# Multiple Heads are Better than One:
# Few-shot Font Generation with Multiple Localized Experts

Song Park[1]    Sanghyuk Chun[2,3]    Junbum Cha[3]    Bado Lee[3]    Hyunjung Shim[1]

[1] School of Integrated Technology, Yonsei University    [2] NAVER AI Lab    [3] NAVER CLOVA

## Abstract

*A few-shot font generation (FFG) method has to satisfy two objectives: the generated images should preserve the underlying global structure of the target character and present the diverse local reference style. Existing FFG methods aim to disentangle content and style either by extracting a universal representation style or extracting multiple component-wise style representations. However, previous methods either fail to capture diverse local styles or cannot be generalized to a character with unseen components, e.g., unseen language systems. To mitigate the issues, we propose a novel FFG method, named Multiple Localized Experts Few-shot Font Generation Network (MX-Font). MX-Font extracts multiple style features not explicitly conditioned on component labels, but automatically by multiple experts to represent different local concepts, e.g., left-side sub-glyph. Owing to the multiple experts, MX-Font can capture diverse local concepts and show the generalizability to unseen languages. During training, we utilize component labels as weak supervision to guide each expert to be specialized for different local concepts. We formulate the component assign problem to each expert as the graph matching problem, and solve it by the Hungarian algorithm. We also employ the independence loss and the content-style adversarial loss to impose the content-style disentanglement. In our experiments, MX-Font outperforms previous state-of-the-art FFG methods in the Chinese generation and cross-lingual, e.g., Chinese to Korean, generation. Source code is available at https://github.com/clovaai/mxfont.*

## 1. Introduction

A few-shot font generation task (FFG) [42, 54, 12, 41, 6, 7, 37] aims to generate a new font library using only a few reference glyphs, *e.g.*, less than 10 glyph images, without additional model fine-tuning at the test time. FFG is especially a desirable task when designing a new font library
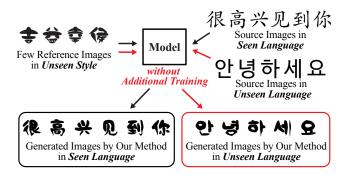


Figure 1. **Cross-lingual few-shot font generation results by MX-Font.** With only four references, the proposed method, MX-Font, can generate a high quality font library. Furthermore, we first show the effectiveness of the proposed method on the *zero-shot cross-lingual* few-shot generation task, *i.e.*, generating unseen Korean glyphs using the Chinese font generation model.

for glyph-rich scripts, *e.g.*, Chinese ($> 50$K glyphs), Korean ($\approx 11$K glyphs), or Thai ($\approx 11$K glyphs). It is because the traditional font design process is very labor-intensive due to the complex characteristics of the font domain. Another real-world scenario of FFG is to extend an existing font design to different language systems. For example, an international multi-media content, such as a video game or movie designed with a creative font, is required to re-design coherent style fonts for different languages.

A high-quality font design is obliged to satisfy two objectives. First, the generated glyph should maintain all the detailed structure of the target character, particularly important for glyph-rich scripts with highly complex structure. For example, even very small damages on a local component of a Chinese glyph can hurt the meaning of the target character. As another objective, a generated glyph should have a diverse local style of the reference glyphs, *e.g.*, serifness, strokes, thickness, or size. To achieve both objectives, existing methods formulate FFG by disentangling the content information and the style information from the given glyphs [42, 54, 12, 6, 37]. They combine the content features from the source glyph and the style features from the reference glyphs to generate a glyph with the reference
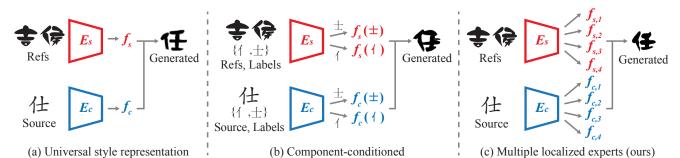
Figure 2. **Comparison of FFG methods.** Three different groups of FFG are shown. All methods combine style representation $f_s$ from a few reference glyphs (Refs) by a style encoder ($E_s$) and content representation $f_c$ from a source glyph (Source) by a content encoder ($E_c$). (a) Universal style representation methods extract only a single style feature for each font. (b) Component-conditioned methods extract component conditioned style features to capture diverse local styles (c) Multiple localized experts method (ours) generates multiple local features without an explicit condition, but attends different local information of the complex input glyph. The generated images in (a), (b) and (c) are synthesized by AGIS-Net [12], LF-Font [37] and MX-Font, respectively.

style. Due to the complex nature of the font domain, the major challenge of FFG is to correctly disentangle the global content structure and the diverse local styles. However, as shown in our experiments, we observe that existing methods are insufficient to capture diverse local styles or to preserve the global structures of unseen language systems.

We categorize existing FFG methods into *universal style representation methods* [42, 54, 34, 12] and *component-conditioned methods* [6, 37]. Universal style representation methods [42, 54, 34, 12] extract only a single style representation for each style – see Figure 2 (a). As glyph images are highly complex, these methods often fail to capture diverse local styles. To address the issue, component-conditioned methods [6, 37] utilize *compositionality*; a character can be decomposed into a number of sub-characters, or *components* – see Figure 2 (b). They explicitly extract component-conditioned features, beneficial to preserve the local component information. Despite their promising performances, their encoder is tightly coupled with specific component labels of the target language domain, which hinders processing the glyphs with unseen components or conducting a cross-lingual font generation.

In this paper, we propose a novel few-shot font generation method, named **M**ultiple Localized e**X**perts Few-shot **Font** Generation Network (MX-Font), which can capture multiple local styles, but not limited to a specific language system. MX-Font has a multi-headed encoder, named *multiple localized experts*. Each localized expert is specialized for different local sub-concepts from the given complex glyph image. Unlike component-conditioned methods, our experts are not explicitly mapped to a specific component, but each expert implicitly learns different local concepts by weak supervision *i.e.* component and style classifiers. To prevent that different experts learn the same local component, we formulate the component label allocation problem as a graph matching problem, optimally solved by the Hungarian algorithm [29] (Figure 4). We also employ the *in-*

*dependence loss* and the *content-style adversarial loss* to enforce the content-style disentanglement by each localized expert. Interestingly, with only weak component-wise supervision (*i.e.* image-level not pixel-level labels), we observe that each localized expert is specialized for different local areas, *e.g.*, attending the left-side of the image (Figure 7). While we inherit the advantage of component-conditioned methods [6, 37] by introducing the multiple local features, our method is not limited to a specific language by removing the explicit component dependency in extracting features. Consequently, MX-Font outperforms the state-of-the-art FFG in two scenarios: *In-domain transfer scenario*, training on Chinese fonts and generating an unseen Chinese font, and *zero-shot cross-lingual transfer scenario*, training on Chinese fonts and generating a Korean font. Our ablation and model analysis support that the proposed modules and optimization objectives are important to capture multiple diverse local concepts.

## 2. Related Works

**Style transfer and image-to-image translation.** Few-shot font generation can be viewed as a task that transfers reference font style to target glyph. However, style transfer methods [14, 21, 31, 35, 32, 47] regard the texture or color as a style while in font generation scenario, a style is often defined by a local shape, *e.g.*, stroke, size, or serif-ness. On the other hand, image-to-image translation (I2I) methods [23, 56, 9, 33, 48, 10] learn the mapping between domains from the data instead of defining the style. For example, FUNIT [34] aims to translate an image to the given reference style while preserving the content. Many FFG methods, thus, are based on I2I framework.

