed from left to right, top to down. It can be seen that the skeletons vary greatly in the column of L1 but are quite consistent in those of PCL.

Wasserstein distance. Thus, $\mathcal{L}_p$ can have various forms:

$$\mathcal{L}_{pc-wdl}(\boldsymbol{Y}, \hat{\boldsymbol{Y}}) = \frac{1}{P}\sum_{p=1}^{P}\left\|\frac{\Lambda(\phi_p(\boldsymbol{Y}))}{\sum \phi_p(\boldsymbol{Y})} - \frac{\Lambda(\phi_p(\hat{\boldsymbol{Y}}))}{\sum \phi_p(\hat{\boldsymbol{Y}})}\right\|$$

$$\mathcal{L}_{pc-kl}(\boldsymbol{Y}, \hat{\boldsymbol{Y}}) = \frac{1}{P}\sum_{p=1}^{P}\mathbf{KL}\left(\frac{\phi_p(\boldsymbol{Y})}{\sum \phi_p(\boldsymbol{Y})}, \frac{\phi_p(\hat{\boldsymbol{Y}})}{\sum \phi_p(\hat{\boldsymbol{Y}})}\right),$$
(5)

where **KL** means the KL-divergence and $\Lambda$ denotes the cumsum function, which turns probability density functions to cumulative distribution functions.

To simply verify the performance of PCL, we generate images of the character "Tong" from 240 fonts and measure their similarity by PCL and L1. The closest ten characters to the top-left one found by different metrics are displayed in Fig 5 respectively. It can be seen that the characters retrieved by L1 are quite different on the character skeleton, which is important for fonts. While those selected by PCL are relatively more consistent and it indicates that PCL is more proper for measuring the skeleton.

Adding PCL to the image reconstruction loss term, we have the following overall loss function for training:

$$\mathcal{L} = \mathcal{L}_{adv} + \lambda_{img}(\mathcal{L}_{img} + \lambda_{pcl}\mathcal{L}_{pcl})$$
$$+ \lambda_{cnt}\mathcal{L}_{cnt} + \lambda_{offset}\mathcal{L}_{offset},$$
(6)

where $\lambda_{img}$, $\lambda_{pcl}$, $\lambda_{cnt}$, and $\lambda_{offset}$ are hyperparameters to adjust the weight of each loss function.

### 3.4. Iterative Style-vector Refinement

For target font $t$, a robust style information can be extracted as the latent style vector $\boldsymbol{s}'_t$ by averaging the outputs of $f_{se}$ with a set of character images [42].

$$\boldsymbol{s}'_t = \frac{1}{Q}\sum_{q=1}^{Q}f_{se}(\boldsymbol{I}_t^q),$$
(7)

where $\boldsymbol{I}_t^q$ is an image of character $q$ of font $t$, and Q denotes the reference character number.

Motivated by the "iterative inference" strategy that optimizes input in the inference stage (*e.g.* [44]), we propose iterative style-vector refinement for further optimizing the style feature $\boldsymbol{s}'_t$. As in Fig. 2 (c), in the inference stage, $\boldsymbol{s}'_t$ is first initialized by Eq. 7. Then, using the provided few reference characters of target fonts $\{\boldsymbol{I}_t^q\}_{q=1}^{Q}$ as supervising samples, we refine $\boldsymbol{s}'_t$ for around ten epochs according to the backpropagation of the reconstruction loss. Finally, the optimized style vector is adopted for inference. Worth noting this style vector can be stored as a signature of the target font and reused in referencing all characters of the same font, which makes the proposed ISR efficient in the real system.

## 4. Experiments

We have implemented the CF-Font method on a GPU server with 8 Nvidia Tesla V100 GPUs. After training with our dataset, our method outperforms the state-of-the-art methods on unseen fonts by 5.7% and 5.0% with respect to L1 and FID metrics, respectively. In the following, we report the preparation of dataset, evaluation metrics, and various experimental results to verify the effectiveness of our method.

