

Few-Shot Font Generation by Learning Fine-Grained Local Styles (Supplementary Materials)

Licheng Tang^{1*} Yiyang Cai^{2*} Jiaming Liu^{1†} Zhibin Hong¹ Mingming Gong³
Minhu Fan¹ Junyu Han¹ Jingtuo Liu¹ Errui Ding¹ Jingdong Wang¹
¹Baidu Inc. ²University of California, Berkeley ³University of Melbourne

{tanglicheng, liujiaming03, hongzhibin, fanminhu, liujingtuo, dingerrui, wangjingdong}@baidu.com
frank_cai@berkeley.edu, mingming.gong@unimelb.edu.au

A. Implementation Details

A.1. Reference Selection

In this section, we will discuss the implementation details of reference selection, which can be divided into two steps: **reference set determination** and **content-reference mapping**.

Reference set determination As mentioned in paper’s section 3.4, our target is to search for a reference glyph set which covers the most conspicuous-level components, and our detailed algorithm can be demonstrated by Algorithm 1. In our algorithm, X denotes the entire character list sorted by glyphs’ occurrence frequency, whose total capacity is 20K. T denotes the single level decomposition look-up tables for all characters in X . For one key x_i in T , their corresponding value χ_i is a subset of X , consisting of their decomposing components’ name list and corresponding decomposing form. By exploring this components list, we employ breadth first search to recursively obtain all the components at different levels, which forming the component tree in our work. C denotes the set that contains all conspicuous-level components, whose total capacity is 374. N_{ref} denotes the total capacity of our target reference set, which is fixed to 100 in our experiments.

As an initialization, we set our target reference set \hat{U} and corresponding covering components \hat{C} as \emptyset . Then we begin searching process in X . For every glyph x_i in X , we obtain its conspicuous-level components c_i via the function *searchComponents* in 1. If there is unique components in c_i which is not in current component pool \hat{C} , we regard this glyph x_i as our target reference glyph, and we add it and its unique components to \hat{U} and \hat{C} respectively. In our experiment, we ensure that every latest adding reference have 2 or more unique components. Once the capacity of \hat{U} reaches N_{ref} , we terminate the searching process.

Content-reference mapping Since we have determined

*Equal contribution.

†Corresponding author.

Algorithm 1: Reference set selection

input : $X = \{x_i\}$: Common-used glyph list.
 $T = \{x_i : \chi_i\}$: Single-level decomposition look-up table. $\chi_i \subset X$.
 C : Conspicuous-level component set.
 N_{ref} : Max capacity of Reference set.
 d : Max search depth in component tree.
output: \hat{U} : Reference set.
 \hat{C} : All components which U can cover.

Function *searchComponents* (x_i):

```
queue  $\leftarrow$   $\{x_i\}$ ;  $i \leftarrow 0$ ;  $c_i \leftarrow \emptyset$ 
repeat
  queue*  $\leftarrow$   $\emptyset$ 
  for  $x \in$  queue do
    queue*  $\leftarrow$  queue*  $\cup$   $T(x)$ 
     $c_i \leftarrow T(x) \cap C$ 
  end
  queue  $\leftarrow$  queue*;  $i \leftarrow i + 1$ 
until  $i \geq d$ ;
return  $c_i$ 
```

End Function

Function *Main*():

```
 $\hat{U} \leftarrow \emptyset$ ;  $\hat{C} \leftarrow \emptyset$ 
for  $x_i \in X$  do
   $c_i \leftarrow$  searchComponents( $x_i$ ).
  if  $\exists \hat{c} \in c_i$  s.t.  $\hat{c} \notin \hat{C}$  then  $\hat{C} \leftarrow \hat{C} \cup c_i$ ,
  add  $x_i$  to  $\hat{U}$ ;
  if  $Len(\hat{U}) \geq N_{ref}$  then break;
end
return  $\hat{U}, \hat{C}$ 
```

End Function

the reference set. It’s more convenient for us to find the mapping relationship between a arbitrary glyph and its ref-