**Many-shot font generation methods.** Early font generation methods, such as zi2zi [43], aim to train the mapping between different font styles. A number of font generation methods [24, 13, 22, 45] first learn the mapping func-
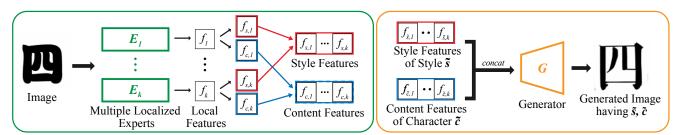
Figure 3. **Overview of MX-Font.** Two modules of MX-Font used for the generation are described. The *multiple localized experts* (green box) consist of $k$ experts. $E_i$ (*i.e.* $i$-th expert) encodes the input image to the $i$-th local feature $f_i$ and the $i$-th style and content feature $f_{s,i}$, $f_{c,i}$ are computed from $f_i$. The right yellow box shows how the generator $G$ generates the target image. When $k$ style features representing the target style $\widetilde{s}$ and $k$ content features representing the target style $\widetilde{c}$ are given, the target glyph having style $\widetilde{s}$ and character $\widetilde{c}$ is generated by passing the element-wisely concatenated style and content features to the $G$.

tion, and fine-tune the mapping function for many reference glyphs, *e.g.* 775 [24]. Despite their remarkable performances, their scenario is not practical because collecting hundreds of glyphs with a coherent style is too expensive. In this paper, we aim to generate an unseen font library without any expensive fine-tuning and collecting a large number of reference glyphs for a new style.

**Few-shot font generation methods.** Since font styles are highly complex and fine-grained, utilizing statistical textures as style transfer is challenging. Instead, the majority of FFG methods aims to disentangle font-specific style and content information from the given glyphs [54, 41, 2, 12, 42, 30]. We categorize existing FFG methods into two different categories. The universal style representation methods, such as EMD [54], AGIS-Net [12], synthesize a glyph by combining the style vector extracted from the reference set, and the content vector extracted from the source glyph. MX-Font employs multiple styles, and does not rely on the font specific loss design, *e.g.*, the local texture refinement loss by AGIS-Net [12]. However, the universal style representation shows limited performances in capturing localized styles and content structures. To address the issue, *component-conditioned methods* such as DM-Font [6], LF-Font [37], remarkably improve the stylization performance by employing localized style representation, where the font style is described multiple localized styles instead of a single universal style. However, these methods require explicit component labels (observed during training) for the target character even at the test time. This property limits practical usages such as cross-lingual font generation. Our method inherits the advantages from component-guided multiple style representations, but does not require the explicit labels at the test time.

## 3. Method

We introduce a novel few-shot font generation method, namely Multiple Localized Experts Few-shot Font Generation Network (MX-Font). MX-Font has a multi-headed encoder called *multiple localized experts*, where $i$-th head (or

expert $E_i$) encodes a glyph image $x$ into a local feature $f_i = E_i(x)$ (§3.1). We induce each expert $E_i$ to attend different local concepts, guided by a set of component labels $U_c$ for the given character $c$ (§3.2). From $f_i$, we compute a local content feature $f_{c,i}$ and a local style feature $f_{s,i}$ (§3.3). Once MX-Font is trained, we generate a glyph $\widetilde{x}$ with a character label $\widetilde{c}$ and a style label $\widetilde{s}$ by combining expert-wise features $f_{\widetilde{c},i}$ and $f_{\widetilde{s},i}$, from the source glyph and the reference glyph, respectively. (§3.5).

### 3.1. Model architecture

Our method consists of three modules; 1) $k$-headed encoder, or localized experts $E_i$, 2) a generator $G$, and 3) style and component feature classifiers $Cls_s$ and $Cls_u$. We illustrate the overview of our method in Figure 3 and Figure 5. We provide the details of the building blocks in the supplementary materials.

The green box in Figure 3 shows how the *multiple localized experts* works. The **localized expert** $E_i$ encodes a glyph image $x$ into a local feature $f_i = E_i(x) \in \mathbb{R}^{d \times w \times h}$, where $d$ is a feature dimension, and $\{w, h\}$ are spatial dimensions. By multiplying two linear weights $W_{i,c}, W_{i,s} \in \mathbb{R}^{d \times d}$ to $f_i$, a local content feature $f_{c,i} = W_{i,c}^\top f_i$ and a local style feature $f_{s,i} = W_{i,s}^\top f_i$ are computed. Here, our localized experts are not supervised by component labels to obtain $k$ local features $f_1, \ldots, f_k$; our local features are not component-specific features. We set the number of the localized experts, $k$, to 6 in our experiments if not specified.

We employ two **feature classifiers**, $Cls_s$ and $Cls_u$ to supervise $f_{s,i}$ and $f_{c,i}$, which serve as weak supervision for $f_i$. The classifiers are trained to predict the style (or component) labels, thereby $E_i$ receives the feedback from the $Cls_s$ and $Cls_u$ that $f_{s,i}$ and $f_{c,i}$ should preserve label information. These classifiers are only used during training but independent to the model inference itself. Following the previous methods [6, 37], we use font library labels for style labels $y_s$, and the component labels $U_c$ for content labels $y_c$. The example of component labels is illustrated in Figure 4. The same decomposition rule used by LF-Font [37]
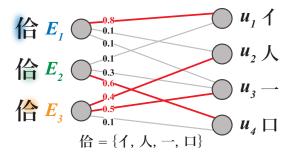
$伲 = \{イ, 人, 一, 口\}$

Figure 4. **An example of localized experts.** The number of experts $k$ is three ($E_1, E_2, E_3$), and the number of target component labels $m$ is four ($u_1, \ldots, u_4$). An edge between an expert $E_i$ and a component $u_j$ means the prediction probability of $u_j$ by $E_i$ using the component classifier $Cls_u$. Our goal is to find a set of edges that maximizes the sum of predictions, where the number of the selected edges are upper bounded by $\max(k, m) = 4$ in this example. The red edges illustrate the optimal solution.

is adopted. While previous methods only use the style (or content) classifier to train style (or content), we additionally utilize them for the content and style disentanglement by introducing the content-style adversarial loss.

The **generator** $G$ synthesizes a glyph image $\widetilde{x}$ by combining content and style features from each expert:

$$\widetilde{x} = G((f_{s,1} \circ f_{c,1}), \ldots, (f_{s,k} \circ f_{c,k})), \quad (1)$$

where $\circ$ denotes a concatenate operation.

In the following, we describe the details of each module, training settings, and how to generate samples with only a few references.

### 3.2. Learning multiple localized experts with weak local component supervision

Our intuition is that extracting different localized features can help each local feature to represent the detailed local structure and fine-grained local style in a complex glyph image. We utilize the compositionality of the font domain to inherit the advantages of component-conditioned methods [6, 37]. Meanwhile, we intentionally remove the explicit component dependency of the feature extractor for achieving generalizability, which is the weakness of previous methods. Here, we employ a multi-headed feature extractor, named *multiple localized experts*, where each expert can be specialized for different local concepts. A naïve solution is to utilize explicit local supervision, *i.e.*, the pixel-level annotation for each sub-glyph, unable to obtain due to expensive annotation cost. As an alternative, a strong machine annotator can be utilized to obtain local supervision [50], but training a strong model, such as the self-trained EfficientNet L2 with 300M images [46], for the font domain is another challenge that is out of our scope.

Utilizing the compositionality, we have the weak component-level labels for the given glyph image, *i.e.*,

what components the image has but without the knowledge where they are, similar to the multiple instance learning scenario [36, 55]. Then, we let each expert attend on different local concepts by guiding each expert with the component and style classifiers. Ideally, when the number of components $m$ is same as the number of experts, $k$, we expect the $k$ predictions by experts are same as the component labels, and the summation of their prediction confidences is maximized. When $k < m$, we expect the predictions by each expert are "plausible" by considering top-k predictions.

