Compositional Zero-Shot Artistic Font Synthesis

Xiang Li¹, Lei Wu^{2*}, Changshuo Wang¹, Lei Meng^{2*}, Xiangxu Meng²

School of Software, Shandong University

{202035260, 202115242}@mail.sdu.edu.cn, {i_lily, lmeng, mxx}@sdu.edu.cn

Abstract

Recently, many researchers have made remarkable 2 achievements in the field of artistic font synthe-3 sis, with impressive glyph style and effect style 4 in the results. However, due to less exploration 5 in style disentanglement, it is difficult for existing 6 methods to envision a kind of unseen style (glyph-7 effect) compositions of artistic font, and thus can 8 only learn the seen style compositions. To solve 9 this problem, we propose a novel compositional 10 zero-shot artistic font synthesis gan (CAFS-GAN), 11 which allows the synthesis of unseen style compo-12 sitions by exploring the visual independence and 13 joint compatibility of encoding semantics between 14 glyph and effect. Specifically, we propose two 15 contrast-based style encoders to achieve style dis-16 entanglement due to glyph and effect intertwin-17 ing in the image. Meanwhile, to preserve more 18 glyph and effect detail, we propose a generator 19 based on hierarchical dual styles AdaIN to reorga-20 nize content-styles representations from structure 21 to texture gradually. Extensive experiments demon-22 strate the superiority of our model in generating 23 high-quality artistic font images with unseen style 24 compositions against other state-of-the-art meth-25 ods. The source code and data will be publicly 26 available. 27

28 1 Introduction

Artistic fonts are frequently employed in signboards, posters, 29 magazines, and web pages, playing an integral role in cap-30 tivating and sustaining the audience's attention. The com-31 pelling nature of these fonts lies in the fact that designers 32 meticulously craft visually appealing and harmonious glyph 33 and effect styles that suit the occasion and theme. In the 34 course of design, the designers draw upon design theory and 35 aesthetic factors to conceive various style elements, often re-36 quiring only a momentary mental picture. It is worth noting 37 that if we can provide a deep learning model with enough 38 glyph styles and effect styles as prior knowledge, whether the 39



Figure 1: We aim to build an artistic font synthesis model for synthesizing unseen style compositions (e.g., Arial-Cookie) by training with some seen concepts, such as Century-Cookie, Century-Metal, and Arial-Metal.

model can also design a kind of artistic font with unseen integrated style like humans.

40

41

In order to achieve the automatic synthesis of artistic font 42 based on deep learning, some conventional methods [Azadi 43 et al., 2018; Gao et al., 2019; Li et al., 2020a] focus on the 44 integrated style (glyph-effect) transfer and generate an artis-45 tic font library with the existing style. These works treat the 46 style of artistic fonts as a whole and generalize the learned 47 integrated style to any character content. However, they 48 ignore the independence and decoupling of styles, making 49 these methods ineffective in scenarios where glyph and effect 50 styles must be controlled separately. Therefore, the conven-51 tional methods cannot synthesize artistic fonts with unseen 52 style (glyph-effect) compositions. There are also some recent 53 works [Ge et al., 2021; Li et al., 2022a] that propose learning 54 disentangled style representations and synthesizing content-55 glyph-effect controllable artistic font images. Unfortunately, 56 these methods focus on the seen style compositions and must 57 require a large amount of data paired with the three attributes 58

^{*}Corresponding author

of content, glyph, and effect. Due to pixel-level supervision 59 information, these methods inevitably focus on pixel-level re-60 lationship instead of creating the new style compositions, re-61 sulting in the generated images with a messy structure and 62 unclear texture. 63 In this paper, we propose a novel and practical task, called 64 compositional zero-shot artistic font synthesis (CAFS), which 65 focuses on unseen style composition synthesis, see Figure 1. 66 It aims to learn the compositionality of glyphs and effects 67

from the training set and is tasked with generalizing to un from the training set and is tasked with generalizing to un seen style (glyph-effect) compositions on any character. To
 realize this task, we propose a new model, CAFS-GAN, from
 the perspective of style disentanglement and content-styles
 representations reorganization.

For the style disentanglement, we propose two contrast-73 based style encoders, glyph encoder and effect encoder, 74 which implement glyph and effect disentanglement and pre-75 cise style feature extraction. The key idea is that we intro-76 duce glyph style contrastive loss and effect style contrastive 77 loss to learn the style commonalities and differences. For 78 the content-styles representations reorganization, we propose 79 an artistic font generator based on hierarchical dual styles 80 AdaIN, which progressively feeds glyph and effect informa-81 tion to preserve more image details. The key idea is that 82 the hierarchical dual styles AdaIN completes the composi-83 tion of glyph and content in the high-dimensional AdaIN 84 layer, and the composition of effect and content in the low-85 dimensional AdaIN layer. Moreover, to enable the model 86 to synthesize artistic font images with controllable style at-87 tributes, we adopt the well-known GAN [Goodfellow et al., 88 2014] framework and introduce two multi-task discrimina-89 tors, glyph discriminator and effect discriminator that con-90 strain the style of the generated glyphs and effects, respec-91 tively. Finally, to comprehensively evaluate the generated re-92 sults, we propose two evaluation metrics: glyph outline mis-93 alignment (GOLM) and effect perception error (EPE). 94

⁹⁵ In summary, our contributions are as follows:

 We propose a novel compositional zero-shot artistic font synthesis gan (CAFS-GAN) to synthesize unseen style compositions for artistic font images. Meanwhile, our model supports the control of artistic font synthesis from three aspects (i.e., glyph, effect, and content).