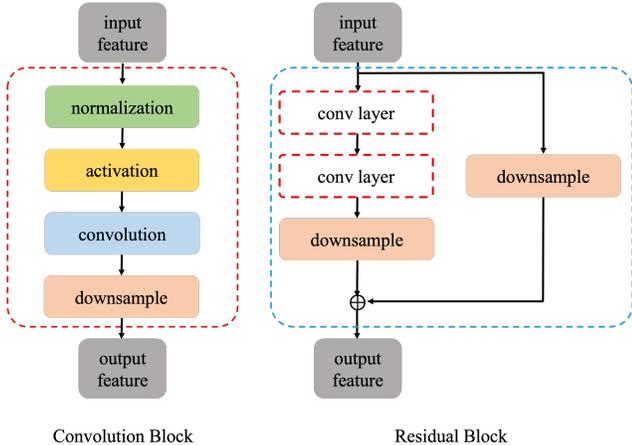


Figure 1. Structures of the Convolution Block and the Residual Block in Tab.1

erences whose conspicuous components are in common. For k -shot reference mapping (one content glyph corresponds to k style references), we do the following searching process for k times: We find the glyph in reference set which share the most conspicuous component, and remove it from the original reference set. If there is multiple choice, we leave the one whose component’s composition form is the same with the content glyph.

We show some of the finding mappings between content and references in Fig. 2.

A.2. Model Architecture

Our entire model can be divided into two parts: the generator and discriminator. Both of them are built up of two typical types of blocks: the Convolution block and the Residual block, which are illustrated in Fig 1. The Residual block contains two identical Convolution blocks. Since the Residual block has its own downsampling operator, its Convolution blocks skip the downsampling step.

Generator As mentioned in our method, the generator model consists of three respective modules: reference encoder E_r , content encoder E_c , and decoder D . Both E_r and E_c are also made up of the following two types of blocks: the Convolution block and Residual block. The detailed architecture is illustrated in Table. 1. In D , every Convolution block is followed by a spectral normalization.

Discriminator The discriminator is also made up of the Convolution block and the Residual block. The detailed architecture is illustrated in Table. 1.

A.3. Training details

We use Adam optimizer to optimize the FSFont’s parameters. The generator and discriminator’s learning rates are 0.0002 and 0.0008, respectively. Kaiming initialization is

applied for the model. During training, by default we set the reference to be 3-shot for training. If one glyph’s reference set’s capacity is less than 3, we duplicate one of reference glyph until the capacity equals to 3 for purpose of batch training. In our proposed SAM, empirically we set number of attention heads to be 8 and batch size to be 32. We train our model with full objective function for 500k iterations.

B. Additional Experimental Results

In Fig.2, we demonstrate more experimental results on unseen fonts. The results show that FSFont can deal with variant font styles, including typewriter fonts, artistic fonts, and handwriting fonts. In addition, we show the content font we used in our experiment in Fig.3.

We also demonstrate more visualization results about the attention maps in Fig. 4. We show the attention maps on characters with variant structures. It is shown that the attention is able to capture the correct correspondence in different structures.

C. User Study Examples

We show the sample images used for the user study in Fig.5 and Fig.6. A content image, a reference set, and shuffled results from five methods, i.e. FUNIT [2], AGIS-Net [1], LF-Font [3], DFFont [5], MX-Font [4], and FSFont, are displayed to users for every query. Users are required to choose the most consistent one according the style of the reference set. As the orders of the methods are shuffled, we also provide the answers for each case.

References

- [1] Jie Chang, Yujun Gu, Ya Zhang, Yan-Feng Wang, and CM Innovation. Chinese handwriting imitation with hierarchical generative adversarial network. In *BMVC*, page 290, 2018. 2
- [2] Ming-Yu Liu, Xun Huang, Arun Mallya, Tero Karras, Timo Aila, Jaakko Lehtinen, and Jan Kautz. Few-shot unsupervised image-to-image translation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10551–10560, 2019. 2
- [3] Song Park, Sanghyuk Chun, Junbum Cha, Bado Lee, and Hyunjung Shim. Few-shot font generation with localized style representations and factorization. *arXiv preprint arxiv:2009.11042*, 2020. 2
- [4] Song Park, Sanghyuk Chun, Junbum Cha, Bado Lee, and Hyunjung Shim. Multiple heads are better than one: Few-shot font generation with multiple localized experts. *arXiv preprint arxiv:2104.00887*, 2021. 2
- [5] Yangchen Xie, Xinyuan Chen, Li Sun, and Yue Lu. Dg-font: Deformable generative networks for unsupervised font generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5130–5140, 2021. 2