To visualize the role of each expert, we illustrate an example in Figure 4. Presuming three multiple experts, they can learn different local concepts such as the left-side (blue), the right-bottom-side (green), and the right-upper-side (yellow), respectively. Given a glyph composed of four components, the feature from each expert can predict one ($E_1, E_2$) or two ($E_3$) labels as shown in the figure. Because we do not want that an expert is explicitly assigned to a component label, *e.g.*, strictly mapping "人" component to $E_1$, we solve an automatic allocation algorithm, finding the optimal expert-component matching as shown in Figure 4. Specifically, we formulate the component allocation problem as the Weighted Bipartite B-Matching problem, which can be optimally solved by the Hungarian algorithm [29].

From a given glyph image $x$, each expert $E_i$ extracts the content feature $f_{c,i}$. Then, the component feature classifier $Cls_u$ takes $f_{c,i}$ as input and produces the prediction probability $p_i = Cls_u(f_{c,i})$, where $p_i = [p_{i0}, \ldots, p_{im}]$ and $p_{ij}$ is the confidence scalar value of the component $j$. Let $U_c = \{u_1^c, \ldots, u_m^c\}$ be a set of component labels of the given character $c$, and $m$ be the number of the components. We introduce an allocation variable $w_{ij}$, where $w_{ij} = 1$ if the component $j$ is assigned to $E_i$, and $w_{ij} = 0$ otherwise. We optimize the binary variables $w_{ij}$ to maximize the summation over the selected prediction probability such that the number of total allocations is $\max(k, m)$. Now, we formulate the component allocation problem as:

$$\max_{w_{ij} \in \{0,1\} | i=1 \ldots k, j \in U_c} \sum_{i=1}^{k} \sum_{j \in U_c} w_{ij} p_{ij},$$

$$\text{s.t.} \quad \sum_{i=1}^{k} w_{ij} \geq 1 \text{ for } \forall j, \quad \sum_{j \in U_c} w_{ij} \geq 1 \text{ for } \forall i, \quad (2)$$

$$\sum_{i=1}^{k} \sum_{j \in U_c} w_{ij} = \max(k, m),$$

where (2) can be reformulated to the Weighted Bipartite B-Matching (WBM) problem, and can be solved by the Hungarian algorithm in a polynomial time $O((m + k)^3)$. We describe the connection between (2) and WBM in the supplementary materials. Now, using the estimated variables $w_{ij}$ in (2), we optimize auxiliary component classification
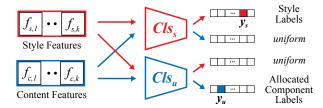
Figure 5. **Feature classifiers.** Two feature classifiers, $Cls_s$ and $Cls_u$ are used during the training. $Cls_s$ classifies the style features to their style label $y_s$ while $Cls_u$ predicts the uniform probability from them. Similarly, $Cls_u$ classifies the content features to their allocated component labels $y_u$ while $Cls_s$ is fooled by them. The details are described in § 3.2 and § 3.3.

loss $\mathcal{L}_{cls,c}$ with the cross entropy loss (CE) as follows:

$$\mathcal{L}_{cls,c,i}(f_{c,i}, U_c) = \sum_{j \in U_c} w_{ij} \text{CE}(Cls_u(f_{c,i}), j). \quad (3)$$

Here, we expect that each localized expert is specialized for a specific local concept so that it facilitates the content-style disentanglement. Because the feedback from (3) encourages the local features to be better separated into the style and content features, we expect that each expert automatically attends local concepts. We empirically observe that each expert is involved to different local areas without explicit pixel-level supervision (Figure 7).

We additionally formulate the independence between each expert by the Hilbert-Schmidt Independence Criterion [16] which has been used in practice for statistical testing [16, 17], feature similarity measurement [28], and model regularization [38, 51, 3]. HSIC is zero if and only if two inputs are independent of each other. Since HSIC is non-negative, the independence criterion can be achieved by minimizing HSIC. Under this regime, we use HSIC and lead the local feature $f_i$ extracted by $E_i$ independent to the other local features $f_{i'}$ as follows:

$$\mathcal{L}_{\text{indp exp},i} = \sum_{i'=1, i' \neq i}^{k} \text{HSIC}(f_i, f_{i'}). \quad (4)$$

We leave the detailed HSIC formulation is in the supplementary materials.

### 3.3. Content and style disentanglement

To achieve perfect content and style disentanglement, the style (or content) features should include the style (or content) domain information but exclude the content (or style) domain information. We employ two objective functions for this: *content-style adversarial loss* and *independent loss*.

The **content-style adversarial loss**, motivated by the domain adversarial network [11], enforces the extracted features for style (or content) is useless to classify content (or style). Thus, a style feature $f_{s,i}$ is trained to satisfy (1) correctly classify a style label $y_s$ by the style classifier $Cls_s$

with the cross entropy loss (CE) and (2) fooling the content labels predicted by the component classifier $Cls_u$. Specifically, we maximize the entropy ($H$) of the predicted probability to enforce the uniform prediction. Formally, we define our objective function for a style feature $f_{s,i}$ as follows:

$$\mathcal{L}_{s,i}(f_{s,i}, y_s) = \text{CE}(Cls_s(f_{s,i}), y_s) - H(Cls_u(f_{s,i})). \quad (5)$$

We define $\mathcal{L}_{c,i}$ as the objective function for a content feature $f_{c,i}$ employs $\mathcal{L}_{cls,c,i}$ (3) instead of the cross entropy of $y_c$ as follows:

$$\mathcal{L}_{c,i}(f_{c,i}, U_c) = \mathcal{L}_{cls,c,i}(f_{c,i}, U_c) - H(Cls_s(f_{c,i})). \quad (6)$$

We also employ the independence loss between content and style local features, $f_{c,i}$ and $f_{s,i}$ for the disentanglement of content and style representations. That is:

$$\mathcal{L}_{\text{indp},i} = \text{HSIC}(f_{s,i}, f_{c,i}). \quad (7)$$

### 3.4. Training

We train our model to synthesize a glyph image from the given content and style labels using the Chinese font dataset (details in §4.2). More specifically, we construct a mini-batch, where $n$ glyphs share the same content label $y_c$ (from random styles), and $n$ glyphs share the same style label $y_s$ (from random contents). Then, we let the model generate a glyph with the content label $y_c$ and the style label $y_s$. In our experiments, we set $n = 3$ and synthesize 8 different glyphs in parallel, *i.e.*, the mini-batch size is 24.

We employ a discriminator module $D$ and the generative adversarial loss [15] to achieve high-quality visual samples. In particular, we use the hinge generative adversarial loss $\mathcal{L}_{adv}$ [52], feature matching loss $\mathcal{L}_{fm}$, and pixel-level reconstruction loss $\mathcal{L}_{recon}$ by following the previous high fidelity GANs, *e.g.*, BigGAN [4], and state-of-the-art font generation methods, *e.g.*, DM-Font [6] or LF-Font [37]. The details of each objective function are in the supplementary materials.

