



# Diff-Font: Diffusion Model for Robust One-Shot Font Generation

Haibin He<sup>1</sup> · Xinyuan Chen<sup>2,3</sup> · Chaoyue Wang<sup>1</sup> · Juhua Liu<sup>1</sup> · Bo Du<sup>1</sup> · Dacheng Tao<sup>4</sup> · Qiao Yu<sup>3,5</sup>

Received: 7 May 2023 / Accepted: 31 May 2024

© The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2024

## Abstract

Font generation presents a significant challenge due to the intricate details needed, especially for languages with complex ideograms and numerous characters, such as Chinese and Korean. Although various few-shot (or even one-shot) font generation methods have been introduced, most of them rely on GAN-based image-to-image translation frameworks that still face (i) unstable training issues, (ii) limited fidelity in replicating font styles, and (iii) imprecise generation of complex characters. To tackle these problems, we propose a unified one-shot font generation framework called Diff-Font, based on the diffusion model. In particular, we approach font generation as a conditional generation task, where the content of characters is managed through predefined embedding tokens and the desired font style is extracted from a one-shot reference image. For glyph-rich characters such as Chinese and Korean, we incorporate additional inputs for strokes or components as fine-grained conditions. Owing to the proposed diffusion training process, these three types of information can be effectively modeled, resulting in stable training. Simultaneously, the integrity of character structures can be learned and preserved. To the best of our knowledge, Diff-Font is the first work to utilize a diffusion model for font generation tasks. Comprehensive experiments demonstrate that Diff-Font outperforms prior font generation methods in both high-fidelity font style replication and the generation of intricate characters. Our method achieves state-of-the-art results in both qualitative and quantitative aspects.

**Keywords** Font generation · One-shot image generation · Diffusion model-based framework · Conditional generation

---

Communicated by Seon Joo Kim.

---

Haibin He and Xinyuan Chen have contributed equally to this work.

---

✉ Chaoyue Wang  
chaoyue.wang@outlook.com

✉ Juhua Liu  
liujuhua@whu.edu.cn

Haibin He  
haibinhe@whu.edu.cn

Xinyuan Chen  
xychen9191@gmail.com

Bo Du  
dubo@whu.edu.cn

Dacheng Tao  
dacheng.tao@ntu.edu.sg

Qiao Yu  
yu.qiao@siat.ac.cn

<sup>1</sup> School of Computer Science, National Engineering Research Center for Multimedia Software, Institute of Artificial Intelligence, and Hubei Key Laboratory of Multimedia and Network Communication Engineering, Wuhan University, Wuhan, China

## 1 Introduction

Words are omnipresent in our everyday lives, appearing on book covers, signboards, advertisements, mobile phones, and even clothing. As a result, font generation holds significant commercial value and potential for application. However, designing a font library could be an extremely challenging task, particularly for glyph-rich languages with complex structures, such as Chinese (with over 60,000 glyphs) and Korean (with over 11,000 glyphs). Recently, the progress made in deep generative models, known for their capability to produce high-quality images, has indicated the feasibility of automatically generating diverse font libraries.

<sup>2</sup> School of Computer Science, Wuhan University, Wuhan, China

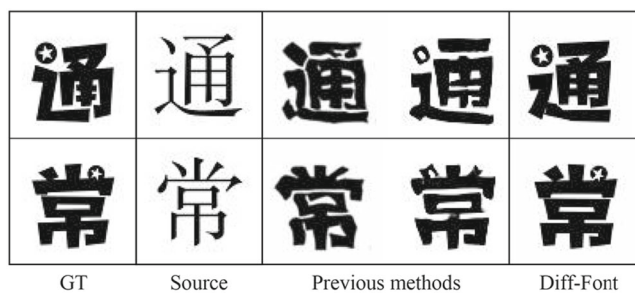
<sup>3</sup> Shanghai AI Laboratory, Shanghai 202150, China

<sup>4</sup> The College of Computing & Data Science at Nanyang Technological University, #32 Block N4 #02a-014, 50 Nanyang Avenue, Singapore 639798, Singapore

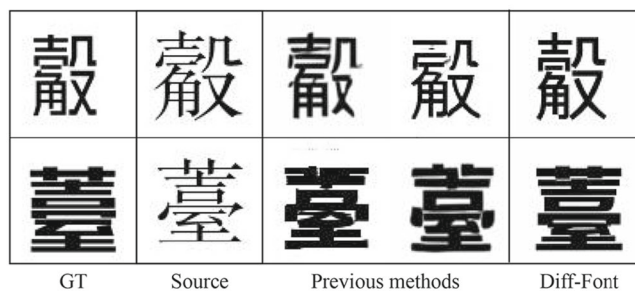
<sup>5</sup> Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, Shenzhen 518055, China

“Zi2zi” Tian (2017) is the first to adopt Generative Adversarial Networks (GANs) Goodfellow et al. (2020) to automatically generate a Chinese font library by learning a mapping from one style font to another, however, it needs paired data which is labor-reliant and expensive to collect. To facilitate the automatic synthesis of new fonts in an easy manner, numerous Few-shot (or even one-shot) Font Generation (FFG) methods have been proposed. These methods use a character image as the content and a few (or one) target characters to supply the font style, then their models are trained to generate the content character’s image with the target font style. Most existing FFG methods are built upon the GAN-based image-to-image translation framework. Some works follow unsupervised methods to obtain content and style features separately, and then fuse them in a generator to generate new characters Zhang et al. (2018b), Gao et al. (2019), Xie et al. (2021). Meanwhile, some other works exploit auxiliary annotations (e.g., strokes, components) to make the models aware of the specific structure and details about glyphs Jiang et al. (2019), Cha et al. (2020), Park et al. (2021a, b, 2022), Kong et al. (2022), Tang et al. (2022).

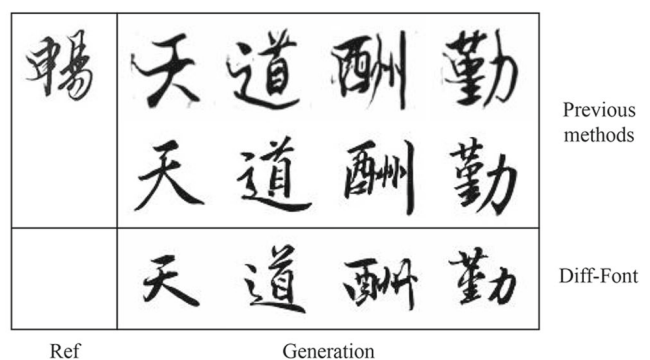
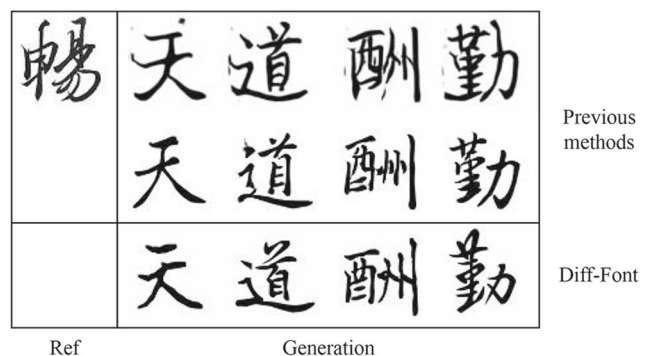
Although previous GAN-based methods have made significant progress and achieved impressive visual quality, font generation remains an extremely challenging long-tail task due to its stringent requirements for intricate details. Most of these methods still grapple with one or more of the following three challenges. Firstly, GAN-based methods employing adversarial training schemes may suffer from training instability and convergence difficulties, particularly with large datasets. While some tricks can alleviate this issue to some extent, they do not completely solve the problem. Secondly, GAN-based methods generally treat font generation as a style transfer problem between source and target image domains, often failing to separately model content and font style of characters. Consequently, neither significant font style transfers (i.e., drastic style changes) yield satisfactory results, nor subtle variations between two similar fonts are properly modeled. Last but not the least, when source characters are complex, these methods may struggle to ensure the integrity of the generated character structure. A qualitative illustration of problems arising from gaps in font style and complicated characters can be found in Fig. 1.



(a) Example for significant font style changes.



(c) Example for incorrect generation of complicated characters.



(b) Example for subtle font style variations.