### 4.1. Dataset and Evaluation Metrics

We collect 300 Chinese fonts to build a dataset (including printed and handwriting fonts) to verify our method for the Chinese font generation task. Our character set (6446 in total) covers almost the full standard Chinese character set (6763 in total) of GB/T 2312 [27], and 317 characters not supported by comparison methods are removed. The training part contains 240 fonts, and each font has 800 characters. The test part consists of (a) 229 seen fonts with 5646 unseen characters; (b) the remaining 60 unseen fonts with 5646 unseen characters, to verify the generalization ability of models. Note that we exclude 11 of the 240 training fonts when testing on seen fonts. They are basis fonts (including *Song*) in CFM and a font *Kai*, in which *Song* and *Kai* are commonly used as source fonts in font generation [20, 31, 42, 47]. Besides, for few-shot font generation, reference images of target fonts in the test are with 16 randomly picked characters from the training part.

We leverage pixel-level and perceptual metrics for evaluation, following [42]. Specifically, pixel-level metrics are L1, root mean square error (RMSE), and structural similarity index measure (SSIM). They focus on per-pixel consistency between generated images and ground-truth ones. Perceptual metrics include FID [11] and LPIPS [46], both of which measure the similarity of features and are closer to human vision.

Table 1. Comparison with state-of-the-art methods on seen/unseen fonts. Bold and underlined numbers denote the best and the second best respectively. The numbers in the last row represent our improvement over the second-best scores.

| Methods | Seen Fonts | | | | | Unseen Fonts | | | | | User Study % |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | L1↓ | RMSE↓ | SSIM↑ | LPIPS↓ | FID↓ | L1↓ | RMSE↓ | SSIM↑ | LPIPS↓ | FID↓ | |
| FUNIT | 0.08591 | 0.2529 | 0.6661 | 0.1169 | **11.66** | 0.09377 | 0.2686 | 0.6432 | 0.1427 | 28.10 | 11.74 |
| LF-Font | 0.08098 | 0.2435 | 0.6829 | 0.1226 | 27.73 | 0.09037 | 0.2620 | 0.6534 | 0.1448 | 38.46 | 13.01 |
| MX-Font | 0.07470 | 0.2319 | 0.7038 | 0.1034 | 18.75 | 0.08171 | 0.2468 | 0.6830 | 0.1193 | 27.91 | 10.86 |
| Fs-Font | 0.08214 | 0.2519 | 0.6657 | 0.1502 | 45.33 | 0.08917 | 0.2657 | 0.6467 | 0.1647 | 55.21 | 12.03 |
| CG-GAN | 0.07977 | 0.2409 | 0.6883 | 0.1117 | 23.93 | 0.08639 | 0.2549 | 0.6690 | 0.1303 | 37.22 | 16.67 |
| DG-Font | 0.06251 | 0.2105 | 0.7437 | 0.0846 | 17.10 | 0.07841 | 0.2442 | 0.6853 | 0.1198 | 27.98 | 14.11 |
| CF-Font | **0.05997** | **0.2053** | **0.7538** | **0.0836** | 13.13 | **0.07394** | **0.2354** | **0.7007** | **0.1182** | **26.51** | **21.58** |
| | (4.1%) | (2.5%) | (1.4%) | (1.1%) | (-) | (5.7%) | (3.6%) | (2.3%) | (0.92%) | (5.0%) | (29.5%) |



Figure 6. Qualitative comparison with state-of-the-art methods on Chinese poems. As mentioned earlier, we use multiple source fonts and pick the best results for these comparison methods for fairness. Here we just plot font *Song* as an example of source fonts for convenience. We mark erroneous skeletons with red boxes and other mismatch styles, such as stroke style, joined-up style, and body frame [26], with blue boxes.