Now we describe our full objective function. The entire model is trained in an end-to-end manner with the weighted sum of all losses, including (4), (5), (6), and (7).

$$\begin{aligned} \mathcal{L}_D &= \mathcal{L}_{adv}^D, \\ \mathcal{L}_G &= \mathcal{L}_{adv}^G + \lambda_{recon}\mathcal{L}_{recon} + \mathcal{L}_{fm} \\ \mathcal{L}_{exp} &= \sum_{i=1}^{k} [\mathcal{L}_{s,i} + \mathcal{L}_{c,i} + \mathcal{L}_{\text{indp},i} + \mathcal{L}_{\text{indp exp},i}] \end{aligned} \quad (8)$$

As conventional GAN training, we alternatively update $\mathcal{L}_D$, $\mathcal{L}_G$, and $\mathcal{L}_{exp}$. The control parameter $\lambda_{recon}$ is set to 0.1 in our experiments. We use Adam optimizer [26], and run the optimizer for 650k iterations. We additionally provide the detailed training setting is in the supplementary materials.

### 3.5. Few-shot generation

When the source and a few reference glyphs are given, MX-Font extract the content features from the source glyphs and the style features from the reference glyphs. Assume we have $n_r$ number of reference glyphs $x_1^r, \ldots, x_{n_r}^r$ with a coherent style $y_{s^r}$. First, our multiple experts $\{E_1, \ldots, E_k\}$ extract localized style feature $[f_{s^r,i}^1, \ldots, f_{s^r,i}^{n^r}]$ for $i = 1 \ldots k$ from the reference glyphs. Then, we take an average over the localized features to represent a style representation, $i.e.$, $f_{s^r,i} = \frac{1}{n^r} \sum_{j=1}^{n^r} f_{s^r,i}^j$ for $i = 1 \ldots k$. Finally, the style representation is combined with the content representation extracted from the known source glyph to generate unseen style glyph.

## 4. Experiments

In this section, we describe the evaluation protocols, and experimental settings. We extend previous FFG benchmarks to unseen language domain to measure the generalizability of a model. MX-Font is compared with four FFG methods on the proposed extended FFG benchmark via both the qualitative and quantitative evaluations. Experimental results demonstrate that MX-Font outperforms existing methods in the most of evaluation metrics. The ablation and analysis study helps understand the role and effects of our multiple experts and objective functions.

### 4.1. Comparison methods

**Universal style representation methods. EMD** [54] adopts content and style encoders that extract universal content and style features from a few reference glyphs. **AGIS-Net** [12] proposes the local texture refinement loss to handle unbalance between the number of positive and negative samples. **FUNIT** [34] is not directly proposed for FFG task, but we employ the modified version of FUNIT as our comparison method following previous works [6, 37].
**Component-conditioned methods. DM-Font** [6] learns two embedding codebooks (or the dual-memory) conditioned by explicit component labels. When the target character contains a component either unseen during training or not in the reference set, DM-Font is unable to generate a glyph. As these drawbacks are impossible to be fixed with only minor modifications, we do not compate DM-Font to MX-Font. **LF-Font** [37] relaxes the restriction of DM-Font by estimating missing component features via factorization module. Although LF-Font is still not applicable to generate a character with unseen components, we slightly modify LF-Font (as described in the supplementary materials) and compare the modified version with other methods.

### 4.2. Evaluation protocols

To show the generalizability to the unseen language systems, we propose an extended FFG scenario; training a FFG model on a language system and evaluating the model on the other language system. In this paper, we first train FFG models on the Chinese font dataset, and evaluate them on both Chinese generation (*in-domain transfer scenario*) and Korean generation (*zero-shot cross-lingual scenario*).

**Dataset.** We use the same Chinese font dataset collected by Park *et al*. [37] for training. The dataset contains 3.1M Chinese glyph images with 467 different styles, and 19,514 characters are covered. We also use the same decomposition rule as Park *et al*. [37] to extract component labels. We exclude 28 fonts, and 214 Chinese characters from the training set, and use them to evaluation. For the Korean FFG evaluation, we use the same test characters with Cha *et al*. [6], 245 characters. To sum up, we evaluate the methods by using 28 font styles with 214 Chinese and 245 Korean characters.

**Evaluation metrics.** Due to the style of the font domain is defined by a local fine-grained shape, *e.g*., stroke, size, or serif-ness, measuring the visual quality with a unified metric is a challenging problem. A typical challenge is the multiplicity of the font styles; because the font style is defined locally, there could be multiple plausible glyphs satisfying our objectives. However, we only have one "ground truth" glyphs in the test dataset. Furthermore, for the Korean generation task with Chinese references, we even do not have "ground truth" Korean glyphs with the reference styles. Thus, we need to employ evaluation metrics that does not require ground truth, and can evaluate plausibility of the given samples. We therefore use four different evaluation metrics to measure the visual quality in various viewpoints.

Following previous works [6, 37], we train evaluation classifiers that classifies character labels (content-aware) and font labels (style-aware). Note that these classifiers are only used for evaluation, and trained separately to the FFG models. We train three classifiers, the style classifier on the Chinese test fonts, the content classifier on the Chinese test characters, and the content classifier on the Korean test characters. The details of the evaluation classifiers are in the supplementary materials. Using the classifiers, we measure the **classification accuracies** for style and content labels. We also report the accuracy when both classifiers are correctly predicted.

We conduct a **user study** for quantifying the subjective quality. The participants are asked to pick the three best results, considering the style, the content, and the most preferred considering both the style and the content. All 28 test styles with 10 characters are shown to the participants. For each test style, we show Chinese and Korean samples separately to the users. *I.e*., a participant picks $28 \times 3 \times 2 = 168$ results. We collect the responses from 57 participants. User study samples are in the supplementary materials.

We also report LPIPS [53] scores to measure the dissimilarity between the generated images and their corresponding ground truth images, thus it is only reported for Chinese

| | | Acc (S) % | Acc (C) % | Acc (B) % | User (S) % | User (C) % | User (B) % | LPIPS ↓ | FID (H) ↓ |
|---|---|---|---|---|---|---|---|---|---|
| CN → CN | EMD (CVPR'18) | 6.6 | 51.3 | 4.6 | 0.7 | 0.1 | 0.3 | 0.212 | 79.7 |
| | AGIS-Net (TOG'19) | 25.5 | **99.5** | 25.4 | 22.4 | **34.2** | 26.8 | 0.124 | 19.2 |
| | FUNIT (ICCV'19) | 34.0 | 94.6 | 31.8 | 22.9 | 21.6 | 22.2 | 0.147 | 19.2 |
| | LF-Font (AAAI'21) | 58.7 | 96.9 | 57.0 | 19.5 | 12.3 | 15.6 | **0.119** | **14.8** |
| | MX-Font (proposed) | **78.9** | **99.5** | **78.7** | **34.5** | 31.8 | **35.2** | 0.120 | 21.8 |
| KR → CN | EMD (CVPR'18) | 4.6 | 15.4 | 0.8 | 0.8 | 0.1 | 0.1 | - | 150.1 |
| | AGIS-Net (TOG'19) | 13.3 | 32.1 | 3.1 | 1.8 | 0.6 | 0.6 | - | 146.5 |
| | FUNIT (ICCV'19) | 11.3 | 66.4 | 6.6 | 12.0 | 17.3 | 9.1 | - | 176.0 |
| | LF-Font (AAAI'21) | 47.6 | 28.7 | 12.8 | 10.6 | 0.7 | 1.0 | - | 148.7 |
| | MX-Font (proposed) | **66.3** | **75.9** | **50.0** | **74.6** | **81.3** | **89.2** | - | **84.1** |

Table 1. **Performance comparison on few-shot font generation scenario.** The performances of five few-shot font generation methods with four reference images are compared. We report accuracy measured by style-aware (Acc (S)) and content-aware (Acc (C)) classifiers and accuracy considering both the style and content labels (Acc (B)). The summarized results of the user study are also reported. The User preference on considering style (User (S)), content (User (C)), both of them (User (B)) are shown. LPIPS shows a perceptual dissimilarity between the ground truth and the generated glyphs. The harmonic mean (H) of style-aware and content-aware FID is also reported. Note that the FIDs are computed differently in two FFG scenarios. All numbers are average of 50 runs with different reference glyphs.