We propose two new evaluation metrics, called glyph outline misalignment (GOLM) and effect perception error (EPE), which enrich the evaluation methods from the unique attribute of the artistic font.

 Extensive experiments demonstrate the effectiveness and superiority of our model in synthesizing unseen style compositions in Chinese standard, creative, handwriting, calligraphy artistic fonts and English artistic fonts.

110 2 Related Work

111 2.1 Artistic Font Generation

Early artistic font generation approaches are based on the high regularity of the spatial distribution for effects. T-Effects [Yang *et al.*, 2016] and DynTypo [Men *et al.*, 2019] focus on texture and special effects for synthesizing complex and re-115 alistic artistic font images. TET-GAN [Yang et al., 2019a] 116 and ShapeMatching-GAN [Yang et al., 2019b] establish the 117 mapping between the original shape and the effect, using the 118 CNN (Convolutional Neural Network) to realize the text ef-119 fect transfer. Then, AGIS-Net [Gao et al., 2019] and FET-120 GAN [Li et al., 2020a] attempt the synchronous style transfer 121 of glyphs and effects of arbitrary characters or symbols. Re-122 cently, DSE-Net [Li et al., 2022a] and GZS-Net [Ge et al., 123 2021] have conducted separate studies on the glyph structure 124 and effects of artistic fonts. Although these methods sepa-125 rately encodes artistic font glyph and effects, they still have 126 a significant data dependency on paired data. These mod-127 els learn to synthesize artistic fonts by training on paired 128 seen style combinations. Therefore, the optimization pro-129 cess for the model parameters is based on the pixel-level er-130 ror between the generated and real images, which causes the 131 model to focus excessively on pixel-level mapping relation-132 ships. This makes it difficult for the models to create new 133 style combinations. 134

2.2 Disentangled Representation Learning

Disentangled representation learning aims to infer latent fac-136 tors for a given object in the real world, where each latent 137 factor is responsible for generating a semantic feature [Han 138 et al., 2021; Yang et al., 2021; Saini et al., 2022]. Following 139 VAE, [Higgins et al., 2017] introduces β -VAE to discover in-140 terpretable latent factor representations in a completely un-141 supervised manner. [Chen *et al.*, 2018] improved β -VAE, 142 and further proposed a principled classifier-free measure of 143 disentanglement. Recently, a large amount of works [Zhang 144 et al., 2018; Li et al., 2020b; Luo et al., 2022] have made 145 great contributions to disentangled shape and texture, unfor-146 tunately, they are unable to generate novel combinations not 147 witnessed during training. 148

135

149

2.3 Compositional Zero-Shot Learning

Compositional zero-shot learning stands at the intersection of 150 compositionality and zero-shot learning and focuses on state 151 and object relations. Compositionality [Naeem et al., 2021] 152 can loosely be defined as the ability to decompose an observa-153 tion into its primitives. Zero-shot learning [Gao et al., 2018; 154 Hong et al., 2022; Feng et al., 2022; Lin et al., 2022] aims 155 at recognizing or generating novel classes that are not ob-156 served during training. Recently, [Yang et al., 2022] present 157 a novel decomposable causal view that characterizes how 158 compositional concepts are formed. [Karthik et al., 2022; 159 Mancini et al., 2021] propose to address the problem of open-160 world compositional zero-shot learning. [Li et al., 2022b] 161 propose a novel siamese contrastive embedding network to 162 excavate discriminative prototypes of state and object. 163

In this paper, we propose a compositional zero-shot artistic font synthesis, and use the artistic font's glyph and effect style as attribute primitives. More importantly, our method is the first to estimate the unseen style compositions, and uses the joint compatibility and differences between the two styles to synthesize and optimize the detailed characteristics of the image styles.



Figure 2: The overview of proposed CAFS-GAN. The encoding process of CAFS-GAN has three input channels (1) effect sample s_{x_i} with effect attribute x_i , (2) glyph sample s_{y_i} with glyph attribute y_j , and (3) content sample s_{z_k} with content attribute z_k . For the style encoders, we additionally input positive samples $s_{x_i}^+$ and $s_{y_j}^+$ of style reference images. SSA integrates the style features of style samples and their positive samples from style encoders. The generator utilizes the hierarchical dual styles AdaIN architecture to reorganize the input content, effects, and glyph signals. The discriminator outputs a one-shot vector. The outputs of the discriminators in different channels indicate whether the generated image comes from the domain corresponding to this channel.