## 4.2. Implementation Details

We train our model using Adam [19] with $\beta_1 = 0.9$ and $\beta_2 = 0.99$ for the style encoder, and RMSprop [12] with $\alpha = 0.99$ for the content encoder. The learning rate and weight decay are both set as $10^{-4}$. The hyper-parameters for loss are $\lambda_{img} = \lambda_{cnt} = 0.1$, $\lambda_{offset} = 0.5$, and $\lambda_{pcl} = 0.01$ (0.05 for PC-KL). For PCL, we orthographically project a character image onto 12 straight lines, which cross at the image center and divide the 2D space evenly. We resize all images to $80 \times 80$ and train the model with a batch size of 32. The whole training takes about 15 hours. We first train the DGN for 180k iterations to obtain the learned content features. Then we cluster these content features into ten groups and select basis fonts by the distance to cluster centers. After that, the model with CFM is further

trained for another 20k iterations. For fairness, the models without CFM in ablations are trained for 200k iterations.

## 4.3. Comparison with State-Of-The-Art Methods

We compare our model with six state-of-the-art methods, including an image-to-image translation method (FUNIT [24]), four component-related methods (LF-Font [31], MX-Font [32], CG-GAN [20], FsFont [36]), and DG-Font [42]. We slightly modify the network of CG-GAN to fit the input image size and the few-shot setting. To be fair, we try each of our basis fonts and font *Kai* as the source font for these comparison methods and report their best results in the following part (see details in our supplementary).

As Tbl. 1 illustrates, our method outperforms other methods, especially on unseen fonts. DG-Font leads other

Table 2. Ablation studies on different components. The first row is the result of DGN. P, C and S represent PC-WDL, CFM and ISR respectively. N means using retrieval strategy, *i.e.* picking the closet font from basis fonts (if marked with a star *, from the whole training set expect the target font itself) as the source according to the similarity between content features.

| Methods | | | | Seen Fonts | | | | | Unseen Fonts | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| P | C | S | N | L1↓ | RMSE↓ | SSIM↑ | LPIPS↓ | FID↓ | L1↓ | RMSE↓ | SSIM↑ | LPIPS↓ | FID↓ |
| | | | | 0.06251 | 0.2105 | 0.7437 | 0.0846 | 17.10 | 0.07841 | 0.2442 | 0.6853 | 0.1198 | 27.98 |
| ✓ | | | | 0.06261 | 0.2103 | 0.7434 | 0.0853 | 16.17 | 0.07803 | 0.2435 | 0.6868 | 0.1202 | 26.79 |
| ✓ | | | ✓ | 0.06727 | 0.2221 | 0.7240 | 0.0957 | 17.02 | 0.08009 | 0.2489 | 0.6786 | 0.1259 | 27.12 |
| ✓ | | | ✓* | **0.05952** | **0.2001** | **0.7552** | 0.0856 | 23.34 | 0.07519 | 0.2359 | 0.6984 | 0.1224 | 34.83 |
| ✓ | ✓ | | | 0.06056 | 0.2071 | 0.7506 | 0.0865 | 16.08 | 0.07574 | 0.2399 | 0.6940 | 0.1199 | 27.01 |
| ✓ | ✓ | ✓ | | 0.05997 | 0.2053 | 0.7538 | **0.0836** | **13.13** | **0.07394** | **0.2354** | **0.7007** | **0.1182** | **26.51** |



Figure 7. Qualitative results in the ablation on different components. P, C, and S are the same notations as Tbl. 2. We mark erroneous skeletons with red circles and other mismatch styles with blue circles.



Figure 8. Comparison between content fusion and retrieval strategy. B represents the baseline (DG-Font), and other notations are the same as Tbl. 2.

comparison methods except on perception metrics. But when added our proposed modules, its LPIPS and FID scores get a significant boost and catch up with others both on seen and unseen fonts. Although FUNIT achieves the best FID score on seen fonts, it performs worse on other metrics. Fig. 6 displays the qualitative comparison. Characters generated by ours are of high quality in terms of style consistency and structural correctness. The results of FU-NIT, LF-Font, MX-Font, Fs-Font, and CG-GAN often have structural errors and incompleteness. Fs-Font select several reference characters from a reference collection through a content-reference mapping, the relationship between a character and its references with common conspicuous components. The reference collection contains around 100 characters and covers almost all components. However, our reference characters are randomly selected and fixed for all source characters, with poor component coverage. Thus, the performance of Fs-Font is not perfectly shown in our few-shot setting. The outputs of DG-Font are great overall but suffer from artifacts and incomplete style transfer.