Figure 6. **Generated Samples.** The generated images by five different models are shown. We also provide the reference and the source images used for the generation in the top two rows. The available ground truth images (GT) are shown in the bottom row. We highlight the samples that reveal the drawback of each model with colored boxes; green for AGIS-Net, red for FUNIT, and yellow for LF-Font.

FFG task. Using the style and content classifiers, Frechét inception distance (FID) [20] between the generated images and real images are computed and their harmonic mean is reported (FID(H)). We describe the details in the supplementary materials.

### 4.3. Experimental results

**Quantitative evaluation.** Table 1 shows the FFG performances by MX-Font and competitors. The reported values are the average of 50 different experiments, where four reference images per style are used for font generation in each experiment. In the table, we observe that MX-Font outperforms other methods in the both in-domain transfer scenario and zero-shot cross-lingual generation scenario with the most of evaluation metrics. Especially, MX-Font remarkably outperforms other methods in the cross-lingual task. In the in-domain transfer scenario, ours exceeds others in the classification accuracies and the user study. We observe that

MX-Font perform worse than others in the Chinese FID, where FID is known to sensitive to noisy or blur images, regardless of the image quality itself [39]. Our method shows the remarkably better performances in more reliable evaluation, user study in all criterions.

**Qualitative evaluation.** We illustrated the generated samples in Figure 6. We show four reference images to extract each style in the top row, and the source images in the second row where each source image is used to extract the content. In the green box in Figure 6, we observe that AGIS-Net often fails to reflect the reference style precisely and generate local details. FUNIT generally shows similar trends with AGIS-Net, while FUNIT often produces shattered glyphs when the target glyph and the source glyph have significantly different structures (red box). At a glance, LF-Font seems to capture the detailed local styles well. However, it often misses important detailed local component such as dot and stroke, as shown in the yellow box. Comparing to other
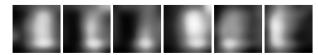
Figure 7. **Each localized expert attends different local areas.** We show the variance of Class Activation Maps (CAMs) on training images for each expert. The brighter intensity indicates that the variance of CAMs is higher in that region.



Figure 8. **Generated samples of the models having different number of heads.** The samples generated with four reference glyphs by the single-headed model and multi-headed model are shown. We highlight the defects in red dotted circles that appeared in the images generated by the single-expert model. $k$ denotes the number of experts.

|  | Acc (S) ↑ | Acc (C) ↑ | Acc (B) ↑ | LPIPS ↓ |
|---|---|---|---|---|
| Ours ($k = 1$) | 72.2 | 98.7 | 71.4 | 0.133 |
| Ours ($k = 6$) | 78.9 | 99.5 | 78.7 | 0.120 |

Table 2. **Impact of the number of experts $k$.** Single-expert model ($k = 1$) and multiple-experts model ($k = 6$, proposed) are compared on in-domain Chinese transfer benchmark.

methods, MX-Font synthesizes the better detailed structures both in content and style, owing to the strong representation power of locally specialized experts. The advantage of MX-Font is highlighted in the cross-lingual FFG. All existing models often generate unrecognizable characters under the cross-lingual scenario. Nevertheless, MX-Font preserves both the detailed local style and content and generates the plausible and recognizable images consistently. Such a noticeable gap in visual quality explains the large performance leap of MX-Font in the user study.

### 4.4. Analyses

**Learned local concepts by different experts.** We show the local concepts learned by each expert by visualizing where each expert attends on. We extract the Class Activation Maps (CAMs) of the training samples using the component classifier $Cls_u$ on each local feature. Then, we visualize the variance of CAMs in Figure 7. In Figure 7, the region of each image with bright intensity than the surrounding indicates the region where each expert pays more attention. Interestingly, without any explicit pixel-level annotation, our localized experts attend different local areas of the images. These maps support that each expert of MX-Font tends to

|  | Acc (S) ↑ | Acc (C) ↑ | Acc (B) ↑ | LPIPS ↓ |
|---|---|---|---|---|
| Ours ($Cls_u$) | 78.9 | 99.5 | 78.7 | 0.120 |
| Ours ($Cls_c$) | 94.8 | 0.04 | 0.04 | 0.214 |

Table 3. **Comparing the component classifier and the character classifier as weak supervision.** We compare two auxiliary classifiers as content supervision. Ours ($Cls_u$) denotes MX-Font using the component classifier and Ours ($Cls_c$) denotes the model replaced the component classifier to the character classifier.

| $\mathcal{L}_{indp,i}$ | $\mathcal{H}_{c,s}$ | $\mathcal{L}_{c,s}$ | Acc (S) | Acc (C) | Acc (B) |
|---|---|---|---|---|---|
| ✔ | ✔ | ✔ | **59.0** | **95.9** | **56.8** |
| ✗ | ✔ | ✔ | 52.0 | 95.8 | 50.0 |
| ✗ | ✗ | ✔ | 51.6 | 95.5 | 49.4 |
| ✗ | ✗ | ✗ | 27.8 | 89.1 | 24.7 |
| LF-Font [37] |  |  | 38.5 | 95.2 | 36.5 |

Table 4. **Impact of loss functions.** We compare models by ablating the proposed object functions trained and tested on Korean-handwriting dataset. The results show that the content-style adversarial loss $\mathcal{L}_{c,s}$ and the maximizing entropy term $\mathcal{H}_{c,s}$ and independent loss $\mathcal{L}_{indp,i}$ are all important components.

cover different local areas of the input image. Summarizing, these experimental studies demonstrate that *multiple localized experts* capture different local areas of the input image as we intend, and employing *multiple localized experts* helps us to enhance the quality of generated images by preserving the local details during the style-content disentanglement.

**Multiple experts vs. single expert.** We compare the performances of the single expert model ($k = 1$) with our multiple expert model ($k = 6$) on benchmark in-domain transfer scenario. The results are shown in Table 2 and Fig 8. We observe that using multiple heads is better than a single head in the classification accuracies. We also observe that the generated images by the single-headed model fails to preserve the local structures delicately, *e.g.* important strokes are missing, while the multi-headed model captures local details well.

**Character labels vs. local component labels.** We assume that component supervision is beneficial to learn experts with different local concepts. We replace the component supervision (multiple image-level sub-concepts) to the character supervision (single image-level label). Table 3 shows that utilizing character supervision incurs a mode collapse. We speculate that two reasons caused the collapse, (1) the number of characters ($\approx 19k$) is too large to learn, while the number of components is reasonably small (371), and (2) our best allocation problem prevents the experts from collapsing into the same values, while the character supervised model has no restriction to learn different concepts.

**Loss ablations.** We investigate the effect of our loss function design by the models trained and tested on Korean

handwritten dataset. The evaluation results are reported in Table 4. The detailed training settings are in supplementary materials. $\mathcal{H}_{c,s}$ denotes the maximizing entropy terms in the content-style adversarial loss, and $\mathcal{L}_{c,s}$ denotes the content-style adversarial loss. Table 4 shows that all the proposed loss functions for the style-content disentanglement are effective to enhance the overall performances.