**Fig. 1** Illustration for the problems caused by the gap in font style and complicated characters. **a** Example of significant font style changes: When the styles between the source and target glyphs differ significantly, methods based on an image-to-image translation framework may generate images with losing local details (column 3 and 4); **b**

Example for subtle font style variations: Our proposed Diff-Font can well capture the subtle variations between two fonts with similar styles while previous methods cannot; **c** Example for incorrect generation of complicated character: Image-to-Image translation framework may not perform well in generating characters with complicated structure

To tackle the aforementioned challenges, we introduce a novel diffusion model-based framework called Diff-Font for one-shot font generation. Instead of treating font generation as a style/domain transfer between a source font domain and a target font domain, the proposed Diff-Font approach considers font generation as a conditional generation task. Specifically, different character content is preprocessed into unique tokens, in contrast to the image inputs employed by previous methods which could cause confusion in similar glyphs. Regarding font styles, we utilize a pre-trained style encoder to extract style features as our conditional inputs. Moreover, to mitigate imprecise generation issues associated with glyph-rich characters, we incorporate a more fine-grained condition signal to help Diff-Font better model character structures. For Chinese fonts, we use stroke conditions, as strokes represent the smallest units that make up Chinese characters. Likewise, the components of Korean characters serve as the additional conditional input for Korean font generation. Instead of using the one-bit encoding employed in StrokeGAN Zeng et al. (2021), we employ count encoding to represent stroke (component) attributes, which more accurately reflects the character's stroke (component) properties. Consequently, the proposed Diff-Font effectively decouples the content and styles of characters, yielding high-quality generation results for complex characters. Simultaneously, thanks to the conditional generation pipeline and diffusion process, Diff-Font can be trained on large-scale datasets while exhibiting improved training stability compared to previous GAN-based methods. Lastly, we assemble a stroke-aware dataset for Chinese font generation and a component-aware dataset for Korean font generation.

In summary, the main contributions of this paper are as follows:

- We present Diff-Font, a unified generative network for robust one-shot font generation based on the diffusion model. In comparison to GAN-based methods, Diff-Font offers the advantages of stable training and the ability to be effectively trained on large datasets. To the best of our knowledge, this is the first attempt to develop a diffusion model for font generation.
- The proposed Diff-Font tackles the font generation task by employing a multi-attribute conditional diffusion model instead of the image-to-image translation framework. Character content and styles are processed as conditions, and the diffusion model utilizes these conditions to generate corresponding character images. Furthermore, a more fine-grained condition, such as stroke or component condition, is employed to enhance the generation of scripts with complex structures. Extensive experiments demonstrate the efficacy of our Diff-Font for one-shot font generation in comparison to previous state-of-the-art methods.

- We have compiled and annotated a stroke-wise dataset for Chinese and a component-wise dataset for Korean, which we believe can enhance font generation performance from the perspective of strokes and components. The source code, pre-trained models, and datasets are available at <https://github.com/Hxyz-123/Font-diff>.

The rest of this paper is organized as follows. In Sect. 2, we briefly review the related works. In Sect. 3, we introduce our proposed method in detail. Section 4 reports and discusses our experimental results. Lastly, we conclude our study in Sect. 5.

## 2 Related Work

### 2.1 Image-to-Image Translation

The task of image-to-image translation involves learning a mapping function that can transform source domain images into corresponding images that preserve the content of the original images while exhibiting the desired style characteristics of the target domain. Generating fonts can be achieved by means of the image-to-image translation models, which can be used to generate any desired font styles from a given content font image. Image-to-image translation using generative adversarial networks (GANs) has been a classical problem in the field of computer vision. Many works have been proposed to address this problem. Conditional GAN-based methods Mirza and Osindero (2014), such as Pix2Pix Isola et al. (2017), require paired data to guide the generation process. To eliminate the dependency on paired data, unsupervised methods have been proposed, including cycle-consistency-based approaches Zhu et al. (2017a), Yi et al. (2017), Kim et al. (2017), Kancharagunta and Dubey (2019) and the UNIT Liu et al. (2017) framework that leverages CoGAN Liu and Tuzel (2016) and VAE An and Cho (2015). BicycleGAN Zhu et al. (2017b) enables one-to-many domain translation by building a bijection between latent coding and output modes. For many-to-many domain translation, methods such as MUNIT Huang et al. (2018), CD-GAN Yang et al. (2018) and FUNIT Liu et al. (2019) disentangle the content and style representations using two encoders and couple them. Recently, due to the impressive results of the diffusion model, many diffusion model-based methods Saharia et al. (2022a), Sasaki et al. (2021), Zhao et al. (2022), Li et al. (2022), Wolleb et al. (2022) are proposed to tackle image-to-image tasks. However, controlling the generated output using diffusion model-based methods remains a challenge, and further exploration and development are needed, especially in the context of font generation.

Existing image-to-image translation methods generally focus on transforming object pose, texture, color, and style

while preserving the content structure, which may not be directly applicable to font generation. Unlike natural images, font styles are primarily defined by variations in shape and specific stroke rules rather than texture and style information. As a result, content structure information may also change during the font generation process. Therefore, applying image-to-image translation methods directly cannot produce satisfactory results.

## 2.2 Few-Shot Font Generation

Few-shot font generation aims to generate an entire font library with thousands of characters with only a few reference-style images as input. Existing few-shot font generation methods are predominantly based on the image-to-image translation framework, which transfers the source style of content characters to the reference style. To incorporate font-specific prior information into the method or the labels for careful design, various approaches have been proposed, demonstrating the potential of integrating such knowledge to improve the quality and diversity of generated fonts. DG-Font Xie et al. (2021) implements effective style transfer by replacing the traditional convolutional blocks with deformable convolutional blocks in an unsupervised framework TUNIT Baek et al. (2021). ZiGAN Wen et al. (2021) projects the same character features of different styles into Hilbert space to learn coarse-grained content knowledge. Some methods employ extra information to enhance training, e.g., strokes and components. SC-Font Jiang et al. (2019) uses stroke-level data to improve the correctness of structure and reduce stroke errors in generated images. DM-Font Cha et al. (2020) employs a dual-memory architecture to disassemble glyphs into stylized components and reassemble them into new glyphs. Its extension version LF-Font Park et al. (2021a, 2022) designs component-wise style encoder and factorization modules to capture local details in rich text design. MX-Font Park et al. (2021b) has a multi-headed encoder for specializing in different local sub-concepts, such as components, from the given image. FS-Font Tang et al. (2022) proposes a Style Aggregation Module (SAM) and an auxiliary branch to learn the component styles from references and the spatial correspondence between the content and reference glyphs. CG-GAN Kong et al. (2022) proposes a component discriminator to supervise the generator decoupling content and style at a fine-grained level. However, all methods mentioned above are based on GANs, which suffer from instability during training due to their adversarial objective and are prone to mode collapse, leading to suboptimal results especially for font styles with significant or subtle variations. As a result, there remains potential for improvement in the quality of font generation.

## 2.3 Diffusion Model

Diffusion Model is a new type of generative model that leverages the iterative reverse diffusion process to generate high-quality images and model complex distributions. It provides state-of-the-art performance in terms of image quality and can generate diverse outputs without mode collapse. Specifically, It employs a Markov chain to convert the Gaussian noise distribution to the real data distribution. Sohl-Dickstein et al. (2015) first clarify the concept of diffusion probabilistic model and denoising diffusion probabilistic models (DDPM) Ho et al. (2020) improves the theory and proposed to use a UNet to predict the noise added into the image at each diffusion time step. Dhariwal and Nichol (2021) propose a classifier-guidance mechanism that adopts a pre-trained classifier to provide gradients as guidance toward generating images of the target class. Ho and Salimans (2022) propose a technique that jointly trains a conditional and an unconditional diffusion model without using a classifier named classifier-free guidance. DDIM Song et al. (2020) extends the original DDPM to non-Markovian cases and is able to make accurate predictions with a large step size that reduces the sampling steps to one of the dozens. Glide Nichol et al. (2021), DALL-E2 Ramesh et al. (2022), Imagen Saharia et al. (2022b) and Stable Diffusion Rombach et al. (2022) introduce a pre-trained text encoder to generate semantic latent spaces and achieve exceptional results in a text-to-image task. Although the above methods have shown amazing results in image generation, they often focus on generating a specific category of objects or concept-driven generation guided by text prompts, with limited controllability.

