

# Typeface Completion with Generative Adversarial Networks

Yonggyu Park, Junhyun Lee, Yookyung Koh, Inyeop Lee, Jinhyuk Lee, Jaewoo Kang

Department of Computer Science and Engineering

Korea University

{yonggyu, ljhyun33, ykko603, blackin77, jinhyuk\_lee, kangj}@korea.ac.kr

## Abstract

The mood of a text and the intention of the writer can be reflected in the typeface. However, in designing a typeface, it is difficult to keep the style of various characters consistent, especially for languages with lots of morphological variations such as Chinese. In this paper, we propose a Typeface Completion Network (TCN) which takes one character as an input, and automatically completes the entire set of characters in the same style as the input characters. Unlike existing models proposed for image-to-image translation, TCN embeds a character image into two separate vectors representing typeface and content. Combined with a reconstruction loss from the latent space, and with other various losses, TCN overcomes the inherent difficulty in designing a typeface. Also, compared to previous image-to-image translation models, TCN generates high quality character images of the same typeface with a much smaller number of model parameters.

We validate our proposed model on the Chinese and English character datasets, which is paired data, and the CelebA dataset, which is unpaired data. In these datasets, TCN outperforms recently proposed state-of-the-art models for image-to-image translation. The source code of our model is available at <https://github.com/yonggyu/TCN>.

## Introduction

Typeface is a set of one or more fonts, each consisting of glyphs that share common design features.<sup>1</sup> Effective typeface not only allows writers to better express their opinions, but also helps convey the emotions and moods of their text. However, there is a small number of typefaces to choose from because there are several difficulties in designing typography. The typeface of all characters should be the same without compromising readability. As a result, it takes much effort to make a typeface for languages with a large character set such as Chinese which contains more than twenty thousand characters.

To deal with this difficulty, we aim to build a model that takes one character images as an input, and generates all the remaining characters in the same typeface of input characters, which is illustrated in Figure 1.

<sup>1</sup><https://en.wikipedia.org/wiki/Typeface>

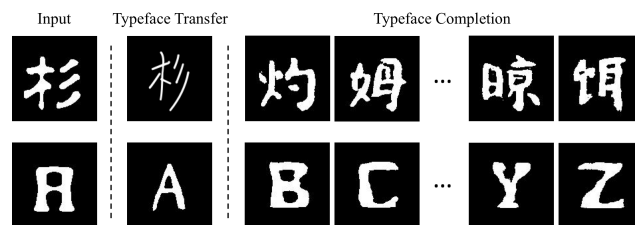


Figure 1: Comparison of the output of TCN with that of the typeface transfer model. Unlike typeface transfer model that changes the style of an input, typeface completion model generates the contents of all character sets in the same style as the input. Typeface completion outputs are the results of TCN.

In the field of computer vision, the typeface completion task has not been much studied. Generating character images in the same typeface could be seen as a image-to-image translation problem. Existing image-to-image translation tasks often refer to extracting a style feature from a desired image, and combines the style feature while keeping the content features of desire images. In case of typeface completion task, after extracting a style feature from a desired image, we combine the style feature while changing the content features of same desire images. For the typeface completion task, we use the terms style feature and typeface feature, interchangeably.

And, as existing single image-to-image translation models learn only a single domain translation (Gatys, Ecker, and Bethge 2016; Huang and Belongie 2017; Li et al. 2017; Li et al. 2018), we need to train  $N(N-1)/2$  models on a set of  $N$  characters for the typeface completion task. While learning all the single models and keeping them for typeface completion is computationally infeasible, recent work of Choi et al. (2017) has addressed this inefficient but fails to produce high quality character images for typeface completion.

In typeface completion, a large number of classes such as Chinese characters should be considered. Image-to-Image translation models designed for small number of classes fail to generalize in the typeface completion task due to the large number of classes in character sets. In the existing model, the input channel becomes  $1+n$  where 1 refers to grey-scale

channels for a character image, and  $n$  refers to the number of classes. However, concatenating one-hot encoded labels directly to the image tensor makes the model ineffective when  $n$  is large. That's because most of the input value will be zero.