Reference Encoder E_r							
Layer Type	Normalization	Activation	Padding	Kernel Size	Stride	Downsample	Feature maps
Convolution block	IN	ReLU	1	3	1	-	32
Convolution block	IN	ReLU	1	3	1	AvgPool	64
Convolution block	IN	ReLU	1	3	1	AvgPool	128
Residual block	IN	ReLU	1	3	1	-	128
Residual block	IN	ReLU	1	3	1	-	128
Residual block	IN	ReLU	1	3	1	AvgPool	256
Residual block	IN	ReLU	1	3	1	-	256
Output Layer	-	Sigmoid	-	-	-	-	256

Content Encoder E_c							
Layer Type	Normalization	Activation	Padding	Kernel Size	Stride	Downsample	Feature maps
Convolution block	IN	ReLU	1	3	1	-	32
Convolution block	IN	ReLU	1	3	2	-	64
Convolution block	IN	ReLU	1	3	2	-	128
Convolution block	IN	ReLU	1	3	2	-	256
Convolution block	IN	ReLU	1	3	1	-	256

Decoder D							
Layer Type	Normalization	Activation	Padding	Kernel Size	Stride	Upsample	Feature maps
Residual block	IN	ReLU	1	3	1	-	256
Residual block	IN	ReLU	1	3	1	-	256
Residual block	IN	ReLU	1	3	1	-	256
Convolution block	IN	ReLU	1	3	1	Nearest	128
Convolution block	IN	ReLU	1	3	1	Nearest	64
Convolution block	IN	ReLU	1	3	1	Nearest	32
Convolution block	IN	ReLU	1	3	1	-	1
Output Layer	-	Sigmoid	-	-	-	-	1

Discriminator D							
Layer Type	Normalization	Activation	Padding	Kernel Size	Stride	Downsample	Feature maps
Convolution block	IN	-	1	3	2	-	32
Residual block	IN	ReLU	1	3	1	AvgPool	64
Residual block	IN	ReLU	1	3	1	AvgPool	128
Residual block	IN	ReLU	1	3	1	AvgPool	256
Residual block	IN	ReLU	1	3	1	-	256
Residual block	IN	ReLU	1	3	1	AdaAvgPool	512
Output layer	Multi-task embedding						

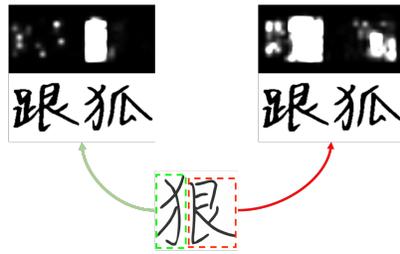
Table 1. **Architecture of the generator modules E_r , E_c , D and the discriminator.** Convolution block and Residual block denote the module mentioned in Fig 1. IN denotes instance normalization. All padding operation is zero-padding. AvgPool and AdaptiveAvgPool denotes average pooling and adaptive average pooling. In discriminator’s output layer, we use two embedding operators to embed the output feature map into two prediction vector of the font style and the character’s name.