Some other works explore the use of multiple conditions to guide the generation of diffusion models. SDG Liu et al. (2021) designs a sampling strategy, which adds multi-modal semantic information to the sampling process of the unconditional diffusion model for achieving language guidance and image guidance generation. ILVR Choi et al. (2021) uses a reference image at each time step during sampling to guide the generation. Diss Cheng et al. (2022) uses stroke images and sketch images as multi-conditions to train a conditional diffusion model to generate images from hand-drawings. Liu et al. (2022) consider the diffusion model as a combination of energy-based models and propose two compositional operators, conjunction and negation, to achieve zero-shot combinatorial generalization to a larger number of objects. Nair et al. (2022) guides the generation of diffusion model by calculating the comprehensive condition scores of multiple modes to solve the problem of multi-modal image generation. ControlNet Zhang and Agrawala (2023) introduces an extra conditional control module to enable a pre-trained diffusion model to be applied to specific tasks. This work is further extended by the multi-attribute conditional diffusion model

which introduces composite-wise and stroke-wise attributes conditional for better training and attribute-wise diffusion guidance strategy for stroke-aware or component-aware font generation.

### 3 Methodology

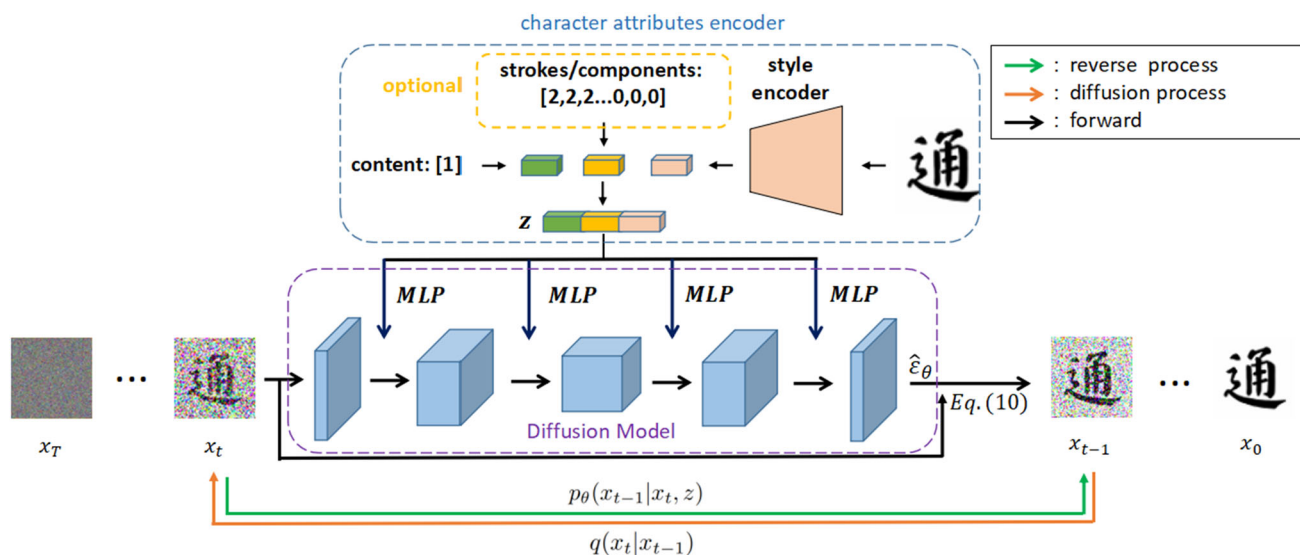
In this section, we introduce the details of Diff-Font. We first illustrate the framework of our model by incorporating the attributes of content, style, strokes and components (Sect. 3.1). Then, we elucidate the training process by formulating our multi-attributes conditional diffusion model (Sect. 3.2). Lastly, we present the adopted strategy to achieve attribute-wise guidance that can set the guidance level of attribute conditions separately during the generation process (Sect. 3.3).

#### 3.1 The Framework of Diff-Font

The framework of our proposed Diff-Font is illustrated in Fig. 2. As shown, Diff-Font consists of two modules: a character attributes encoder, which encodes the attributes of a character (i.e., content, style, strokes, components) into a latent variable, and a diffusion generation model, which uses the latent variable as a condition to generate the character image from Gaussian noise. The character attributes encoder is designed to process the attributes (content, style, strokes, components) of a character image separately.

In the character attributes encoder  $f$ , the content (denoted as  $c$ ), style (denoted as  $s$ ), and optional condition (like strokes or components, denoted as  $op$ ) are encoded as the latent variable:  $z = f(c, s)$ . If using the optional condition, then  $z = f(c, s, op)$ . Unlike previous font generation methods based on image-to-image translation that use the images from the source domain to obtain the content representations, we regard different content characters as different tokens. As practices commonly used in the NLP community (Devlin et al., 2018; Cui et al., 2021; Touvron et al., 2023), we adopt a randomly initialized embedding layer to convert different tokens of characters into different content representations. Specifically, different character content is first tokenized, and then the embedding layer is employed to transform these tokens into unique content embeddings. The content embedding layer is updated together with the diffusion generator. There are three reasons why we chose a content embedding layer instead of a content encoder. Firstly, characters are usually a finite set, making it possible to use countable tokens to represent the content of character and encode tokens as content embedding by a content embedding layer. Secondly, a content embedding layer consumes less computing resources than a content encoder. Lastly, using an embedding layer to encode the tokenized content can avoid the confusion of similar glyphs when using content encoder.

The style representation is extracted by a pre-trained style encoder. A trained style encoder in DG-Font is used as our pre-trained style encoder and its parameters are frozen in our diffusion model training. As for strokes (or components), we encode each character into a 32-dimensional vector. Each

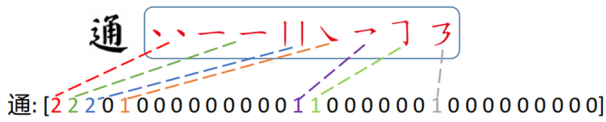


**Fig. 2** Overview of our proposed method. In the diffusion process, we gradually add noise to image  $x_0$ , and make it become approximately a Gaussian noise after time step  $T$ . For the reverse diffusion process, we use a latent variable  $z$ , which contains the content, style, and other

optional attributes semantic information of  $x_0$ , as a condition to train a diffusion model (based on UNet architecture) to predict the added noise at each time step in the diffusion process

No	Stroke	Name	example	No	Stroke	Name	example
1	丶	Dian(点)	立	17	乙	HengZheWanGou(横折弯钩)	九
2	一	Heng(横)	丛	18	㇇	HengPieWanGou(横撇弯钩)	那
3	丨	Shu(竖)	十	19	㇇	HengZheZheZheGou(横折折折钩)	乃
4	ノ	Pie(撇)	八	20	㇇	ShuZheZheGou(竖折折钩)	马
5	㇇	Na(捺)	人	21	㇇	ShuWan(竖弯)	四
6	㇇	Ti(提)	习	22	㇇	HengZheWan(横折弯)	没
7	㇇	PieDian(撇点)	女	23	㇇	HengZhe(横折)	口
8	丨	ShuTi(竖提)	长	24	㇇	ShuZhe(竖折)	山
9	㇇	HengZheTi(横折提)	认	25	㇇	PieZhe(撇折)	云
10	㇇	WanGou(弯钩)	狗	26	㇇	HengPie(横撇)	水
11	㇇	ShuGou(竖钩)	小	27	㇇	HengZheZhePie(横折折撇)	及
12	㇇	ShuWanGou(竖弯钩)	儿	28	㇇	ShuZhePie(竖折撇)	专
13	㇇	XieGou(斜钩)	我	29	㇇	HengXieGou(横斜钩)	飞
14	㇇	WoGou(卧钩)	心	30	㇇	ShuZheZhe(竖折折)	鼎
15	㇇	HengGou(横钩)	买	31	㇇	HengZheZhe(横折折)	凹
16	㇇	HengZheGou(横折钩)	用	32	㇇	HengZheZheZhe(横折折折)	凸

(a) 32 basic strokes of Chinese characters.



(b) Strokes and stroke count encoding vector of Chinese character ‘Tong’.

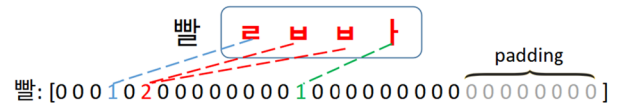
**Fig. 3** **a** 32 basic strokes of Chinese characters. The first and sixth columns are the dimensional locations of the basic strokes in the stroke vector. **b** Strokes and stroke count encoding vector of Chinese character ‘Tong’. Each dimension of the encoding vector represents the counts of corresponding basic stroke it contains

dimension of the vector represents the number of corresponding basic strokes (or components) it contains (shown in Figs. 3 and 4). This count encoding can better represent the stroke (or component) attribute of a character than one-bit encoding used in StrokeGAN Zeng et al. (2021). Thereafter, a stroke (or component) vector can be expanded into a vector consistent with the dimension of the content embedding. Using this method, we can obtain attribute representations of a character image and then concatenate them as a condition  $z$  for later conditional diffusion model training.