**User study.** We conduct a user study to further compare our model with other methods. We randomly selected 40 font styles (30 seen fonts and 10 unseen fonts) from the test set, and for each style, 5 test characters were randomly se-
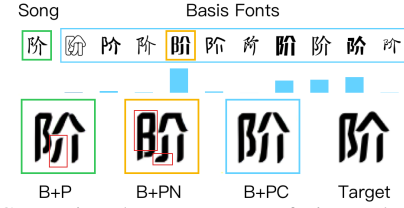
lected. Corresponding character images are generated with our method and the other 6 comparison methods. 20 participants who use Chinese characters every day are asked to pick the best group (5 character images yielded by one method) for one test style. Here, the order of these groups is randomly shuffled and we allow multiple choices since the participants might think several synthesized characters are of comparable quality. The results of user study are shown in the last column of 1, which present that our CF-Font gains the highest user preference 21.58%, surpassing the second place CG-GAN 16.67% by a large margin.

### 4.4. Ablation Studies

This subsection shows the effects of all proposed components and discusses how CFM and PCL work in font generation.

**Effectiveness of different components.** We separate the proposed modules and sequentially add them to DGN to observe the effects of each. The quantitative results can be seen in Tbl. 2, verifying that PCL, CFM, and ISR all can help improve the quality of generated images. These modules bring not only a numerical improvement but also a noticeable improvement in the visual aspect of geometric structures and stylistic strokes, as displayed in Fig. 7. In the fourth-to-last line, PCL shows its ability to improve character semantics and skeletons. Moreover, CFM makes the generated results a big step closer to the target in human perception. In the penultimate line, ISR further refines the detail of results by enhancing the stylistic representation.
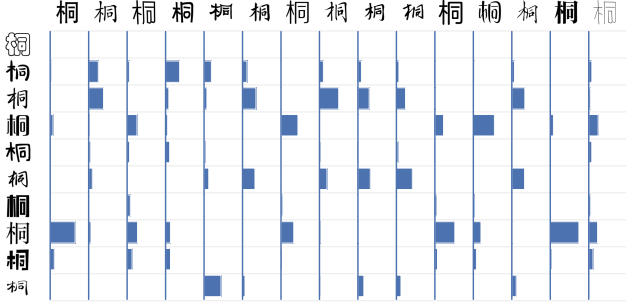
Figure 9. Visualization of weights on basis fonts. We take the character "Tong" for example. The left column represents the basis fonts, and the top row shows a part of training fonts. The weight on basis fonts of one training font are displayed as a vertical histogram.

**Comparison between content fusion and retrieval strategy.** Among these modules, CFM is the most efficient one. We further analyze where the gain of CFM comes from through a comparison with the retrieval strategy, *i.e.* during the test, we select the closet font for every target font as input from basis fonts and the whole training set (except the target font itself, *i.e.* 239/240 fonts in total for each seen/unseen target font) respectively. The quantitative result is shown in the second to fourth row from the bottom of Tbl. 2. It indicates that the result of inputting the closet basis font is much worse than that of content fusion, or even worse than the baseline (using a stand font *Song* all the time). Meanwhile, retrieving the closest font from the whole set gets a comparable results with CF-Font on seen fonts, but not good as it on unseen fonts and FID metrics. As Fig. 8 illustrates, the closet font may still be very different on the character skeleton from the target one and will introduce some noises (parts mismatched to the target skeleton). With these observations, we claim that content fusion matters rather than retrieving a close font in CFM.

**Variations of PCL.** We use two variations of PCL, PC-WDL, and PC-KL, to train a model respectively. Tbl. 3 shows the result on unseen fonts and demonstrates that not only PC-WDL, PC-KL can also improve the network performance. PC-KL and PC-WDL have similar improvements on pixel-level metrics, but PC-WDL has obvious advantages in FID while PC-KL performs better on LPIPS. We attribute this to that benefit from character projection, both of the distribution distance metrics can focus on the global properties, such as skeleton topology.