171 **3 Method**

172 **3.1 Problem Define**

Compositional zero-shot artistic font synthesis (CAFS) aims 173 to predict an unseen style composition, namely to synthesize 174 glyph-effect compositions that do not exist in the training set 175 and map it to any character to obtain a complete artistic font 176 library. Let us denote with $\mathcal{X} = \{x_i\}_{i=1}^{N_x}$ the set of effect attributes, with $\mathcal{Y} = \{y_j\}_{j=1}^{N_y}$ the set of glyph attributes, with 177 178 $\mathcal{Z} = \{z_k\}_{k=1}^{N_z}$ the set of characters, and with $\mathcal{C} = \mathcal{X} \times \mathcal{Y}$ 179 the set of all their possible compositions. $\mathcal{T} = \{\mathcal{Z}_t, \mathcal{C}_t\}$ is 180 a training set where \mathcal{Z}_t is a character set seen during train-181 ing $(\mathcal{Z}_t \subseteq \mathcal{Z})$ and \mathcal{C}_t is a style compositions set seen during 182 training ($C_t \subseteq C$). When the glyph and effect elements in C_t 183 covers all elements in \mathcal{X} and \mathcal{Y} , \mathcal{T} can be used to train the 184 model $f : \{\mathcal{Z}_t, \mathcal{C}_t\} \to \{\mathcal{Z}_t, \mathcal{C}_u\}$ synthesizing the artistic 185 font images with unseen style combinations where $\mathcal{C}_u \subset \mathcal{C}$ 186 denote the unseen style compositions and $C_t \cup C_u = C$. 187

The difficulty of the CAFS task varies depending on the 188 proportion of the C_t . If the style compositions in C_t covers 189 all compositions and $C_u \equiv \emptyset$, the task definition is the same 190 as the conventional artistic font generation task, where the 191 model only needs to predict the seen style combination on ar-192 bitrary character content. In the case of $C_t \subset C$, since the 193 model only learns jointly compatibility of encoding seman-194 tics between glyph and effect in seen style compositions, it 195 is very challenging to predict unseen style combinations. It 196

is worth noting that as the C_t shrinks, the training data can provide the model with fewer data on the joint compatibility relationship of glyph and effect. In this case, the shrink of composition information hinders the recognizability of glyph and effect, making it difficult for the model to predict unseen style combinations. Regarding this hypothesis, we verified it in Experiment 5.5.

204

3.2 Overview of CAFS-GAN

The CAFS-GAN consists of the following modules: two style 205 encoders E_x and E_y , two style similarity attention modules, 206 a content encoder E_z , an artistic font generator G, and two 207 style discriminators D_x and D_y , as shown in Figure 2. First, 208 E_x and E_y represent effect style encoder and glyph style en-209 coder, respectively, which are used to disentangle and extract 210 glyph and effect style features. At the end of the two style 211 encoders, we add a style similarity attention (SSA) module, 212 which uses the similarity of style attributes to enhance the 213 model's perception of various glyphs or effects. The struc-214 ture details of E_x and E_y are similar to VGG11 [Simonyan 215 and Zisserman, 2014]. Unlike E_x and E_y , our E_z adds sev-216 eral padding layers to increase the sampling times for the font 217 strokes at the image's edge. This operation protects the in-218 tegrity of the character structure. In addition, since the con-219 tent information of characters belongs to high-dimensional 220 semantic information, we add resblocks at the end of the con-221 tent encoder to retain more content information. Lastly, our 222



Figure 3: Two contrast-based style encoders. s_{x_i} and s_{y_j} represent effect and glyph samples. $s_{x_i}^+$ and $s_{y_j}^+$ have the same visual attributes as s_{x_i} and s_{y_j} in the corresponding attributes, respectively, and vice versa for effect and glyph negative samples.

223 D_x and D_y are two multi-task discriminators consisting of 224 FRN (Filter Response Normalization) [Singh and Krishnan, 2020] and convolutional layer, which consists of multiple out-226 put branches. Each branch learns a binary classification de-227 termining whether an artistic font has real glyph style or real 228 effect style.

In the next sections, we will look at the two aspects: style (glyph-effect) disentanglement (in sections 3.3 and 3.4) and content-styles representations reorganization (in section 3.5).

233 3.3 Contrast-Based Style Encoders

In the process of achieving the CAFS task, the style encoders 234 need to provide the generator with disentangled glyph fea-235 tures and effect features. However, the actual situation is 236 that the visual elements of effect, glyph, and content are en-237 tangled, and the commonly used data enhancement methods 238 cannot eliminate or highlight a certain visual element. There-239 fore, we introduce a contrastive learning [He et al., 2020] 240 strategy to encourage encoders to identify deep similarities 241 and differences between the two style attributes. Taking the 242 pipeline of effect extraction as an example, we define $s^+_{x_i}$ and $\{s^-_{1,x_i}, s^-_{2,x_i}, ..., s^-_{N_x-1,x_i}\}$ as the positive sample and nega-243 244 tive sample set of the original input s_{x_i} , respectively. N_x 245 denotes the number of all kinds of effect styles, one of which 246 is the effect of positive samples, and $N_x - 1$ is the number 247 of all kinds of negative effects. The positive pair $(s_{x_i}, s_{x_i}^+)$ 248 only shares the same effect, and the negative pair (s_{x_i}, s_{r,x_i}^-) 249 have different effects $(1 \le r \le N_x - 1)$, as shown in Figure 250 3. At this time, we utilize the effect style contrastive loss to 251 enhance the effect similarity between positive pairs and the 252 dissimilarity between negative pairs: 253