In the diffusion process, we add random gaussian noise to the real image  $x_0$  slowly to obtain a long Markov chain from the real image  $x_0$  to noise  $x_T$ . We adopt UNet architecture as our diffusion model and follow Dhariwal and Nichol (2021) to learn the reverse diffusion process. The reverse diffusion process generates characters images from gaussian noise by using multi-attributes condition latent variable  $z$ . This conditional generation is designed to mitigate the impact of the distinction in font style.

1	2	3	4	5	6	7	8
㇇	㇇	㇇	㇇	㇇	㇇	㇇	㇇
9	10	11	12	13	14	15	16
㇇	㇇	㇇	㇇	㇇	㇇	㇇	㇇
17	18	19	20	21	22	23	24
㇇	㇇	㇇	㇇	㇇	㇇	㇇	㇇

(a) 24 basic Korean components.



(b) Components and count encoding vector of example Korean character.

**Fig. 4** **a** 24 basic components of Korean characters. **b** Components and count encoding vector of example Korean character. We encode Korean components in the same way as Chinese strokes. Since Korean has only 24 basic components, we pad into 32 dimensions with 0

### 3.2 Multi-Attributes Conditional Diffusion Model

In our method, we regard each raw image of the character which is determined by its content ( $c$ ), style ( $s$ ) (and optional conditions ( $op$ )) attributes as a sample in the whole training data distribution, and denote the sample as  $x_0 \sim q(x_0 | f(c, s))$ . If using the optional condition, then,  $x_0 \sim q(x_0 | f(c, s, op))$ . Like the thermal motion of molecules, we add random Gaussian noise to the image thousands of times to gradually transform it from a stable state to a chaotic state. This process is called diffusion process and can be defined as:

$$q(x_{1:T} | x_0) = \prod_{t=1}^T q(x_t | x_{t-1}), \tag{1}$$

where

$$q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t \mathbf{I}), \tag{2}$$

$$t = 1, \dots, T,$$

and  $\beta_1 < \dots < \beta_T$  is a variance schedule following Ho et al. (2020). According to the Eq. 2,  $x_t$  can be rewritten as:

$$x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon_{t-1}, \quad \epsilon_{t-1} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \tag{3}$$

$$= \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, \quad \bar{\alpha}_t = \prod_{i=1}^t \alpha_i, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \tag{4}$$

$$\sim \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I}) \tag{5}$$

where  $\alpha_t = 1 - \beta_t$ , and  $\alpha_t$  is negatively correlated with  $\beta_t$ , therefore  $\alpha_1 > \dots > \alpha_T$ . When the  $T \rightarrow \infty$ ,  $\bar{\alpha}_T$  close to 0,  $x_T$  nearly obeys  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  and the posterior  $q(x_{t-1} | x_t)$  is also a Gaussian. So in the reverse process, we can sample a noisy image  $x_T$  from an isotropic Gaussian and generate the designated character image by denosing the  $x_T$  in the long Markov chain with a multi-attributes condition  $z = f(c, s)$  (if using the optional condition, then,  $z = f(c, s, op)$ ) that contains the semantic meaning of character. Since the posterior  $q(x_{t-1} | x_t)$  is hard to estimate, we use  $p_\theta$  to approximate the posterior distribution which can be denoted as:

$$p_\theta(x_{0:T} | z) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1} | x_t, z), \quad (6)$$

$$p_\theta(x_{t-1} | x_t, z) = \mathcal{N}(\mu_\theta(x_t, t, z), \Sigma_\theta(x_t, t, z)), \quad (7)$$

Following DDPM Ho et al. (2020), we set  $\Sigma_\theta(x_t, t, z)$  as constants and the diffusion model  $\epsilon_\theta(x_t, t, z)$  learns to predict the noise  $\epsilon$  added to  $x_0$  in diffusion process from  $x_t$  and condition  $z$  for easier training. Through these simplified operations, we can adopt a standard MSE loss to train our multi-attributes-conditional diffusion model:

$$L_{simple} = \mathbb{E}_{x_0 \sim q(x_0), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), z} [\| \epsilon - \epsilon_\theta(x_t, t, z) \|^2]. \quad (8)$$

### 3.3 Attribute-wise Diffusion Guidance Strategy

For glyph-rich scripts (e.g., Chinese and Korean), we adopt a two-stage training strategy to improve the generation effect. Based on the multi-attributes conditional training (i.e., first training stage), we also design a fine-tuning strategy (second training stage) that randomly discards content attribute or stroke (or component) attribute vectors with a 30% probability. If the content and stroke (or component) are discarded at the same time, the style attribute vector also be discarded. Such strategy has two advantages: first, it can enable our model to be more sensitive to these three attributes, and second, it can reduce the number of hyperparameters for we only need two guidance scales instead of three. In our case, we use zero vectors to replace the discarded attribute vectors, denoted as  $\mathbf{0}$ . When sampling, we modify the predicted noise to  $\hat{\epsilon}_\theta$ :

$$\begin{aligned} \hat{\epsilon}_\theta(x_t, t, f(c, s, op)) &= \epsilon_\theta(x_t, t, \mathbf{0}) \\ &+ s_1 * (\epsilon_\theta(x_t, t, f(c, s, \mathbf{0})) - \epsilon_\theta(x_t, t, \mathbf{0})) \\ &+ s_2 * (\epsilon_\theta(x_t, t, f(\mathbf{0}, s, op)) - \epsilon_\theta(x_t, t, \mathbf{0})), \end{aligned} \quad (9)$$

where  $s_1$  and  $s_2$  are the guidance scales of content and strokes. Then we adopt DDIM Song et al. (2020) to sample on a subset of diffusion steps  $\{\tau_1, \dots, \tau_S\}$  and set the variance weight parameter  $\eta = 0$  to speed up the generation process. So, we

can obtain  $x_{\tau_{i-1}}$  from  $x_{\tau_i}$  by the following equation:

$$x_{\tau_{i-1}} = \sqrt{\bar{\alpha}_{\tau_{i-1}}} \left( \frac{x_{\tau_i} - \sqrt{1 - \bar{\alpha}_{\tau_i}} \hat{\epsilon}_\theta}{\sqrt{\bar{\alpha}_{\tau_i}}} \right) + \sqrt{1 - \bar{\alpha}_{\tau_{i-1}}} \hat{\epsilon}_\theta. \quad (10)$$

The final character image  $x_0$  can be obtained by iterating through the above formula.

## 4 Experiments

In this section, we evaluate the performance of the proposed method on the one-shot font generation task by comparing it with state-of-the-art methods. In Sect. 4.1, we first introduce the datasets and evaluation metrics used to conduct experiments. The implementation details are described in Sect. 4.2. The results of qualitative and quantitative comparisons between Diff-Font and previous SOTA methods on different script generation are listed in Sects. 4.3, 4.4, 4.5 and 4.6. Limitations are discussed in Sect. 4.7.

### 4.1 Datasets and Evaluation Metrics

#### 4.1.1 Chinese Font Datasets

We collect 410 fonts (styles) including handwritten fonts and printed fonts as our whole dataset. Each font has 6625 Chinese characters that cover almost all commonly used Chinese characters. To evaluate the capacity of methods for different scale datasets, we use a small dataset and a large dataset for experiments. For the small dataset, the training set contains 400 fonts and 800 randomly selected characters, and the testing set contains the remaining 10 fonts with the same characters as the training set. For the large dataset, we use the same 400 fonts but all 6625 characters in training. The testing set consists of the remaining 10 fonts and 800 characters with complex structures and multiple strokes. In our experiment, the number of small dataset is set consistent with previous methods Xie et al. (2021). For fair comparison, the image size is also the same as the previous methods Xie et al. (2021), Zhang et al. (2018b), which is set as  $80 \times 80$ .

#### 4.1.2 Evaluation Metrics

In order to quantitatively compare our method with other advanced methods, we use the common evaluation metrics in image generation task, e.g., SSIM Wang et al. (2004), RMSE, LPIPS Zhang et al. (2018a), FID Heusel et al. (2017). SSIM (Structural Similarity) imitates the human visual system to compare the structural similarity between two images from three aspects: luminance, contrast and structure. RMSE (Root Mean Square Error) evaluates the similarity between

two images by calculating the root mean square error of their pixel values. Both of them are pixel-level metrics. LPIPS (Learned Perceptual Image Patch Similarity), a perceptual-level metric, measures the distance between two images in a deep feature space.