### 4.5. Evaluation of Basis Selection.

We visualize the basis fonts and the corresponding weights of content fusion here. Taking the character "Tong" as an example, in Fig. 9, ten images of basis fonts are shown in the left column, fifteen target images with randomly se-

Table 3. Quantitative evaluation using variations of PCL.

| method | L1 ↓ | RMSE ↓ | SSIM ↑ | LPIPS ↓ | FID ↓ |
|---|---|---|---|---|---|
| Baseline | 0.07841 | 0.2442 | 0.6853 | 0.1198 | 27.98 |
| +PC-KL | **0.07802** | **0.2434** | **0.6872** | **0.1191** | 27.72 |
| +PC-WDL | 0.07803 | 0.2435 | 0.6868 | 0.1202 | **26.79** |



Figure 10. Failure case.

lected fonts are listed in the top row, and the weights of content fusion are plotted in the form of a vertical histogram. We can observe that (a) the basis fonts selected by clustering are indeed visually different from each other (they also can be chosen manually), which means that they are capable of building a space for content fusion; (b) the greater the weight value, the corresponding basis font is more similar to the target font and this proves that content fusion is physically meaningful; (c) the values of these weights are scattered rather than concentrated in a particular basis font, which can also be a reason why the retrieval strategy fails as described in subsection 4.4.

### 4.6. Failure Cases and Limitations

Fig. 10 illustrates a case of generated images of complex characters with many strokes and a tight layout. Although our method works relatively well, many structural errors appear in the first row and some strokes are missed in the second row.

### 5. Conclusion

In this paper, we design a content fusion module and a projected character loss to improve the quality of skeleton transfer in few-shot font generation. We also propose a iterative style-vector refinement strategy to find a better font-level style representation. Experiments demonstrate that our method can outperform existing state-of-the-art methods, and each of the proposed novel modules is effective.

In the future, we may try vector font generation because vector characters are scale-invariant and more convenient for practical applications. It would be interesting to investigate whether the content fusion strategy can help solve the problem of complex vector font generation.

### Acknowledgments

# References

[1] Kyungjune Baek, Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Hyunjung Shim. Rethinking the truly unsupervised image-to-image translation. In *Int. Conf. Comput. Vis.*, pages 14134–14143. IEEE, 2021. 2

[2] Sagie Benaim and Lior Wolf. One-sided unsupervised domain mapping. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Adv. Neural Inform. Process. Syst.*, pages 752–762, 2017. 2

[3] Junbum Cha, Sanghyuk Chun, Gayoung Lee, Bado Lee, Seonghyeon Kim, and Hwalsuk Lee. Few-shot compositional font generation with dual memory. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Eur. Conf. Comput. Vis.*, volume 12364 of *Lecture Notes in Computer Science*, pages 735–751. Springer, 2020. 2

[4] Jie Chang, Yujun Gu, Ya Zhang, and Yan-Feng Wang. Chinese handwriting imitation with hierarchical generative adversarial network. In *Brit. Mach. Vis. Conf.*, page 290. BMVA Press, 2018. 2

[5] Xinyuan Chen, Chang Xu, Xiaokang Yang, Li Song, and Dacheng Tao. Gated-gan: Adversarial gated networks for multi-collection style transfer. *IEEE Trans. Image Process.*, 28(2):546–560, 2019. 2

[6] Yunjey Choi, Min-Je Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 8789–8797. Computer Vision Foundation / IEEE Computer Society, 2018. 2

[7] Yiming Gao and Jiangqin Wu. Gan-based unpaired chinese character image translation via skeleton transformation and stroke rendering. In *AAAI*, volume 34, pages 646–653, 2020. 2

[8] Yiming Gao and Jiangqin Wu. Gan-based unpaired chinese character image translation via skeleton transformation and stroke rendering. In *AAAI*, pages 646–653. AAAI Press, 2020. 2