## 5. Conclusion

We propose a novel few-shot font generation method, namely MX-Font. Our goal is to achieve both the rich representation for the local details and the generalizability to the unseen component and language. To this end, MX-Font employ *multi-headed encoder*, trained by weak local component supervision, *i.e.* style and content feature classifiers. Based on interactions between these *feature classifiers* and localized experts, MX-Font learns to disentangle the style and content successfully by developing localized features. Finally, the proposed model generates the plausible font images, which preserve both local detailed style of the reference images and precise characters of the source images. Experimental results show that MX-Font outperforms existing methods in in-domain transfer scenario and zero-shot cross-lingual transfer scenario; especially large performance leap in the cross-lingual scenario.

## Acknowledgements

## References

[1] Faez Ahmed, John P Dickerson, and Mark Fuge. Diverse weighted bipartite b-matching. *arXiv preprint arXiv:1702.07134*, 2017. 11

[2] Samaneh Azadi, Matthew Fisher, Vladimir G Kim, Zhaowen Wang, Eli Shechtman, and Trevor Darrell. Multi-content gan for few-shot font style transfer. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018. 3

[3] Hyojin Bahng, Sanghyuk Chun, Sangdoo Yun, Jaegul Choo, and Seong Joon Oh. Learning de-biased representations with biased representations. In *ICML*, 2020. 5, 13

[4] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018. 5

[5] Yue Cao, Jiarui Xu, Stephen Lin, Fangyun Wei, and Han Hu. GCNet: Non-local networks meet squeeze-excitation networks and beyond. In *IEEE International Conference on Computer Vision Workshops*, 2019. 11

[6] Junbum Cha, Sanghyuk Chun, Gayoung Lee, Bado Lee, Seonghyeon Kim, and Hwalsuk Lee. Few-shot compositional font generation with dual memory. In *Eur. Conf. Comput. Vis.*, 2020. 1, 2, 3, 4, 5, 6, 13

[7] Junbum Cha, Sanghyuk Chun, Gayoung Lee, Bado Lee, Seonghyeon Kim, and Hwalsuk Lee. Toward high-quality few-shot font generation with dual memory. *IEEE Conf. Comput. Vis. Pattern Recog. Worksh.*, 2020. 1

[8] Cheng Chen, Lan Zheng, Venkatesh Srinivasan, Alex Thomo, Kui Wu, and Anthony Sukow. Conflict-aware weighted bipartite b-matching and its application to e-commerce. *IEEE Transactions on Knowledge and Data Engineering*, 28(6):1475–1488, 2016. 11

[9] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. StarGAN: Unified generative adversarial networks for multi-domain image-to-image translation. In *CVPR*, 2018. 2

[10] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. StarGAN v2: Diverse image synthesis for multiple domains. In *CVPR*, 2020. 2

[11] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–2030, 2016. 5

[12] Yue Gao, Yuan Guo, Zhouhui Lian, Yingmin Tang, and Jianguo Xiao. Artistic glyph image synthesis via one-stage few-shot learning. *ACM Trans. Graph.*, 2019. 1, 2, 3, 6, 14

[13] Yiming Gao and Jiangqin Wu. Gan-based unpaired chinese character image translation via skeleton transformation and stroke rendering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020. 2

[14] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *CVPR*, 2016. 2

[15] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, 2014. 5

[16] Arthur Gretton, Olivier Bousquet, Alex Smola, and Bernhard Schölkopf. Measuring statistical dependence with hilbert-schmidt norms. In *International conference on algorithmic learning theory*, pages 63–77. Springer, 2005. 5, 13

[17] Arthur Gretton, Kenji Fukumizu, Choon H Teo, Le Song, Bernhard Schölkopf, and Alex J Smola. A kernel statistical test of independence. In *Advances in neural information processing systems*, pages 585–592, 2008. 5, 13

[18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2016. 13

[19] Byeongho Heo, Sanghyuk Chun, Seong Joon Oh, Dongyoon Han, Sangdoo Yun, Gyuwan Kim, Youngjung Uh, and Jung-Woo Ha. Adamp: Slowing down the slowdown for momentum optimizers on scale-invariant weights. In *Int. Conf. Learn. Represent.*, 2021. 13

[20] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*, 2017. 7, 14

[21] Xun Huang and Serge J Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, 2017. 2

[22] Yaoxiong Huang, Mengchao He, Lianwen Jin, and Yongpan Wang. Rd-gan: Few/zero-shot chinese character style transfer via radical decomposition and rendering. In *ECCV*, 2020. 2

[23] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. 2

[24] Yue Jiang, Zhouhui Lian, Yingmin Tang, and Jianguo Xiao. SCFont: Structure-guided chinese font generation via deep stacked networks. In *AAAI*, 2019. 2, 3

[25] Hanjoo Kim, Minkyu Kim, Dongjoo Seo, Jinwoong Kim, Heungseok Park, Soeun Park, Hyunwoo Jo, KyungHyun Kim, Youngil Yang, Youngkwan Kim, et al. NSML: Meet the MLaaS platform with a real-world case study. *arXiv preprint arXiv:1810.09957*, 2018. 9

[26] Diederick P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 5, 13

[27] Peter Kleinschmidt and Heinz Schannath. A strongly polynomial algorithm for the transportation problem. *Mathematical Programming*, 68(1):1–13, 1995. 11

[28] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *International Conference on Machine Learning*, 2019. 5

[29] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955. 2, 4, 11

[30] Chenhao Li, Yuta Taniguchi, Min Lu, and Shin'ichi Konomi. Few-shot font style transfer between different languages. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 433–442, 2021. 3

[31] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. Universal style transfer via feature transforms. In *Advances in Neural Information Processing Systems*, 2017. 2

[32] Yijun Li, Ming-Yu Liu, Xueting Li, Ming-Hsuan Yang, and Jan Kautz. A closed-form solution to photorealistic image stylization. In *ECCV*, 2018. 2

[33] Alexander H Liu, Yen-Cheng Liu, Yu-Ying Yeh, and Yu-Chiang Frank Wang. A unified feature disentangler for multi-domain image translation and manipulation. In *Advances in neural information processing systems*, 2018. 2

[34] 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.*, 2019. 2, 6, 13, 14

[35] Fujun Luan, Sylvain Paris, Eli Shechtman, and Kavita Bala. Deep photo style transfer. In *CVPR*, 2017. 2

[36] Oded Maron and Tomás Lozano-Pérez. A framework for multiple-instance learning. *Advances in neural information processing systems*, pages 570–576, 1998. 4

[37] Song Park, Sanghyuk Chun, Junbum Cha, Bado Lee, and Hyunjung Shim. Few-shot font generation with localized style representations and factorization. In *AAAI*, 2021. 1, 2, 3, 4, 5, 6, 8, 13, 14

[38] Novi Quadrianto, Viktoriia Sharmanska, and Oliver Thomas. Discovering fair representations in the data domain. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8227–8236, 2019. 5

[39] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. *arXiv preprint arXiv:2102.12092*, 2021. 7

[40] Le Song, Alex Smola, Arthur Gretton, Justin Bedo, and Karsten Borgwardt. Feature selection via dependence maximization. *Journal of Machine Learning Research*, 13(May):1393–1434, 2012. 13

[41] Nikita Srivatsan, Jonathan Barron, Dan Klein, and Taylor Berg-Kirkpatrick. A deep factorization of style and structure in fonts. In *Conference on Empirical Methods in Natural Language Processing*, 2019. 1, 3

[42] Danyang Sun, Tongzheng Ren, Chongxuan Li, Hang Su, and Jun Zhu. Learning to write stylized chinese characters by reading a handful of examples. In *IJCAI*, 2018. 1, 2, 3

[43] Yuchen Tian. zi2zi: Master chinese calligraphy with conditional adversarial networks, 2017. 2