For computing the FID, we utilize an Inception-v3 model pre-trained on the ImageNet dataset Heusel et al. (2017) in accordance with Xie et al. (2021), Kong et al. (2022). Then calculating the Fréchet Distance between the final average pooling features of generated images and real images which are extracted by the Inception-v3 model. Following previous works Park et al. (2021a, b, 2022), we trained content classifier on test characters and style classifier on test fonts to classify character labels (content-aware) and font labels (style-aware). The architecture of classifiers are consistent with the setting in MX-font Park et al. (2021b).  $ACC_C$ ,  $ACC_S$  and  $ACC_B$  represent the classification accuracy of content labels, style labels, and combining content and style labels, respectively, as measured by these two trained classifiers. Moreover, we follow the similar idea in MX-font to conduct user study for human testing.

## 4.2 Implementation Details

### 4.2.1 Character Attributes Encoder

Character attributes encoder in Diff-Font consists of a content embedding layer, a style encoder, a style embedding layer, and an optional embedding layer. The architecture of our style encoder is the same as the style encoder in DG-Font, and the dimensions of the output feature maps are set to 128. Specifically, we adopt an embedding layer for the content attribute and optional attribute respectively, and an MLP for the style attribute. If using the optional attribute, the dimensions of the content, style and optional attribute vectors are set to 128, 128 and 256, respectively. Otherwise, the dimensions of both the content and style vectors are set to 256. Finally, they are concatenated as a 512 dimensions conditional latent vector  $z$  for training.

### 4.2.2 Multi-attributes Conditional Diffusion Model

Our multi-attributes conditional diffusion model is based on DDPM architecture. We list the hyperparameters setting for our training in Table 1. For sampling, we set 25 sampling steps to speed up the generation process.

### 4.3 Comparison with State-of-the-art Methods

Due to the complexity of the dataset we use, in addition to the natural image generation method FUNIT, we choose two advanced font generation methods, MX-Font and DG-Font, for Chinese one-shot font generation comparison: (1) **FUNIT**

**Table 1** Hyperparameters setting for multi-attributes conditional diffusion model

	Small dataset	Large dataset
Images trained	320K	2.65M
Batch size	24	64
Channels	128	128
Res. blocks num	3	3
Channel multiplier	1, 2, 3, 4	1, 2, 3, 4
Attention resolution	[40, 20, 10]	[40, 20, 10]
Diffusion steps	1000	1000
Noise Schedule	Linear	Linear
Conditional training iters	300k	420k
Fine-tuning iters	300k	380k
Learning rate	1e−4	1e−4
Optimizer	Adam with no weight decay	
Loss	MSE	MSE

Liu et al. (2019): FUNIT is a few-shot image-to-image translation framework that disentangles content and style representations by two different encoders and uses AdaIN Huang and Belongie (2017) to couple them. (2) **MX-Font** Park et al. (2021b): MX-Font extracts different local sub-concepts by employing multi-headed encoders. (3) **DG-Font** Xie et al. (2021): DG-Font uses the deformable convolution to replace the traditional convolution in an unsupervised framework. All these methods are based on GANs.

We use both datasets described in Sect. 4.1 to retrain models of FUNIT, MX-Font and DG-Font. During the generation process, only one reference character image with the target font is used. When evaluating these GANs-based methods, we choose the Song font commonly used in the font generation task as the source font Xie et al. (2021), Park et al. (2021b).

### 4.3.1 Quantitative Comparison

Table 2 shows the quantitative comparison results between our method and other previous state-of-the-art methods. In the experiments on both small and large datasets, Diff-Font achieves the best performance on all evaluation metrics of SSIM, RMSE, LPIPS and FID. In particular, our method has a great improvement over the second-best method in terms of FID indicators, 22.4% for the small dataset and 39.2% for the large dataset. The excellent performance on two scale datasets demonstrates the effectiveness and advantage of our Diff-Font. As for classification results, Diff-Font outperforms other methods in terms of  $ACC_C$ ,  $ACC_S$  and  $ACC_B$ , both on small and large datasets.



**Table 2** Quantitative comparison results on two different scale datasets

Methods	SSIM ( $\uparrow$ )	RMSE ( $\downarrow$ )	LPIPS ( $\downarrow$ )	FID ( $\downarrow$ )	ACC <sub>C</sub> ( $\uparrow$ )	ACC <sub>S</sub> ( $\uparrow$ )	ACC <sub>B</sub> ( $\uparrow$ )
<i>Quantitative comparison on small dataset</i>							
FUNIT	0.700	0.303	0.166	35.20	98.81	81.06	80.25
MX-Font	0.721	0.283	0.151	37.15	97.68	81.28	79.36
DG-Font	0.729	0.280	0.137	43.44	98.29	82.28	81.05
Diff-Font(ours)	<b>0.742</b>	<b>0.271</b>	<b>0.124</b>	<b>27.30</b>	<b>99.36</b>	<b>93.05</b>	<b>92.46</b>
<i>Quantitative comparison on large dataset</i>							
FUNIT	0.682	0.311	0.166	26.70	75.71	78.59	60.34
MX-Font	0.692	0.298	0.138	26.64	95.01	78.30	74.20
DG-Font	0.709	0.292	0.112	28.63	95.25	92.81	88.24
Diff-Font(ours)	<b>0.722</b>	<b>0.277</b>	<b>0.104</b>	<b>16.20</b>	<b>95.78</b>	<b>96.55</b>	<b>92.48</b>

ACC<sub>C</sub>, ACC<sub>S</sub> and ACC<sub>B</sub> respectively indicate the classification accuracy of content labels, style labels, and combining content and style labels. The best performance is marked in **bold**

### 4.3.2 Qualitative Comparison

The qualitative comparison results are shown in Fig. 5. For qualitative comparison, we define style and content based on the difficulty of implementation as follows. The target styles similar to the source font are regarded as easy styles, otherwise as difficult styles. The characters with the number of strokes less than or equal to 10 are defined as easy contents, and the characters with the number of strokes more than or equal to 15 as difficult contents. We make qualitative comparisons under the three settings of ESEC (easy styles and easy contents), ESDC (easy styles and difficult contents), and DSDC (difficult styles and difficult contents), respectively. As shown in Fig. 5, FUNIT often generates incomplete characters, and when the character structure is more complex, it would produce distorted structures. MX-Font could maintain the shape of characters to a certain extent, but it tends to generate vague characters and unclear backgrounds. DG-Font performs well in ESEC task, but losses some important stroke detailed local components in ESDC and DSDC tasks. Compared to these previous methods, our proposed Diff-Font could generate high quality character images in all three tasks.

In addition, Fig. 6 shows more qualitative comparison results on four chosen art fonts to better illustrate the effectiveness and advantages of Diff-Font. As these comparison results, when there is significant stylistic difference between the source and target font, GAN-based image-to-image translation frameworks would lead to worse structural distortion and loss of details, and our proposed Diff-Font based on conditional diffusion model could effectively reduce the occurrence.

### 4.3.3 Human Testing

We conducted a user study with 10 test fonts, as specified in Sect. 4.1.

Each method was applied to generate a line of ancient Chinese poetry on each font, and 64 participants were asked to evaluate the results based on content, style and both of them, respectively. Participants chosen their favorite output, so we obtained  $64 \times 10 \times 3 = 1920$  results and calculated the percentage of scores for each method. Some visualization of generation examples are shown in Fig. 7, and study results are presented in Table 3.

As can be seen, our proposed Diff-Font achieves the best score in human testing among the three evaluation criteria, which also verifies the effectiveness of our proposed framework.

### 4.4 Ablation Studies

In this part, we further conduct ablation studies to evaluate the effectiveness of the stroke count encoding, and discuss the impact of guidance scales.