[9] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In Zoubin Ghahramani, Max Welling, Corinna Cortes, Neil D. Lawrence, and Kilian Q. Weinberger, editors, *Adv. Neural Inform. Process. Syst.*, pages 2672–2680, 2014. 1, 2

[10] Zheng Gu, Wenbin Li, Jing Huo, Lei Wang, and Yang Gao. Lofgan: Fusing local representations for few-shot image generation. In *Int. Conf. Comput. Vis.*, pages 8443–8451. IEEE, 2021. 2

[11] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Adv. Neural Inform. Process. Syst.*, pages 6626–6637, 2017. 5

[12] Geoffrey Hinton. Neural networks for machine learning. https://www.coursera.org/learn/neural-networks/home/welcome. 6

[13] Xun Huang and Serge J. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Int. Conf. Comput. Vis.*, pages 1510–1519. IEEE Computer Society, 2017. 2, 3

[14] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with conditional adversarial networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5967–5976. IEEE Computer Society, 2017. 2

[15] Yue Jiang, Zhouhui Lian, Yingmin Tang, and Jianguo Xiao. Dcfont: an end-to-end deep chinese font generation system. In Diego Gutierrez and Hui Huang, editors, *SIGGRAPH Asia 2017 Technical Briefs, Bangkok, Thailand, November 27 - 30, 2017*, pages 22:1–22:4. ACM, 2017. 2

[16] Yue Jiang, Zhouhui Lian, Yingmin Tang, and Jianguo Xiao. Scfont: Structure-guided chinese font generation via deep stacked networks. In *AAAI*, pages 4015–4022. AAAI Press, 2019. 2

[17] Hadi Kazemi, Seyed Mehdi Iranmanesh, and Nasser Nasrabadi. Style and content disentanglement in generative adversarial networks. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 848–856. IEEE, 2019. 2

[18] Taeksoo Kim, Moonsu Cha, Hyunsoo Kim, Jung Kwon Lee, and Jiwon Kim. Learning to discover cross-domain relations with generative adversarial networks. In Doina Precup and Yee Whye Teh, editors, *Int. Conf. Mach. Learn.*, volume 70 of *Proceedings of Machine Learning Research*, pages 1857–1865. PMLR, 2017. 2

[19] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *Int. Conf. Learn. Represent.*, 2015. 6

[20] Yuxin Kong, Canjie Luo, Weihong Ma, Qiyuan Zhu, Shenggao Zhu, Nicholas Yuan, and Lianwen Jin. Look closer to supervise better: One-shot font generation via component-based discriminator. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 13482–13491, 2022. 2, 5, 6

[21] Gihyun Kwon and Jong Chul Ye. Diagonal attention and style-based gan for content-style disentanglement in image generation and translation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 13980–13989, 2021. 2

[22] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation applied to handwritten zip code recognition. *Neural Computation*, 1(4):541–551, 1989. 1

[23] Jae Hyun Lim and Jong Chul Ye. Geometric GAN. *CoRR*, abs/1705.02894, 2017. 3

[24] Ming-Yu Liu, Xun Huang, Arun Mallya, Tero Karras, Timo Aila, Jaakko Lehtinen, and Jan Kautz. Few-shot unsupervised image-to-image translation. In *Int. Conf. Comput. Vis.*, pages 10550–10559. IEEE, 2019. 2, 6

[25] Wei Liu, Fangyue Liu, Fei Ding, Qian He, and Zili Yi. Xmpfont: Self-supervised cross-modality pre-training for few-shot font generation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 7905–7914, 2022. 2

[26] Xiaoqing Lu and Ting Tang. Elements of chinese typeface design. In *Digital Fonts and Reading*, pages 109–130. World Scientific, 2016. 6

[27] Ken Lunde. *CJKV Information Processing: Chinese, Japanese, Korean & Vietnamese Computing*. O'Reilly, 1999. 5

[28] Pengyuan Lyu, Xiang Bai, Cong Yao, Zhen Zhu, Tengteng Huang, and Wenyu Liu. Auto-encoder guided gan for chinese calligraphy synthesis. In *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, volume 1, pages 1095–1100. IEEE, 2017. 2