$$\mathcal{L}_{sty}^{E_x} = -\log \frac{\exp(f_{x_i} \cdot f_{x_i}^+ / \tau)}{\sum_{r=1}^{N_x - 1} \exp(f_{x_i} \cdot f_{r, x_i}^- / \tau)},$$
(1)

where $f_{x_i}, f_{x_i}^+, f_{r,x_i}^-$ are effect features obtained by s_x, s_x^+, s_{r,x_i}^- through E_x . Similarly, we also impose the glyph style contrastive loss $\mathcal{L}_{sty}^{E_y}$ to improve the glyph encoder. Furthermore, the total style contrastive loss can be defined as: _____ 258

$$\mathcal{L}_{sty} = \mathcal{L}_{sty}^{E_x} + \mathcal{L}_{sty}^{E_y}.$$
 (2)

3.4 Style Similarity Attention

To make full use of the style similarity between positive samples and original samples as auxiliary information for synthesizing disentangled style features, we introduce a style similarity attention module at the end of the style encoders. Specifically, we use the style features of positive samples as K and V, and use the style features of original images as Q. Style similarity attention can be expressed as: 260

$$SSA(Q, K, V) = softmax(\frac{f \cdot f^{+T}}{\sigma})f^{+}, \qquad (3)$$

where f, f^+ are style features from the original image and positive sample, and σ factor follows Attention Mechanism [Vaswani *et al.*, 2017] to prevent the magnitude of the dot product from growing extreme. 270

Overall, our proposed contrast-based style encoders encourage the encoders to have more robust style disentanglement ability. The SSA enhances the prominent glyph-effect characteristics by amplifying the specific style signal strength to obtain a pure glyph or effect representation. 275

3.5 Hierarchical Dual Styles AdaIN

Since neural networks are easier to retain abstract informa-277 tion in high-dimensional layers and easier to retain color in-278 formation in low-dimensional layers [Gatys et al., 2016], we 279 propose an artistic font generator based on hierarchical dual 280 styles AdaIN. Specifically, we pass the disentangled glyph 281 features and effect features through a fully connected layer 282 (FC) to obtain high- and low-dimensional glyph style sig-283 nals, respectively. Here, we input the glyph signal into the 284 AdaIN layer [Huang and Belongie, 2017] of the generator 285 and fuse the content information through high-dimensional 286 connections, so that the generator can determine the overall 287 outline and structural pattern in the early stage of genera-288 tion. Furthermore, the effect signal is input to the genera-289 tor through low-dimensional connections to render the color 290 and texture details of the artistic font based on the established 291 glyph. Formally, we use the style encoders and SSA to ex-292 tract the effect feature f_{x_i} and glyph features f_{y_i} , and input 293 them to the fully connected layer. The fully connected layer 294 aims to align f_{x_i} and f_{y_j} with the channel means and variances of the content inputs f_{z_k} , and to use f_{x_i} and f_{y_j} as the 295 296 adaptive affine parameters of the AdaIN layer (i.e., w and 297 b). Ultimately, we achieve a progressive reorganization of the 298 content with glyph and effect using hierarchical dual styles 299 AdaIN: 300

$$f_{z_k}^{l+1} = \begin{cases} & w_{y_j}(\frac{f_{z_k}^l - \mu}{\sigma}) + b_{y_j} , \quad l \le h \\ & w_{x_i}(\frac{f_{z_k}^l - \mu}{\sigma}) + b_{x_i} , \quad l > h \end{cases}$$
(4)

where l denotes the current layer number and h denotes the threshold for dividing the high-dimensional AdaIN layers and the low-dimensional AdaIN layers. 303

259

276

Methods	Disentangled Style	Training	$L_1 \log \downarrow$	FID \downarrow	SSIM \uparrow	GOLM \downarrow	EPE \downarrow		
Non-zero-shot methods for synthesizing seen style compositions									
AGIS-Net [Gao et al., 2019]	×	paired	0.2277	107.01	0.4313	81.025	4.3981		
FET-GAN [Li <i>et al.</i> , 2020a]	×	paired	0.2005	100.56	0.4474	68.820	7.5113		
StarGANv2 [Choi et al., 2020] ×		unpaired	0.2997	72.24	0.3647	82.934	3.7708		
Zero-shot methods for synthesizing unseen style compositions									
GZS-Net [Ge et al., 2021]	\checkmark	paired	0.2460	140.35	0.3648	87.328	7.2335		
DSE-Net [Li et al., 2022a]	\checkmark	paired	0.1754	72.19	0.4428	83.345	3.7332		
Ours	\checkmark	unpaired	0.1271	64.79	0.5883	73.225	3.0734		

Table 1: Quantitative comparison of the CAFS-GAN and the existing state-of-the-art methods.