#### 4.4.1 Effectiveness of the Stroke Count Encoding

We train three Diff-Font separately on the small dataset, one does not use the stroke condition, one uses the one-bit encoding stroke condition and the remaining one uses the count encoding stroke condition. As is shown in Table 4, using count encoding stroke condition achieves the best quantitative results in all evaluation metrics among the three models and we can observe that adding the one-bit encoding stroke condition (Fig. 8) even causes a decline in model performance. In the visualization result of columns 2 and 3 in Fig. 9, we find that other characters with the same basic strokes are generated when using the one-bit encoding. And according

Reference:	喻	喻	喻	喻	喻
Source:	蚘 蚘 蚤 虬	根 裕 裛 衿	讫 讫 讫 讫	祝 贯 责 豈	邳 邯 邳 邳
FUNIT:	蚘 蚘 蚤 虬	根 裕 裛 衿	讫 讫 讫 讫	祝 贯 责 豈	邳 邯 邳 邳
MX-Font:	蚘 蚘 蚤 虬	根 裕 裛 衿	讫 讫 讫 讫	祝 贯 责 豈	邳 邯 邳 邳
DG-Font:	蚘 蚘 蚤 虬	根 裕 裛 衿	讫 讫 讫 讫	祝 贯 责 豈	邳 邯 邳 邳
Diff-Font:	蚘 蚘 蚤 虬	根 裕 裛 衿	讫 讫 讫 讫	祝 贯 责 豈	邳 邯 邳 邳
GT:	蚘 蚘 蚤 虬	根 裕 裛 衿	讫 讫 讫 讫	祝 贯 责 豈	邳 邯 邳 邳

(a) Easy styles and easy contents.

Reference:	喻	喻	喻	喻	喻
Source:	藁 藁 薰 藪	虢 靡 藜 藿	睿 毅 覲 覲	躡 躡 躡 躡	鑣 鑣 鑣 鑣
FUNIT:	藁 藁 薰 藪	虢 靡 藜 藿	睿 毅 覲 覲	躡 躡 躡 躡	鑣 鑣 鑣 鑣
MX-Font:	藁 藁 薰 藪	虢 靡 藜 藿	睿 毅 覲 覲	躡 躡 躡 躡	鑣 鑣 鑣 鑣
DG-Font:	藁 藁 薰 藪	虢 靡 藜 藿	睿 毅 覲 覲	躡 躡 躡 躡	鑣 鑣 鑣 鑣
Diff-Font:	藁 藁 薰 藪	虢 靡 藜 藿	睿 毅 覲 覲	躡 躡 躡 躡	鑣 鑣 鑣 鑣
GT:	藁 藁 薰 藪	虢 靡 藜 藿	睿 毅 覲 覲	躡 躡 躡 躡	鑣 鑣 鑣 鑣

(b) Easy styles and difficult contents.

Reference:	喻	喻	喻	喻	喻
Source:	蠹 蠹 蠹 蠹	鑫 整 整 整	鞞 鞞 鞞 鞞	鍤 鍤 鍤 鍤	鼙 鼙 鼙 鼙
FUNIT:	蠹 蠹 蠹 蠹	鑫 整 整 整	鞞 鞞 鞞 鞞	鍤 鍤 鍤 鍤	鼙 鼙 鼙 鼙
MX-Font:	蠹 蠹 蠹 蠹	鑫 整 整 整	鞞 鞞 鞞 鞞	鍤 鍤 鍤 鍤	鼙 鼙 鼙 鼙
DG-Font:	蠹 蠹 蠹 蠹	鑫 整 整 整	鞞 鞞 鞞 鞞	鍤 鍤 鍤 鍤	鼙 鼙 鼙 鼙
Diff-Font:	蠹 蠹 蠹 蠹	鑫 整 整 整	鞞 鞞 鞞 鞞	鍤 鍤 鍤 鍤	鼙 鼙 鼙 鼙
GT:	蠹 蠹 蠹 蠹	鑫 整 整 整	鞞 鞞 鞞 鞞	鍤 鍤 鍤 鍤	鼙 鼙 鼙 鼙

(c) Difficult styles and difficult contents.

Fig. 5 Example generation results on large test dataset. Easy style means the style of the reference font is similar to the source font. The characters with 10 or fewer strokes are easy contents, and those with 15 or more are difficult contents

Ref	GT	覃 巽 痛 球 瘦 稍 俟 硝 粟 痞 森 情 啼 阮 深
	MX-Font	覃 巽 痛 球 瘦 稍 俟 硝 粟 痞 森 情 啼 阮 深
	DG-Font	覃 巽 痛 球 瘦 稍 俟 硝 粟 痞 森 情 啼 阮 深
	Diff-Font	覃 巽 痛 球 瘦 稍 俟 硝 粟 痞 森 情 啼 阮 深
Ref	GT	淞 椅 毯 越 然 清 棉 涸 喃 崴 翁 婿 邵 跚 掏
	MX-Font	淞 椅 毯 越 然 清 棉 涸 喃 崴 翁 婿 邵 跚 掏
	DG-Font	淞 椅 毯 越 然 清 棉 涸 喃 崴 翁 婿 邵 跚 掏
	Diff-Font	淞 椅 毯 越 然 清 棉 涸 喃 崴 翁 婿 邵 跚 掏
Ref	GT	晴 祁 浙 惟 棠 梁 喔 液 舜 善 媒 晰 跽 渊 腆
	MX-Font	晴 祁 浙 惟 棠 梁 喔 液 舜 善 媒 晰 跽 渊 腆
	DG-Font	晴 祁 浙 惟 棠 梁 喔 液 舜 善 媒 晰 跽 渊 腆
	Diff-Font	晴 祁 浙 惟 棠 梁 喔 液 舜 善 媒 晰 跽 渊 腆
Ref	GT	惋 茱 探 茸 跹 舒 喂 胸 赧 悻 推 茹 猛 殖 散
	MX-Font	惋 茱 探 茸 跹 舒 喂 胸 赧 悻 推 茹 猛 殖 散
	DG-Font	惋 茱 探 茸 跹 舒 喂 胸 赧 悻 推 茹 猛 殖 散
	Diff-Font	惋 茱 探 茸 跹 舒 喂 胸 赧 悻 推 茹 猛 殖 散

Fig. 6 Example generation results of MX-Font, DG-Font, Diff-Font on four art fonts. It can be seen that the structure of the characters generated by MX-Font is severely distorted and the characters generated by DG-Font may contain artifacts

to column 4 and column 5 in Fig. 9, when in the case of generating a difficult structure character, Diff-Font without stroke condition and Diff-Font with one-bit encoding may generate characters with stroke errors since the number of basic strokes is not explicitly encoded. These reveals that count encoding is effective for improving the quality by preserving a completed number of strokes.

#### 4.4.2 Impact of Guidance Scales

We further discuss the impact of content and stroke on the generation by setting different content scales ( $s_1$ ) and stroke scales ( $s_2$ ). Our experiments are conducted on the test set in

	妙 联 横 生 贴 门 前
聊	妙 联 横 生 贴 门 前 ①
啤	妙 联 横 生 贴 门 前 ②
宿	妙 联 横 生 贴 门 前 ③
	妙 联 横 生 贴 门 前 ④

Fig. 7 An example for human testing. The first column shows three characters with the reference target style, and the first row lists characters with source content

**Table 3** Results of Human testing

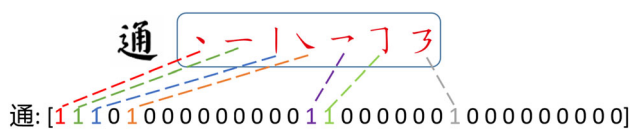
	FUNIT (%)	MX-Font (%)	DG-Font (%)	Diff-Font (%)
Content	6.41	6.41	26.09	<b>61.09</b>
Style	13.91	17.65	19.38	<b>49.06</b>
Both	13.44	13.28	4.53	<b>48.75</b>

Best results is marked in **bold**

**Table 4** Effectiveness of the stroke count encoding form versus one-bit stroke encoding

Methods	SSIM ( $\uparrow$ )	RMSE ( $\downarrow$ )	LPIPS ( $\downarrow$ )	FID ( $\downarrow$ )	ACC <sub>C</sub> ( $\uparrow$ )	ACC <sub>S</sub> ( $\uparrow$ )	ACC <sub>B</sub> ( $\uparrow$ )
w/o strokes	0.740	0.275	0.127	28.83	99.15	89.14	88.56
One-bit encoding	0.739	0.277	0.131	30.44	98.38	88.38	87.40
Count encoding	<b>0.742</b>	<b>0.271</b>	<b>0.124</b>	<b>27.30</b>	<b>99.36</b>	<b>93.05</b>	<b>92.46</b>

Best results is marked in **bold**



**Fig. 8** One-bit stroke encoding in StrokeGAN Zeng et al. (2021). Each dimension of the encoding vector indicates whether the character contains the corresponding basic stroke

ground truth	太	士	谁	直
w/o stroke	太	士	谁	直
one-bit encoding	太	士	谁	直
count encoding	太	士	谁	直

**Fig. 9** Qualitative results of ablation studies using different stroke condition. The first row is the ground truth, and from the second to the fourth row are results of Diff-Font without stroke condition, with one-bit stroke encoding, with stroke count encoding, respectively

large dataset mentioned in Sect. 4.1. In Table 5, we obtain that using the setting  $s_1 = 3, s_2 = 3$  can get the best quality generated images.