Content	Reference	Generate	GT	Content	Reference	Generate	GT
剧	呢别卖	剧	剧	剧	呢别卖	剧	剧
咪	说和	咪	咪	忘	就想甩	忘	忘
连	转还	连	连	徂	组得	徂	徂
约	组的	约	约	意	是想亲	意	意
焦	焦你	焦	焦	弦	就张累	弦	弦
媪	好去累	媪	媪	描	把累菜	描	描
底	的度舐	底	底	绀	组疝	绀	绀
完	下家祝	完	完	轻	没转差	轻	轻
估	你喻装	估	估	叠	没带组	叠	叠
昙	是去那	昙	昙	褪	还跟被	褪	褪
备	峰累	备	备	歌	啊次灯	歌	歌
悖	累快请	悖	悖	静	请你建	静	静
餌	鲫饭	餌	餌	药	的菜组	药	药
谓	请累	谓	谓	逛	还弄狐	逛	逛
洳	啊好没	洳	洳	魏	和好魁	魏	魏
锱	啊亲钟	锱	锱	炸	酢灯	炸	炸
猓	岁罪狐	猓	猓	砍	次码	砍	砍
超	起啊那	超	超	保	你和	保	保
房	咋下放	房	房	杰	和焦	杰	杰
鞋	去鞭	鞋	鞋	狗	的啊狐	狗	狗

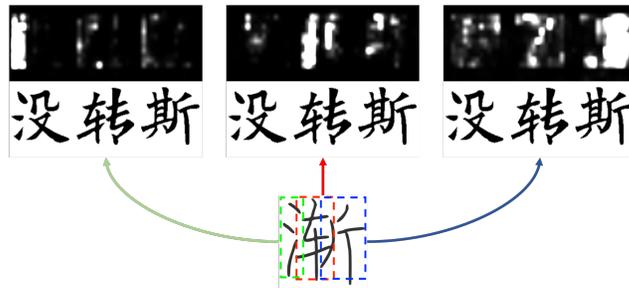
Figure 2. More Results. This Figure shows the Content-Reference mapping and generated results on all test fonts.

图抢娅钛伫鳞赅萦阎么骠纭仓计缙缝傲动迟脍
 垒项诤掩鹄焕鹄灏屈乐飙匭垆铈均邝噤层镌驪
 鸚变矧尘躐项饺罚谤丝荡敌鹅僚莠摺颇晓鲟谀
 脓郟矸恢窍滦册嵘貳见鵞议铈掳缙缙维钦趺榉
 驷圻噉鸱辽胼枉轳盐苜苒诔苇鸢啮蛲滚汧鳖镍铸
 绉熠头狎浑带喜匀黠洼嫫毳钹农掷涛鲑臬拣淮轩饼瞩
 桥腾减陕钩奋戆辄颁駉紧农掷涛鲑臬拣淮轩饼瞩
 焜骤鸯纒专论吓葜务佻鹏鲨呕冥赏癃签潇呛赢砗
 问纹咻炅栌蒞态碇暮馗鸬瓢临谜蔚荆渗媛产饩
 卜猫执鸣巅将盗帜钼隍簧椽酝预鰈浅菘鏢导馁
 犹员酱掬脏驻寝汤拔鹜钙蹕搥渍祸阔坠许涡戈
 吟蛄诚鲛恼惭启萆样说麸抡溲躡戩详蚩癞尧录
 贿帅剪鰌碎滤叹兽猕鳧嚶断诩玃缓讶媪赖锰捞

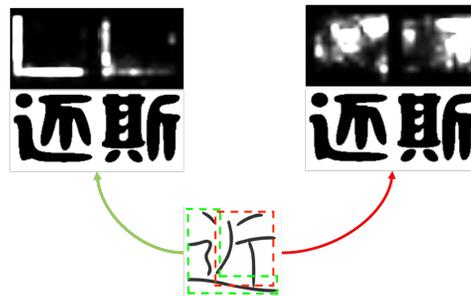
Figure 3. The content font we used in our experiment.