[44] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *ECCV*, 2018. 11

[45] Shan-Jean Wu, Chih-Yuan Yang, and Jane Yung-jen Hsu. Calligan: Style and structure-aware chinese calligraphy character generator. *AI for Content Creation Workshop. CVPR Workshop*, 2020. 2

[46] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10687–10698, 2020. 4

[47] Jaejun Yoo, Youngjung Uh, Sanghyuk Chun, Byeongkyu Kang, and Jung-Woo Ha. Photorealistic style transfer via wavelet transforms. In *ICCV*, 2019. 2

[48] Xiaoming Yu, Yuanqi Chen, Shan Liu, Thomas Li, and Ge Li. Multi-mapping image-to-image translation via learning disentanglement. In *Advances in Neural Information Processing Systems*, 2019. 2

[49] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Int. Conf. Comput. Vis.*, 2019. 14

[50] Sangdoo Yun, Seong Joon Oh, Byeongho Heo, Dongyoon Han, Junsuk Choe, and Sanghyuk Chun. Re-labeling imagenet: from single to multi-labels, from global to localized labels. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 4

[51] Changqing Zhang, Yeqinq Liu, Yue Liu, Qinghua Hu, Xinwang Liu, and Pengfei Zhu. Fish-mml: Fisher-hsic multi-view metric learning. In *International Joint Conference on Artificial Intelligence*, pages 3054–3060, 2018. 5

[52] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *ICML*, 2019. 5, 13

[53] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 6

[54] Yexun Zhang, Ya Zhang, and Wenbin Cai. Separating style and content for generalized style transfer. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018. 1, 2, 3, 6, 14

[55] Zhi-Hua Zhou, Min-Ling Zhang, Sheng-Jun Huang, and Yu-Feng Li. Multi-instance multi-label learning. *Artificial Intelligence*, 176(1):2291–2320, 2012. 4

[56] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. 2

# Appendix

We describe additional experimental results to complement the main paper (§A). The implementation details are in §B. Finally, we provide the detailed evaluation protocols (§C).

## A. Additional experimental results

### A.1. More visual examples

We show more generated glyphs in Figure A.2. MX-Font correctly synthesizes the strokes, dot, thickness and size of the ground truth glyphs. In the cross-lingual FFG, MX-Font can produce promising results in that they are all readable. Meanwhile, all other competitors provide inconsistent results, which are often impossible to understand. These results show a similar conclusion as our main paper.

### A.2. Impact of the number of experts

In Table A.1, we report the performances by varying the number of experts, $k$. We observe that larger $k$ brings better performances until $k = 6$, but larger $k$, *e.g.*, 8, shows slightly worse performance than $k = 6$. We presume that this is because there are no sufficient data having more than or equal to eight components for training all the eight experts to capture different concepts. Figure A.1 illustrates the frequency of the number of components. From this graph, we find that the most characters have less than 8 components in our Chinese dataset. Moreover, larger $k$ means the number of parameters are increased, resulting in more training and inference runtime. Hence, in the paper, we choose $k = 6$ for all experiments.

## B. Implementation details

### B.1. Network architecture

Each localized expert $E_i$ has 11 layers including convolution, residual, global-context [5], and convolutional block attention (CBAM) [44] blocks. The multiple localized experts share the weights of their first five blocks. The two feature classifiers $Cls_s$ and $Cls_u$ have the same structure;



Figure A.1. **The distribution of number of components.** The left shows the percentage of characters with different number of components and the right shows the cumulative summation of the left.

| $k$ | Acc (S) $\uparrow$ | Acc (C) $\uparrow$ | Acc (B) $\uparrow$ | LPIPS $\downarrow$ |
|---|---|---|---|---|
| 1 | 72.2 | 98.7 | 71.4 | 0.133 |
| 2 | 79.0 | 99.3 | 78.5 | 0.128 |
| 6 | 78.9 | 99.5 | 78.7 | 0.120 |
| 8 | 75.5 | 99.5 | 75.2 | 0.123 |

Table A.1. **Impact of the number of experts** $k$. The models with different number of heads are compared on in-domain Chinese transfer benchmark.

a linear block following two residual blocks. The weights of the first two residual blocks are shared. The generator $G$ consists of convolution and residual blocks. Please refer our code for the detailed architecture.

### B.2. Component allocation problem to weighted bipartite B-matching problem

Given a bipartite graph $G = (V, E)$, where $V$ is a set of vertices, $E$ is a set of edges and $W$ is the weight values for each edge $e \in E$, the weighted bipartite B-matching (WBM) problem [27] aims to find subgraph $H = (V, E')$ maximizing $\sum_{e \in E} W(e)$ with every vertex $v \in V$ adjacent to at most the given budget, $B(v)$, edges. WBM problem can be solved by the Hungarian algorithm [29], a typical algorithm to solve combinatorial optimization in a polynomial time, in $O(|V||E|) = O(|V|^3)$. For curious readers, we refer recent papers solving variants of WBM problems [8, 1].

We recall the component allocation problem described in the main paper:

$$\max_{w_{ij} \in \{0,1\} | i=1\ldots k, j \in U_c} \sum_{i=1}^{k} \sum_{j \in U_c} w_{ij} p_{ij},$$

$$\text{s.t.} \quad \sum_{i=1}^{k} w_{ij} \geq 1 \text{ for } \forall j, \quad \sum_{j \in U_c} w_{ij} \geq 1 \text{ for } \forall i,$$

$$\sum_{j \in U_c} w_{ij} \leq \max\left(1, \left\lceil \frac{m}{k} \right\rceil\right) \text{ for } \forall i$$

$$\sum_{i=1}^{k} w_{ij} \leq \max\left(1, \left\lceil \frac{k}{m} \right\rceil\right) \text{ for } \forall j. \tag{B.1}$$

Figure A.2. **Generation samples.** We provide more generated glyphs with four reference glyphs.

We replace the last condition, $\sum_{i=1}^{k} \sum_{j \in U_c} w_{ij} = \max(k, m)$ to the upper bound condition where $\lceil \cdot \rceil$ denotes the ceiling function. For example, if $k = 3$ and $m = 4$, the budget for each expert is 2, while the budget for each component is 1. We build a bipartite graph where the vertex set contains all experts and all valid components, and the edge weights are the prediction probability $p_{ij}$. Now (B.1) can be re-formulated by the WBM problem.

## B.3. HSIC Formulation

When training MX-Font, we let the two feature outputs from different experts, or content and style features independent of each other. To measure the independence between content feature and style feature, we first assume that the content features $f_c$ and the style features $f_s$ are drawn from two different random variables, $Z_c$ and $Z_s$, *i.e.*, $f_c \sim Z_c$ and $f_s \sim Z_s$. We employ Hilbert Schmidt independence criterion (HSIC) [16] to measure the independence between two random variables. For two random variables $Z_c$ and $Z_s$, HSIC is defined as $\mathrm{HSIC}^{k,l}(Z_c, Z_s) := \|C_{Z_c Z_s}^{k,l}\|_{\mathrm{HS}}^2$ where $k$ and $l$ are kernels, $C^{k,l}$ is the cross-covariance operator in the Reproducing Kernel Hilbert Spaces (RKHS) of $k$ and $l$, $\| \cdot \|_{\mathrm{HS}}$ is the Hilbert-Schmidt norm [16, 17]. If we use radial basis function (RBF) kernels for $k$ and $l$, HSIC is zero if and only if two random variables are independent.