[29] Pengyuan Lyu, Xiang Bai, Cong Yao, Zhen Zhu, Tengteng Huang, and Wenyu 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*, pages 1095–1100. IEEE, 2017. 2

[30] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *Int. Conf. Learn. Represent.*, 2018. 3

[31] Song Park, Sanghyuk Chun, Junbum Cha, Bado Lee, and Hyunjung Shim. Few-shot font generation with localized style representations and factorization. In *AAAI*, pages 2393–2402. AAAI Press, 2021. 2, 4, 5, 6

[32] Song Park, Sanghyuk Chun, Junbum Cha, Bado Lee, and Hyunjung Shim. Multiple heads are better than one: Few-shot font generation with multiple localized experts. In *Int. Conf. Comput. Vis.*, pages 13880–13889. IEEE, 2021. 2, 6

[33] Ashish Shrivastava, Tomas Pfister, Oncel Tuzel, Joshua Susskind, Wenda Wang, and Russell Webb. Learning from simulated and unsupervised images through adversarial training. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2242–2251. IEEE Computer Society, 2017. 2

[34] Danyang Sun, Tongzheng Ren, Chongxuan Li, Hang Su, and Jun Zhu. Learning to write stylized chinese characters by reading a handful of examples. In Jérôme Lang, editor, *IJCAI*, pages 920–927. ijcai.org, 2018. 2

[35] Donghui Sun, Qing Zhang, and Jun Yang. Pyramid embedded generative adversarial network for automated font generation. In *Int. Conf. Pattern Recog.*, pages 976–981. IEEE Computer Society, 2018. 2

[36] Licheng Tang, Yiyang Cai, Jiaming Liu, Zhibin Hong, Mingming Gong, Minhu Fan, Junyu Han, Jingtuo Liu, Errui Ding, and Jingdong Wang. Few-shot font generation by learning fine-grained local styles. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 7895–7904, June 2022. 3, 6

[37] Michael C Thrun. The exploitation of distance distributions for clustering. *International Journal of Computational Intelligence and Applications*, 20(03):2150016, 2021. 4

[38] Yuchen Tian. Zi2zi. https://github.com/kaonashi-tyc/zi2zi. 2

[39] Sergei Vassilvitskii and David Arthur. k-means++: The advantages of careful seeding. In *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 1027–1035, 2006. 4

[40] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 8798–8807. Computer Vision Foundation / IEEE Computer Society, 2018. 2

[41] Shan Jean Wu, Chih-Yuan Yang, and Jane Yung-jen Hsu. Calligan: Style and structure-aware chinese calligraphy character generator. *CoRR*, abs/2005.12500, 2020. 2

[42] Yangchen Xie, Xinyuan Chen, Li Sun, and Yue Lu. Dg-font: Deformable generative networks for unsupervised font generation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5130–5140. Computer Vision Foundation / IEEE, 2021. 2, 3, 5, 6

[43] Songhua Xu, Hao Jiang, Tao Jin, Francis CM Lau, and Yunhe Pan. Automatic generation of chinese calligraphic writings with style imitation. *IEEE Intelligent Systems*, 24(02):44–53, 2009. 2

[44] Raymond A Yeh, Chen Chen, Teck Yian Lim, Alexander G Schwing, Mark Hasegawa-Johnson, and Minh N Do. Semantic image inpainting with deep generative models. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5485–5493, 2017. 5

[45] Han Zhang, Ian J. Goodfellow, Dimitris N. Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *Int. Conf. Mach. Learn.*, volume 97, pages 7354–7363, 2019. 3

[46] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 586–595. Computer Vision Foundation / IEEE Computer Society, 2018. 5

[47] Yexun Zhang, Ya Zhang, and Wenbin Cai. Separating style and content for generalized style transfer. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 8447–8455, 2018. 2, 5

[48] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Int. Conf. Comput. Vis.*, pages 2242–2251. IEEE Computer Society, 2017. 2