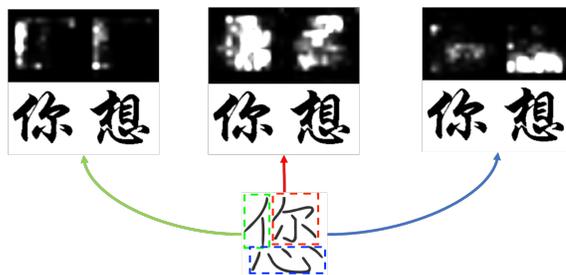
(a) Left-Right structure



(b) Left-Middle-Right structure



(c) Half-Surrounded Structure



(d) Compound Structure (Left-Right + Top-down)

Figure 4. **Visualization on Attention maps.** This Figure shows the attention maps given the different regions in content glyph.

content	强	纸	床	怎	催	崔	吹	客	坑	话	冯	连	法	拍	伟
reference	啊	舐	和	想	你	岁	喻	啊	就	和	次	远	没	的	你
	张	组	度	醉	岁	焦	冷	家	凯	说	吗	转	去	是	韩
	蛮				焦			峰	去	卖				把	
A	强	纸	床	怎	催	崔	吹	客	坑	话	冯	连	法	拍	伟
B	强	纸	床	怎	催	崔	吹	客	坑	话	冯	连	法	拍	伟
C	强	纸	床	怎	催	崔	吹	客	坑	话	冯	连	法	拍	伟
D	强	纸	床	怎	催	崔	吹	客	坑	话	冯	连	法	拍	伟
E	强	纸	床	怎	催	崔	吹	客	坑	话	冯	连	法	拍	伟
F	强	纸	床	怎	催	崔	吹	客	坑	话	冯	连	法	拍	伟
answer	①	③	④	⑥	②	⑤	④	①	②	⑤	⑥	④	②	①	④
	②	④	①	②	⑤	⑥	③	②	⑤	⑥	①	③	③	②	⑤
	③	⑤	⑥	⑤	④	①	⑥	⑤	⑥	④	④	⑥	④	⑥	③
	④	⑥	②	①	③	④	①	⑥	③	①	③	⑤	①	③	②
	⑤	②	⑤	③	①	③	②	③	④	②	⑤	②	⑤	④	⑥
	⑥	①	③	④	⑥	②	⑤	④	①	③	②	①	⑥	⑤	①

①: AGIS-Net ②: FUNIT ③: DG-Font ④: LF-Font ⑤: MX-Font ⑥: Our FS-Font

Figure 5. User studies on UFSC (unseen font seen character) data. This Figure shows how do we conduct the user study.

content	钜	挫	焦	馁	徉	炜	湮	鲤	狸	滋	揉	扛	怯	姣	迂
reference	去	去	焦	好	得	韩	没	去	去	没	把	差	去	皎	还
	钟	个	你	饭	差	灯	甄	累	累	累	鞣	把	快	好	和
		把		菜				酸	狐	噗					菜
A	钜	挫	焦	馁	徉	炜	湮	鲤	狸	滋	揉	扛	怯	姣	迂
B	钜	挫	焦	馁	徉	炜	湮	鲤	狸	滋	揉	扛	怯	姣	迂
C	钜	挫	焦	馁	徉	炜	湮	鲤	狸	滋	揉	扛	怯	姣	迂
D	钜	挫	焦	馁	徉	炜	湮	鲤	狸	滋	揉	扛	怯	姣	迂
E	钜	挫	焦	馁	徉	炜	湮	鲤	狸	滋	揉	扛	怯	姣	迂
F	钜	挫	焦	馁	徉	炜	湮	鲤	狸	滋	揉	扛	怯	姣	迂
answer	③	⑥	③	④	②	①	④	②	⑥	④	②	⑤	①	⑥	⑤
	②	④	④	⑤	⑤	②	③	③	①	③	⑤	⑥	②	②	⑥
	⑥	①	⑤	③	④	⑥	⑥	④	④	⑥	⑥	④	③	⑤	①
	①	②	⑥	②	③	③	①	①	③	⑤	③	①	⑤	①	④
	④	⑤	②	⑥	①	④	②	⑤	⑤	②	④	②	⑥	③	③
	⑤	③	①	①	⑥	⑤	⑤	⑥	②	①	①	③	④	④	②

①: AGIS-Net ②: FUNIT ③: DG-Font ④: LF-Font ⑤: MX-Font ⑥: Our FS-Font

Figure 6. User studies on UFUC (unseen font unseen character) data. This Figure shows how do we conduct the user study.