304 **3.6 Full objective**

305 Our full objective functions can be summarized as follows:

$$\min_{G,E} \max_{D} \lambda_{sty} \mathcal{L}_{sty} + \lambda_{adv} \mathcal{L}_{adv}^{x} + \lambda_{adv} \mathcal{L}_{adv}^{x}, \quad (5)$$

where λ_{sty} and λ_{adv} are hyperparameters. The \mathcal{L}_{adv}^{x} and \mathcal{L}_{adv}^{y} denote two adversarial loss terms for the effect discriminator and glyph discriminator:

$$\mathcal{L}_{adv}^{x} = \mathbb{E}[\log D_{x_{i}}(s_{x_{i}}) + \log(1 - D_{x_{i}}(s_{x_{i},y_{j},z_{k}}))], \quad (6)$$

$$\mathcal{L}_{adv}^{y} = \mathbb{E}[\log D_{y_j}(s_{y_i}) + \log(1 - D_{y_j}(s_{x_i, y_j, z_k}))], \quad (7)$$

where $D_{x_i}(\cdot)$ and $D_{y_j}(\cdot)$ denote the logits from the domainspecific (x_i) effect discriminator and domain-specific (y_j) glyph discriminator. s_{x_i,y_j,z_k} denote the generated artistic font image with three specific attributes.

314 4 Metrics

309

In order to better evaluate the generated glyphs and effects, we propose two kinds of new quantitative measures, GOLM for glyph and EPE for effect. Meanwhile, we also use three classic quantitative measures, such as L_1 , SSIM, and FID.

Glyph outline misalignment (GOLM). GOLM is used to 319 evaluate whether the edge information of the generated artis-320 tic font is correct and complete. Firstly, we convert the images 321 I to its grayscale I_{gray} , and calculate horizontal and vertical 322 directions gradients using the Sobel operator. By summing 323 the root mean square of the gradients in the two directions, 324 we can get the final gradient of each pixel. The formula for 325 GOLM is as follows: 326

$$GOLM = \left| I_{edge} - I'_{edge} \right|, \tag{8}$$

327

$$I_{edge} = \sqrt{(A \cdot I_{gray})^2 + (B \cdot I_{gray})^2}, \tag{9}$$

where I_{edge} and I'_{edge} denote the edge image of the real image and generated image. A and B denote horizontal and vertical Sobel matrixs.

Effect perception error (EPE). The visual communication of effect is often presented in the form of texture in artistic font images. EPE is used to evaluate whether the texture information of the generated image is accurate. First, we use the VGG19 [Simonyan and Zisserman, 2014] network to calculate the feature maps of the image in the deep layers, and then obtain the texture gram matrix [Gatys *et al.*, 2016] ³³⁷ through the inner product operation to represent the texture ³³⁸ features. Then, EPE can be formulated as follows: ³³⁹

$$EPE = \frac{1}{n} \sum_{i=1}^{n} (\mathcal{G}_i - \mathcal{G}'_i)^2,$$
 (10)

344

345

353

365

where n denotes the number of network layers involved in the calculation of feature maps, \mathcal{G}_i and \mathcal{G}'_i denote the gram matrixs calculated in the *i*-layer network of the real image and the generated image. 342

Experiments

5.1 Datasets

5

SSAF Dataset. SSAF [Li *et al.*, 2022a] contains a large number of high-quality Chinese and English artistic images, with annotations for their glyphs, effects, and content. **Fonts Dataset.** Fonts [Ge *et al.*, 2021] is a computer generated RGB font image dataset. It consists of 52 English letters with 5 independent attributes: letter identity, font size, font color, background color, and glyph. 352

5.2 Implementation Details

In our experiments, all images are resized to 128×128 pixels. 354 The hyperparameters are set as: $\lambda_{adv} = 1.0$ and $\lambda_{sty} = 0.1$. 355 In training, we set the batch size as 8 and train 10^5 iterations 356 for Chinese artistic font generation and 2×10^4 iterations for 357 English. The learning rate is set to 0.0001, using Adam op-358 timizer. Regarding the division of all possible style composi-359 tions, we set the proportion of the number of style composi-360 tions in C_u to C_t to be 1: 8. In each category of artistic font, 361 775 Chinese characters and 22 uppercase English letters are 362 used for training. 197 Chinese characters and 4 uppercase 363 English letters are used for testing. 364

5.3 Comparison with SOTA Methods

Quantitative comparison. We compare three non-zero-shot 366 methods, such as AGIS-Net [Gao et al., 2019], FET-GAN [Li 367 et al., 2020a], and StarGANv2 [Choi et al., 2020]. The style 368 (glyph-effect) compositions of the target artistic fonts synthe-369 sized by them are seen in the training. Meanwhile, we also 370 compare two zero-shot methods, such as GZS-Net [Ge et al., 371 2021] and DSE-Net [Li et al., 2022a]. The style composi-372 tions they synthesized are unseen during training. In Table 1, 373



Figure 4: Comparison with state-of-the-art methods. Manual results by human are shown in the last column as ground truth. Six rows of experimental results correspond to (1) Chinese artistic font with normal glyph. (2) Creative glyph. (3) Handwriting glyph. (4) Calligraphy glyph. (5) English artistic font with simple effect. (6) English artistic font with delicate effect.