Since we only have the finite number of samples drawn from the distributions, we need a finite sample estimator of HSIC. Following Bahng *et al.* [3], we employ an unbiased estimator of HSIC, $\mathrm{HSIC}_1^{k,l}(Z_c, Z_s)$ [40] with $m$ samples. Formally, $\mathrm{HSIC}_1^{k,l}(Z_c, Z_s)$ is defined as:

$$
\mathrm{HSIC}_1^{k,l}(Z_c, Z_s) = \frac{1}{m(m-3)} \Bigg[ \mathrm{tr}(\widetilde{Z}_c \widetilde{Z}_s^T) +
$$
$$
\frac{\mathbf{1}^T \widetilde{Z}_c \mathbf{1} \mathbf{1}^T \widetilde{Z}_s^T \mathbf{1}}{(m-1)(m-2)} - \frac{2}{m-2} \mathbf{1}^T \widetilde{Z}_c \widetilde{Z}_s^T \mathbf{1} \Bigg]
$$
(B.2)

where $(i, j)$-th element of a kernel matrix $\widetilde{Z}_c$ is defined as, $\widetilde{Z}_c(i, j) = (1 - \delta_{ij}) k(f_c^i, f_c^j)$, and the $i$-th feature in the mini-batch $f_c^i$, is assumed to be sampled from the $Z_c$, *i.e.*, $\{f_c^i\} \sim Z_c$. We similarly define $\widetilde{Z}_s(i, j) = (1 - \delta_{ij}) l(f_s^i, f_s^j)$.

In practice, we compute $\mathrm{HSIC}_1^{k,l}(Z_c, Z_s)$ in a mini-batch, *i.e.*, $m$ is the batch size. We use the RBF kernel with kernel radius 0.5, *i.e.*, $k(f_c^i, f_c^j) = \exp(-\frac{1}{2}\|f_c^i - f_c^j\|_2^2)$.

## B.4. GAN objective details

We employ two conditional discriminators $D_s$ and $D_c$ which predict a style label $y_s$ and a content label $y_c$, respectively. In practice, we employ a multitask discriminator $D$, and different projection embeddings for content labels and style labels, following the previous methods [34, 6, 37]. The

hinge loss [52] is employed to high fidelity generation:

$$
\mathcal{L}_{adv}^D = \mathbb{E}_{(x, y_c, y_s)} \left[ [1 - D(x, y_s)]_+ + [1 - D(x, y_c)]_+ \right]
$$
$$
+ \mathbb{E}_{(\widetilde{x}, y_c, y_s)} \left[ [1 - D(\widetilde{x}, y_s)]_+ + [1 - D(\widetilde{x}, y_c)]_+ \right]
$$
$$
\mathcal{L}_{adv}^G = -\mathbb{E}_{(\widetilde{x}, y_c, y_s)} \left[ D(\widetilde{x}, y_s) + D(\widetilde{x}, y_c) \right],
$$
(B.3)

where $\widetilde{x}$ is the generated image by combining a content feature extracted from an image with content label $y_c$ and a style feature extracted from an image with style label $y_s$.

The feature matching loss $\mathcal{L}_{fm}$ and the reconstruction loss $\mathcal{L}_{recon}$ are formulated as follows:

$$
\mathcal{L}_{fm} = \mathbb{E}_{(x, \widetilde{x})} \left[ \sum_{l=1}^{L-1} \|D^l(x) - D^l(\widetilde{x})\|_1 \right],
$$
(B.4)
$$
\mathcal{L}_{recon} = \mathbb{E}_{(x, \widetilde{x})} \left[ \|x - \widetilde{x}\|_1 \right],
$$

where $L$ is the number of layers in the discriminator $D$ and $D^l$ denotes the output of $l$-th layer of $D$.

## B.5. Training details

We use Adam [26] optimizer to optimize the MX-Font. The learning rate is set to 0.001 for the discriminator and 0.0002 for the remaining modules. The mini-batch is constructed with the target glyph, style glyphs, and content glyph during training. Specifically, we first pick the target glyph randomly. Then, we randomly select $n$ style glyphs with the same style as the target glyph, and $n$ content glyphs with the same character as the target glyph for each target glyph. Here, the target glyph is excluded from the style and content glyphs selection. We set $n$ to 3 during training. We set the number of heads $k$ to 6 and train the model for 650k iteration with the full objective functions for the Chinese glyph generation. For the Korean, we set the number of heads $k$ to 3 and train the model for 200k iteration with the all objective functions except $\mathcal{L}_{\mathrm{indp\,exp},i}$. We do not employ the $\mathcal{L}_{\mathrm{indp\,exp},i}$ during training for the Korean glyph generation, due to the special characteristic of the Korean script; always decomposed to fixed number of components, *e.g.*, 3.

## C. Evaluation details

### C.1. Classifiers

Three classifiers are trained for the training; the style classifier, the Chinese character classifier, and the Korean character classifier. The style classifier and the Chinese character classifier are trained with the same Chinese dataset, including 209 Chinese fonts and 6428 Chinese characters per font. Besides, we used the Korean dataset that DM-Font [6] provides to train the Korean character classifier. The classifiers have ResNet-50 [18] structure. We optimize the classifiers using AdamP optimizer [19] with learning rate 0.0002 for 20 epochs. During training, the CutMix

(a) User study example (Chinese generation)



(b) User study example (Korean generation)

Figure C.1. **User study examples.** The example images that we provide to the candidates are shown. Each image includes the reference images, source images, and the generated images.

| FIDs | CN → CN | | | CN → KR | | |
|---|---|---|---|---|---|---|
| | S | C | H | S | C | H |
| EMD | 145.5 | 51.1 | 79.7 | 220.3 | 113.8 | 150.0 |
| AGIS-Net | 91.0 | 10.8 | 19.2 | 235.5 | 106.5 | 146.5 |
| FUNIT | 50.6 | 11.8 | 19.2 | 486.4 | 107.4 | 176.0 |
| LF-Font | **43.5** | **9.0** | **14.8** | 187.8 | 123.4 | 148.7 |
| MX-Font | 50.5 | 13.9 | 21.8 | **113.2** | **78.1** | **84.1** |

Table C.1. We provide style-aware (S), content-aware(C) FIDs measured by the style and content classifiers. The harmonic mean (H) of the style-aware and the content-aware FIDs values are identical to the values reported in the main table.

augmentation [49] is adopted and the mini-batch size is set to 64.

## C.2. LF-Font modification

Since LF-Font [37] cannot handle the unseen components in the test time due to its component-conditioned structure, we modify its structure to enable the cross-lingual font generation. We loose the component-condition of LF-Font in the test time only, by skipping the component-condition block when the unseen component is given. Note that, we use original LF-Font structure for the training to reproduce its original performance.

## C.3. User study examples

We show the sample images used for the user study in Figure C.1. Five methods, including EMD [54], AGIS-Net [12], FUNIT [34], LF-Font [37], and MX-Font are randomly displayed to users for every query.

## C.4. FID

We measure the style-aware and content-aware Frechét inception distance (FID) [20] between generated images and rendered images using the style and content classifier. For the Chinese glyphs, the style-aware and content-aware FIDs are measured with the generated glyphs and the corresponding ground truth glyphs. Since the ground truth glyphs of cross-lingual generation do not exist, the style-aware FID is measured the generated glyphs and all the available rendered glyphs having the same style with the generated images. The content-ware FID is measured similar to the style-aware FID. The style-aware (S) and the content-aware (C) FID values and their harmonic mean (H) are reported in Table C.1. Despite that MX-Font shows the slight degradation in FID for Chinese font generation, these results are not consistent with the user study and qualitative evaluation. For quantifying the image quality, we tend to trust the user study more because it better reveals the user's preference.