#### 4.5 Korean Script Generation

Our proposed Diff-Font is language independent, so it provides potential general solution for font generation in different languages by utilizing various attribute conditions. In this section, we evaluate the effectiveness of Diff-Font in

Korean. As illustrated in Fig. 4, the Chinese stroke condition can be substituted with the component condition of Korean.

Specifically, we collect a dataset of 201 Korean fonts, 195 for training, and the remaining 6 for testing. This dataset contains 2350 Korean characters. To evaluate the effectiveness of our proposed method, we conducted comparisons with the DG-Font and MX-Font approaches in generating 800 Korean characters and the results are presented in Table 6 and Fig. 10. We can see that our method also achieves the best results in generating Korean script.

#### 4.6 Other Script Generation

As for some simple scripts without complex structures (e.g., Latin and Greek), we can train a Diff-Font in the first stage by only using content and style attribute conditions without fine-tuning in the second stage. As shown in Fig. 11, our model is also effective in Latin and Greek font generation.

#### 4.7 Limitations

As our proposed Diff-Font is based on the denoising diffusion model, it has the same problem as most existing diffusion models with low inference efficiency. Moreover, our experimental results show that equipping with stroke/component condition for font generation could reduce generation errors, but cannot completely eliminate them. Some characters with extreme intricate structures or uncommon styles that were infrequently encountered in the training set still suffer generation failures. Some failure cases are shown in Fig. 12. In addition, Diff-Font can only generate the characters it has seen before. This limitation arises from the utilization of tokenization processes for character content, as it is now incapable to define tokens for unseen characters. However, the character set is normally finite, the character dictionary used for tokenization can cover almost all commonly used characters, as shown in the experimental setting of Sect. 4.1.

**Table 5** Impact of guidance scales

Scales	SSIM (↑)	RMSE (↓)	LPIPS (↓)	FID (↓)	ACC <sub>C</sub> (↑)	ACC <sub>S</sub> (↑)	ACC <sub>B</sub> (↑)
$s_1 = 1, s_2 = 1$	0.720	0.280	0.108	16.67	93.59	95.16	89.13
$s_1 = 1, s_2 = 3$	0.720	0.281	0.112	16.88	93.26	95.35	89.04
$s_1 = 1, s_2 = 5$	0.716	0.285	0.120	17.16	92.16	95.50	88.23
$s_1 = 3, s_2 = 1$	<b>0.722</b>	0.279	0.105	16.36	95.60	95.28	91.11
$s_1 = 3, s_2 = 3$	<b>0.722</b>	<b>0.277</b>	<b>0.104</b>	<u>16.20</u>	95.78	<b>96.55</b>	<b>92.48</b>
$s_1 = 5, s_2 = 1$	0.720	0.280	0.107	<b>16.18</b>	<b>95.85</b>	93.45	89.49
$s_1 = 5, s_2 = 3$	0.721	<u>0.278</u>	<b>0.104</b>	16.27	<u>95.84</u>	<u>96.43</u>	<b>92.48</b>

The best and second-best results are marked in **bold** and underlined, respectively

**Table 6** Quantitative results on Korean script

Methods	SSIM (↑)	RMSE (↓)	LPIPS (↓)	FID (↓)	ACC <sub>C</sub> (↑)	ACC <sub>S</sub> (↑)	ACC <sub>B</sub> (↑)
MX-Font	0.691	0.278	0.158	47.05	93.52	48.50	44.85
DG-Font	0.771	0.235	0.095	43.36	92.81	80.33	73.52
Diff-Font	<b>0.812</b>	<b>0.196</b>	<b>0.072</b>	<b>10.69</b>	<b>94.83</b>	<b>99.13</b>	<b>94.13</b>

Best results are marked in **bold**

Ref	갱	갱
MX-Font	갱 곳 각 갠	갱 곳 갠 갠
DG-Font	갱 곳 각 갠	갱 곳 갠 갠
Diff-Font	갱 곳 각 갠	갱 곳 갠 갠
GT	갱 곳 각 갠	갱 곳 갠 갠
Ref	갠	갱
MX-Font	갠 겹 겹 갠	갱 끊 꿩 갈
DG-Font	갠 겹 겹 갠	갱 끊 꿩 갈
Diff-Font	갠 겹 겹 갠	갱 끊 꿩 갈
GT	갠 겹 겹 갠	갱 끊 꿩 갈
Ref	갱	갱
MX-Font	갱 뵈 뵈 뵈	갠 속 쉼 쉼
DG-Font	갱 뵈 뵈 뵈	갠 속 쉼 쉼
Diff-Font	갱 뵈 뵈 뵈	갠 속 쉼 쉼
GT	갱 뵈 뵈 뵈	갠 속 쉼 쉼

**Fig. 10** Qualitative results on Korean script

GT	A B C D S T U V
Diff-Font	A B C D S T U V
GT	w x y z l k π ρ
Diff-Font	w x y z l k π ρ

**Fig. 11** Example generation results of Diff-Font on Latin and Greek

GT	鞞 靄 鋹 鑿 醪 醴 醴 醴
Diff-Font	鞞 靄 鋹 鑿 醪 醴 醴 醴
GT	豐 邨 躑 鞞 讖 溲 審 攀
Diff-Font	豐 邨 躑 鞞 讖 溲 審 攀

**Fig. 12** Some failure cases. Characters with extreme complex structures or uncommon styles still suffer generation failures

Therefore, Diff-Font is able to generate a comprehensive set of commonly used characters. Moreover, we have noticed that continual learning can expand the task scope of the model. In the future work, we will investigate leveraging this technology to endow Diff-Font with the ability to generate unseen characters.

## 5 Conclusion

In this paper, we propose a unified method based on the diffusion model, namely Diff-Font, for one-shot font generation task. The proposed Diff-Font has a stable training process and can be well-trained on large datasets. To address the problems of unsatisfactory generation results on large or subtle differences in the style of source font and target font faced by previous GANs-based methods, we regard font generation as a conditional generation task and generate the corresponding character images according to the given character attribute conditions. Furthermore, we introduce stroke- and component-wise information to improve the structural integrity of generated characters and solve the problem of low generation quality of complicated characters for Chinese and Korean generation. The remarkable performance on two datasets with different scales shows the effectiveness of Diff-Font.

**Acknowledgements** This work was supported in part by the National Key Research and Development Program of China under Grant 2023YFC2705700, in part by the National Natural Science Foundation of China under Grants U23B2048, 62076186, 62225113, and 62102150, and in part by the Innovative Research Group Project of Hubei Province under Grant 2024AFA017. The numerical calculations in this paper have been done on the supercomputing system in the Supercomputing Center of Wuhan University.