Figure 5: Ablation study of CAFS-GAN. The Baseline includes three encoders and a generator without two style contrastive losses and SSA, and it receives two style vectors that have been spliced in the basic AdaIN layer. The setup of 5 groups of experiments: (A) adding $\mathcal{L}_{sty}^{E_x}$ to the Baseline, (B) incrementally adding $\mathcal{L}_{sty}^{E_y}$, (C) incrementally adding SSA, (D) replacing AdaIN with a reverse version of hierarchical dual styles AdaIN based on (C). (E) incrementally adding hierarchical dual styles AdaIN based on (C). The setup of experiment (E) denotes the full of CAFS-GAN.

the CAFS-GAN proposed by us has achieved apparent advantages in synthesizing unseen style compositions. Moreover,
the synthesized results by CAFS-GAN are also ahead of the

377 conventional artistic font synthesis methods in five metrics.

Qualitative comparison. In Figure 4, our method has gen-378 erated photo-realistic glyph and effect style and is superior 379 to other methods. We can easily observe that some meth-380 ods work well in the normal glyph, but their performance in 381 creative, handwriting, and calligraphy drops sharply. For En-382 glish, some methods are difficult to generate the correct glyph 383 and effect (e.g., DSE-Net), and the othes are difficult to gen-384 erate the correct character content (e.g., GZS-Net). 385

386 5.4 Ablation Study

We conducted ablation study to validate the effectiveness of the components and loss functions of the model. The experimental results are depicted in Figure 5 and Table 2.

Style contrastive losses. The purpose of style contrastive
 losses is to disentangle the glyph and effect and improve the
 encoder's ability to extract pure glyph and effect features. In

Figure 5(A), after we add $\mathcal{L}_{sty}^{E_x}$, the dark red effect disappears 393 obviously and the correct metal texture effect appears. After 394 we simultaneously add $\mathcal{L}_{sty}^{E_x}$ and $\mathcal{L}_{sty}^{E_y}$, the glyph structure of 395 (B) becomes more accurate than (A). 396

Style similarity attention. The SSA makes use of the style 397 similarity between the positive and original samples to en-398 hance the feature signal of the glyph and effect. We add SSA 399 to the setup of experiment (B). In Figure 5(C), the stroke on 400 the left side of this character has been significantly improved. 401 Hierarchical dual styles AdaIN. This structure helps the 402 model to synthesize artistic fonts from structure to texture 403 through hierarchically input to improve image details. The 404 reverse version of this structure treats the glyph as low-405 dimensional information and the effect as high-dimensional 406 information. We add the reverse version of hierarchical dual 407 styles AdaIN to the setup of experiment (C). Figure 5 (C)(D) 408 shows that the reverse version will lose a lot of effects and 409 glyph details. Then, we add the right version of hierarchical 410 dual styles AdaIN to the setup of experiment (C). Figure 5 411 (C)(E) shows the optimization of image details. 412

	$L_1 \log \downarrow$	$FID\downarrow$	SSIM \uparrow	$\text{GOLM} \downarrow$	$EPE\downarrow$
Baseline	0.2750	261.08	0.3039	189.29	2.6751
(A)	0.2852	257.73	0.2653	187.51	3.9582
(B)	0.2336	201.66	0.3345	183.83	2.0330
(C)	0.2290	178.07	0.3333	182.81	1.2076
(D)	0.2452	262.31	0.3099	185.10	1.6954
(E)	0.2251	179.61	0.3520	179.25	1.0767

Table 2: Quantitative evaluation of ablation study.



Figure 6: Influence of the proportion of seen style compositions. The x-coordinate represents the proportion of C_t to C, and the y-coordinate represents the value of each metric.

413 **5.5** Proportion of the Seen Style Compositions

We also discussed the influence of the proportion of seen style 414 composition C_t to all possible style compositions C on the ex-415 perimental results. We use six different training sets to train 416 CAFS-GAN, each containing the same three effects and three 417 glyphs, but their number of compositions is different. The 418 ratios of style combinations of C_t to C are set to 4/9, 5/9, 419 6/9, 7/9, 8/9, and 9/9. As shown in Figure 6, with the pro-420 portion increase, the model's performance presents an over-421 all improved state. Therefore, we concluded that sufficient 422 glyph-effect joint compatibility relationship will improve the 423 model's ability to understand the artistic font's attributes and 424 425 help the model synthesize unseen style compositions.

426 5.6 Visualization

In order to further demonstrate the style disentanglement ca-427 pability of the E_x and E_y and the ability to recombine content 428 and styles of the generator, we visualize the attention maps 429 generated by style encoders and feature maps generated by 430 the generator. In Figure 7(a), we feed three different effects 431 of the artistic font images to 1 and 1, and the texture part of 432 these images got a lot of attention. In Figure 7(b), the glyph 433 encoder tends to focus on local areas of artistic fonts, which 434 are the unique characteristics of the glyphs, such as curves 435 and corners. In Figure 7(c), the structure of feature maps of 436 fonts are changed firstly (e.g., the lines become clear, and the 437 corners become apparent). Then, there is more pixel filling 438 inside the feature maps of the font. After that, the texture is 439 rendered. 440



Figure 7: The visualization of the style attention maps and generated feature maps. (a)(b) We show the original images (in the first row) and their attention maps of glyph and effect (in the second row) in two style encoders. (c) We show how the generator adjusts the structure and then renders the effect.



Figure 8: Glyph style interpolation and effect style interpolation.

5.7 Style Interpolation

We further demonstrate the flexibility of CAFS-GAN through glyph style interpolation and effect style interpolation. In CAFS-GAN, we can explicitly control the weighting between different glyph or effect representations and decode the integrated representation back to the image space, obtaining the new mixed attributes, see Figure 8. This is meaningful to the diversification of artistic fonts. 442

441

449

5.8 Conclusion

In this paper, we propose a new task called compositional 450 zero-shot artistic font synthesis (CAFS), which allows syn-451 thesizing arbitrary character's artistic font image with un-452 seen style compositions. To achieve this task, we propose 453 the CAFS-GAN model, focusing on style disentanglement of 454 glyph and effect, and hierarchical reorganization of content 455 and styles representations. We also propose two evaluation 456 metrics for a more comprehensive evaluation of artistic font 457 images: glyph outline misalignment and effect perception 458 error. Extensive experiments demonstrate the effectiveness 459 of our model's multi-attributes control and the superiority of 460 generation quality over existing methods. 461

462 **References**

- ⁴⁶³ [Azadi *et al.*, 2018] Samaneh Azadi, Matthew Fisher,
 ⁴⁶⁴ Vladimir G Kim, Zhaowen Wang, Eli Shechtman, and
 ⁴⁶⁵ Trevor Darrell. Multi-content gan for few-shot font style
 ⁴⁶⁶ transfer. In *CVPR*, pages 7564–7573, 2018.
- 467 [Chen et al., 2018] Ricky TQ Chen, Xuechen Li, Roger
- Grosse, and David Duvenaud. Isolating sources of disentanglement in vaes. In *NIPS*, pages 2615–2625, 2018.
- 470 [Choi et al., 2020] Yunjey Choi, Youngjung Uh, Jaejun Yoo,
- and Jung-Woo Ha. Stargan v2: Diverse image synthesis
 for multiple domains. In *CVPR*, pages 8188–8197, 2020.

473 [Feng et al., 2022] Yaogong Feng, Xiaowen Huang, Pengbo

- Yang, Jian Yu, and Jitao Sang. Non-generative generalized
 zero-shot learning via task-correlated disentanglement and
 controllable samples synthesis. In *CVPR*, pages 9346–
 9355, 2022.
- [Gao *et al.*, 2018] Rui Gao, Xingsong Hou, Jie Qin, Li Liu,
 Fan Zhu, and Zhao Zhang. A joint generative model for
 zero-shot learning. In *ECCV*, pages 631–646, 2018.
- [Gao *et al.*, 2019] Yue Gao, Yuan Guo, Zhouhui Lian, Ying min Tang, and Jianguo Xiao. Artistic glyph image synthe sis via one-stage few-shot learning. *ACM Transactions on*
- Graphics (TOG), 38(6):1–12, 2019.
 [Gatys *et al.*, 2016] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional
- neural networks. In *CVPR*, pages 2414–2423, 2016.

⁴⁸⁸ [Ge *et al.*, 2021] Yunhao Ge, Sami Abu-El-Haija, Gan Xin,
⁴⁸⁹ and Laurent Itti. Zero-shot synthesis with group⁴⁹⁰ supervised learning. In *ICLR*, 2021.

⁴⁹¹ [Goodfellow *et al.*, 2014] Ian Goodfellow, Jean Pouget⁴⁹² Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley,
⁴⁹³ Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Gen⁴⁹⁴ erative adversarial networks. In *NIPS*, page 2672–2680,
⁴⁹⁵ 2014.

- ⁴⁹⁶ [Han *et al.*, 2021] Yuxuan Han, Jiaolong Yang, and Ying Fu.
 ⁴⁹⁷ Disentangled face attribute editing via instance-aware la⁴⁹⁸ tent space search. In *IJCAI*, pages 715–721, 2021.
- ⁴⁹⁹ [He *et al.*, 2020] Kaiming He, Haoqi Fan, Yuxin Wu, Sain⁵⁰⁰ ing Xie, and Ross Girshick. Momentum contrast for un⁵⁰¹ supervised visual representation learning. In *CVPR*, pages
 ⁵⁰² 9729–9738, 2020.
- [Higgins *et al.*, 2017] Irina Higgins, Loic Matthey, Arka Pal,
 Christopher Burgess, Xavier Glorot, Matthew Botvinick,
- 505 Shakir Mohamed, and Alexander Lerchner. beta-vae:
- Learning basic visual concepts with a constrained varia-
- tional framework. In *ICLR*, 2017.
- [Hong *et al.*, 2022] Ziming Hong, Shiming Chen, Guo-Sen
 Xie, Wenhan Yang, Jian Zhao, Yuanjie Shao, Qinmu Peng,
- and Xinge You. Semantic compression embedding for
- generative zero-shot learning. In *IJCAI*, pages 956–963,
- 512 2022.
- ⁵¹³ [Huang and Belongie, 2017] Xun Huang and Serge Be ⁵¹⁴ longie. Arbitrary style transfer in real-time with adaptive
 ⁵¹⁵ instance normalization. In *ICCV*, pages 1501–1510, 2017.