## References

- An, J., & Cho, S. (2015). Variational autoencoder based anomaly detection using reconstruction probability. *Special Lecture on IE*, 2(1), 1–18.
- Baek, K., Choi, Y., Uh, Y., Yoo, J., & Shim, H. (2021). Rethinking the truly unsupervised image-to-image translation. In *Proceedings of the IEEE/CVF international conference on computer vision* (pp. 14154–14163).
- Cha, J., Chun, S., Lee, G., Lee, B., Kim, S., & Lee, H. (2020). Few-shot compositional font generation with dual memory. In *European conference on computer vision* (pp. 735–751). Springer.
- Cheng, S. I., Chen, Y. J., Chiu, W. C., Tseng, H. Y., & Lee, H. Y. (2023). Adaptively-realistic image generation from stroke and sketch with diffusion model. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision* (pp. 4054–4062).
- Choi, J., Kim, S., Jeong, Y., Gwon, Y., & Yoon, S. (2021). Ilvr: Conditioning method for denoising diffusion probabilistic models. [arXiv:2108.02938](https://arxiv.org/abs/2108.02938)
- Cui, Y., Che, W., Liu, T., Qin, B., & Yang, Z. (2021). Pre-training with whole word masking for chinese bert. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29, 3504–3514.
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. [arXiv:1810.04805](https://arxiv.org/abs/1810.04805)
- Dhariwal, P., & Nichol, A. (2021). Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34, 8780–8794.
- Gao, Y., Guo, Y., Lian, Z., Tang, Y., & Xiao, J. (2019). Artistic glyph image synthesis via one-stage few-shot learning. *ACM Transactions on Graphics (TOG)*, 38(6), 1–12.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2020). Generative adversarial networks. *Communications of the ACM*, 63(11), 139–144.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., & Hochreiter, S. (2017). GANs trained by a two time-scale update rule converge to a local Nash equilibrium. *Advances in Neural Information Processing Systems*, 30.
- Ho, J., & Salimans, T. (2022). Classifier-free diffusion guidance. [arXiv:2207.12598](https://arxiv.org/abs/2207.12598)
- Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33, 6840–6851.
- Huang, X., & Belongie, S. (2017). Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE international conference on computer vision* (pp. 1501–1510).
- Huang, X., Liu, M. Y., Belongie, S., & Kautz, J. (2018). Multimodal unsupervised image-to-image translation. In *Proceedings of the European conference on computer vision (ECCV)* (pp. 172–189).
- Isola, P., Zhu, J. Y., Zhou, T., Efros, A. A. (2017). Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1125–1134).
- Jiang, Y., Lian, Z., Tang, Y., & Xiao, J. (2019). Sfont: Structure-guided chinese font generation via deep stacked networks. In *Proceedings of the AAAI conference on artificial intelligence* (pp. 4015–4022).
- Kancharugunta, K. B., & Dubey, S. R. (2019). Csgan: Cyclic-synthesized generative adversarial networks for image-to-image transformation. [arXiv:1901.03554](https://arxiv.org/abs/1901.03554)
- Kim, T., Cha, M., Kim, H., Lee, J. K., & Kim, J. (2017). Learning to discover cross-domain relations with generative adversarial networks. In *International conference on machine learning, PMLR* (pp. 1857–1865).
- Kong, Y., Luo, C., Ma, W., Zhu, Q., Zhu, S., Yuan, N., & Jin, L. (2022). Look closer to supervise better: One-shot font generation via component-based discriminator. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 13482–13491).
- Li, B., Xue, K., Liu, B., & Lai, Y. K. (2022). VQBB: Image-to-image translation with vector quantized Brownian bridge. [arXiv:2205.07680](https://arxiv.org/abs/2205.07680)
- Liu, M. Y., & Tuzel, O. (2016). Coupled generative adversarial networks. *Advances in Neural Information Processing Systems*, 29.
- Liu, M. Y., Breuel, T., & Kautz, J. (2017). Unsupervised image-to-image translation networks. *Advances in Neural Information Processing Systems*, 30.
- Liu, M. Y., Huang, X., Mallya, A., Karras, T., Aila, T., Lehtinen, J., & Kautz, J. (2019). Few-shot unsupervised image-to-image translation. In *Proceedings of the IEEE/CVF international conference on computer vision* (pp. 10551–10560).
- Liu, N., Li, S., Du, Y., Torralba, A., & Tenenbaum, J. B. (2022). Compositional visual generation with composable diffusion models. [arXiv:2206.01714](https://arxiv.org/abs/2206.01714)
- Liu, X., Park, D. H., Azadi, S., Zhang, G., Chopikyan, A., Hu, Y., Shi, H., Rohrbach, A., & Darrell, T. (2021). More control for free! image synthesis with semantic diffusion guidance. [arXiv:2112.05744](https://arxiv.org/abs/2112.05744)
- Mirza, M. & Osindero, S. (2014). Conditional generative adversarial nets. In *CoRR*.
- Nair, N. G., Bandara, W. G. C., Patel, V. M. (2022). Image generation with multimodal priors using denoising diffusion probabilistic models. [arXiv:2206.05039](https://arxiv.org/abs/2206.05039)
- Nichol, A., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., McGrew, B., Sutskever, I., & Chen, M. (2021). Glide: Towards photorealistic image generation and editing with text-guided diffusion models. [arXiv:2112.10741](https://arxiv.org/abs/2112.10741)
- Park, S., Chun, S., Cha, J., Lee, B., & Shim, H. (2021a) Few-shot font generation with localized style representations and factorization.

- In *Proceedings of the AAAI conference on artificial intelligence* (pp. 2393–2402).
- Park, S., Chun, S., Cha, J., Lee, B., & Shim, H. (2021b) Multiple heads are better than one: Few-shot font generation with multiple localized experts. In *Proceedings of the IEEE/CVF international conference on computer vision* (pp. 13900–13909).
- Park, S., Chun, S., Cha, J., Lee, B., & Shim, H. (2022). Few-shot font generation with weakly supervised localized representations. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 1–17.
- Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., & Chen, M. (2022). Hierarchical text-conditional image generation with clip latents. [arXiv:2204.06125](https://arxiv.org/abs/2204.06125)
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 10684–10695).
- Saharia, C., Chan, W., Chang, H., Lee, C., Ho, J., Salimans, T., Fleet, D., & Norouzi, M. (2022a). Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 conference proceedings* (pp. 1–10).
- Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E., Ghasemipour, S. K. S., Ayan, B. K., Mahdavi, S. S., Lopes, R. G., Salimans, T., & Ho, J. (2022b). Photorealistic text-to-image diffusion models with deep language understanding. [arXiv:2205.11487](https://arxiv.org/abs/2205.11487)
- Sasaki, H., Willcocks, C. G., & Breckon, T. P. (2021). Unit-ddpm: Unpaired image translation with denoising diffusion probabilistic models. [arXiv:2104.05358](https://arxiv.org/abs/2104.05358)
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., & Ganguli, S. (2015). Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning, PMLR* (pp. 2256–2265).
- Song, J., Meng, C., & Ermon, S. (2020). Denoising diffusion implicit models.
- Tang, L., Cai, Y., Liu, J., Hong, Z., Gong, M., Fan, M., Han, J., Liu, J., Ding, E., & Wang, J. (2022). Few-shot font generation by learning fine-grained local styles. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 7895–7904).
- Tian, Y. (2017). zi2zi: Master Chinese calligraphy with conditional adversarial networks. *Internet* <https://github.com/kaonashi-tyc/zi2zi>, 3.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M. A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., & Rodriguez, A. (2023). Llama: Open and efficient foundation language models. [arXiv:2302.13971](https://arxiv.org/abs/2302.13971)
- Wang, Z., Bovik, A. C., Sheikh, H. R., & Simoncelli, E. P. (2004). Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4), 600–612.
- Wen, Q., Li, S., Han, B., & Yuan, Y. (2021). Zigan: Fine-grained Chinese calligraphy font generation via a few-shot style transfer approach. In *Proceedings of the 29th ACM international conference on multimedia* (pp. 621–629).
- Wolleb, J., Sandkühler, R., Bieder, F., & Cattin, P. C. (2022). The swiss army knife for image-to-image translation: Multi-task diffusion models. [arXiv:2204.02641](https://arxiv.org/abs/2204.02641)
- Xie, Y., Chen, X., Sun, L., & Lu, Y. (2021). Dg-font: Deformable generative networks for unsupervised font generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 5130–5140).
- Yang, X., Xie, D., & Wang, X. (2018). Crossing-domain generative adversarial networks for unsupervised multi-domain image-to-image translation. In *Proceedings of the 26th ACM international conference on multimedia* (pp. 374–382).
- Yi, Z., Zhang, H., Tan, P., & Gong, M. (2017). Dualgan: Unsupervised dual learning for image-to-image translation. In *Proceedings of the IEEE international conference on computer vision* (pp. 2849–2857).
- Zeng, J., Chen, Q., Liu, Y., Wang, M., & Yao, Y. (2021). Strokegan: Reducing mode collapse in Chinese font generation via stroke encoding. In *Proceedings of the AAAI conference on artificial intelligence* (pp. 3270–3277).
- Zhang, L., & Agrawala, M. (2023). Adding conditional control to text-to-image diffusion models. [arXiv:2302.05543](https://arxiv.org/abs/2302.05543)
- Zhang, R., Isola, P., Efros, A. A., Shechtman, E., & Wang, O. (2018a). The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 586–595).
- Zhang, Y., Zhang, Y., & Cai, W. (2018b). Separating style and content for generalized style transfer. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8447–8455).
- Zhao, M., Bao, F., Li, C., & Zhu, J. (2022). Egsde: Unpaired image-to-image translation via energy-guided stochastic differential equations. [arXiv:2207.06635](https://arxiv.org/abs/2207.06635)
- Zhu, J. Y., Park, T., Isola, P., & Efros, A. A. (2017a). Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision* (pp. 2223–2232).
- Zhu, J. Y., Zhang, R., Pathak, D., Darrell, T., Efros, A. A., Wang, O., & Shechtman, E. (2017b). Toward multimodal image-to-image translation. *Advances in Neural Information Processing Systems*, 30.

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.