- [Karthik *et al.*, 2022] Shyamgopal Karthik, Massimiliano
 Mancini, and Zeynep Akata. Kg-sp: Knowledge guided
 simple primitives for open world compositional zero-shot
 learning. In *CVPR*, pages 9336–9345, 2022.
- [Li et al., 2020a] Wei Li, Yongxing He, Yanwei Qi, Zejian520Li, and Yongchuan Tang. Fet-gan: Font and effect transfer521via k-shot adaptive instance normalization. In AAAI, pages5221717–1724, 2020.523
- [Li *et al.*, 2020b] Yuheng Li, Krishna Kumar Singh, Utkarsh Ojha, and Yong Jae Lee. Mixnmatch: Multifactor disentanglement and encoding for conditional image generation. In *CVPR*, pages 8039–8048, 2020. 527
- [Li et al., 2022a] Xiang Li, Lei Wu, Xu Chen, Lei Meng, and Xiangxu Meng. Dse-net: Artistic font image synthesis via disentangled style encoding. In *ICME*, pages 1–6, 2022.
- [Li *et al.*, 2022b] Xiangyu Li, Xu Yang, Kun Wei, Cheng
 Deng, and Muli Yang. Siamese contrastive embedding
 network for compositional zero-shot learning. In *CVPR*,
 pages 9326–9335, 2022.
- [Lin et al., 2022] Chung-Ching Lin, Kevin Lin, Lijuan
 Wang, Zicheng Liu, and Linjie Li. Cross-modal representation learning for zero-shot action recognition. In CVPR, pages 19978–19988, 2022.
- [Luo *et al.*, 2022] Canjie Luo, Lianwen Jin, and Jingdong Chen. Siman: Exploring self-supervised representation learning of scene text via similarity-aware normalization. In *CVPR*, pages 1039–1048, 2022. 542
- [Mancini *et al.*, 2021] Massimiliano Mancini, Muhammad Ferjad Naeem, Yongqin Xian, and Zeynep Akata. Open world compositional zero-shot learning. In *CVPR*, pages 5222–5230, 2021. 546
- [Men *et al.*, 2019] Yifang Men, Zhouhui Lian, Yingmin 547 Tang, and Jianguo Xiao. Dyntypo: Example-based dynamic text effects transfer. In *CVPR*, pages 5870–5879, 549 2019. 550
- [Naeem *et al.*, 2021] Muhammad Ferjad Naeem, Yongqin Xian, Federico Tombari, and Zeynep Akata. Learning graph embeddings for compositional zero-shot learning. In *CVPR*, pages 953–962, 2021. 554
- [Saini et al., 2022] Nirat Saini, Khoi Pham, and Abhinav Shrivastava. Disentangling visual embeddings for attributes and objects. In CVPR, pages 13658–13667, 2022. 557
- [Simonyan and Zisserman, 2014] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv:1409.1556*, 2014. 560
- [Singh and Krishnan, 2020] Saurabh Singh and Shankar Krishnan. Filter response normalization layer: Eliminating batch dependence in the training of deep neural networks. In *CVPR*, pages 11237–11246, 2020. 564
- [Vaswani et al., 2017] Ashish Vaswani, Noam Shazeer, Niki
 Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
 Łukasz Kaiser, and Illia Polosukhin. Attention is all you
 need. NIPS, pages 5998–6008, 2017.

- 569 [Yang et al., 2016] Shuai Yang, Jiaying Liu, Zhouhui Lian,
- and Zongming Guo. Awesome typography: Statisticsbased text effects transfer. *CoRR*, abs/1611.09026, 2016.
- [Yang *et al.*, 2019a] Shuai Yang, Jiaying Liu, Wenjing
 Wang, and Zongming Guo. Tet-gan: Text effects transfer
 via stylization and destylization. In *AAAI*, pages 1238–
 1245, 2019.
- [Yang *et al.*, 2019b] Shuai Yang, Zhangyang Wang,
 Zhaowen Wang, Ning Xu, Jiaying Liu, and Zongming Guo. Controllable artistic text style transfer via
 shape-matching gan. In *CVPR*, pages 4442–4451, 2019.
- [Yang *et al.*, 2021] Mengyue Yang, Furui Liu, Zhitang Chen,
- Xinwei Shen, Jianye Hao, and Jun Wang. Causalvae:
 Disentangled representation learning via neural structural
 causal models. In *CVPR*, pages 9593–9602, 2021.
- [Yang *et al.*, 2022] Muli Yang, Chenghao Xu, Aming Wu,
 and Cheng Deng. A decomposable causal view of compositional zero-shot learning. *IEEE Transactions on Mul- timedia*, pages 1–11, 2022.
- 588 [Zhang et al., 2018] Yuting Zhang, Yijie Guo, Yixin Jin, Yi-
- jun Luo, Zhiyuan He, and Honglak Lee. Unsupervised dis-
- 590 covery of object landmarks as structural representations.
- ⁵⁹¹ In *CVPR*, pages 2694–2703, 2018.