In order to overcome the weakness of existing image-to-image translation models, and to deal with the large number of classes and paired dataset, we propose a Typeface Completion Network (TCN) that generates all characters in a character set from a single model. TCN represents the typefaces and contents of characters as latent vectors, and uses various losses. We show that TCN outperforms on Chinese and English datasets in terms of task-specific quantitative metric and image qualities. This is possible because the character image is paired data, unlike other general images. And for the same reason, we can take advantage to generate more plausible images with additional loss. Furthermore, it also showed better performance than the existing image-to-image translation model in general images that do not take the advantage mentioned above. This shows that TCN is applicable to unpaired dataset.

## Related Works

### Image-to-Image Translation

Generative Adversarial Networks (GAN)(Goodfellow et al. 2014) has been highlighted as one of the hottest research topics in computer vision. GAN generates images using an adversarial loss with a deep convolution architecture (Radford, Metz, and Chintala 2015; Goodfellow 2017). GAN has gained popularity and resulted in a variety of follow-up studies (Mirza and Osindero 2014; Perarnau et al. 2016; Arjovsky, Chintala, and Bottou 2017; Ledig et al. 2017; Chen and Koltun 2017) in the context of image generation, super-resolution, colorization, inpainting, attribute transfer, and image-to-image translation.

The style transfer task, one of the image-to-image translation tasks, involves changing the style of an image while retaining its content. Since most existing style transition models have a fixed pair of input and target style, they cannot receive or generate styles in various domains using a single model. However, the models of (Mirza and Osindero 2014; Choi et al. 2017), can take a target style label as an input and generate an image of the desired style using a single model. This reduces the number of parameters in a task which performs the transition into various style domains.

A task of completing a character set with some subsets can also be treated as a image-to-image task. Our model changes the content of an input to the content of the target while maintaining the style of the input. This is the same as the existing style transfer model where the terms, style and content, are reversed. In addition, since our task requires various content domains, our model is based on a multi-domain transfer model.

### Character Image Generation

Early character image generation models focused on geometric information (Tenenbaum and Freeman 1997; Suvananont and Igarashi 2010; Campbell and Kautz 2014;

Phan, Fu, and Chan 2015). But now, with the development of deep learning, many models focus on character style transfer tasks (Upchurch, Snaveley, and Bala 2016; Baluja 2016). In particular, there have been researches on the style transfer task using GAN in the character image domain (Lyu et al. 2017; Chang and Gu 2017; Azadi et al. 2018; Bhunia et al. 2018).

However, the difference between the existing typeface transfer model with character images and our typeface completion model is that our model considers the content as a style, not a typeface, and transfer it. This changes the number of target labels. Due to languages with thousands of characters such as Chinese, the existing single-domain image-to-image translation model and the multi-domain image-to-image translation model which uses the one-hot vector as a domain label deal with the parameter inefficiency problem. To solve this problem, we express the domain label as a latent vector and propose new losses accordingly.

## Task Definition

The typeface completion task involves completing the remaining characters  $X \setminus x_i$  of a character set  $X = \{x_1, x_2, \dots, x_N\}$  of a single typeface, using one of the  $N$  characters,  $x_i$ . TCN receives a triplet  $(x_i, y_i, y_k)$  as input, where character label  $y_i, y_k \in Y = \{y_1, y_2, \dots, y_N\}$  corresponds to  $x_i, x_k$ . Using the triplet, our model generates a character image  $\hat{x}_k \in \hat{X} = \{\hat{x}_1, \hat{x}_2, \dots, \hat{x}_N\}$ , corresponding to the character of  $y_k$  with the typeface of  $x_i$ . Since TCN generates one character at a time, above generating process is repeated  $N - 1$  times. The goal of this task is to obtain a model parameter  $\theta$  that minimizes the difference between  $\hat{x}_k$  and  $x_k$  while generating all the character sets. The overall formula of the task is as follows.

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \sum_{k=1, k \neq i}^N d(x_k, f(x_i, y_i, y_k, \theta)) \quad (1)$$

where  $d$  is the distance between the images. We mainly used the SSIM index to measure the distance between the images. From the above optimal  $\theta^*$ , we can obtain the  $\hat{x}_k$  that is similar to  $x_k$ .

$$x_k \approx \hat{x}_k = f(x_i, y_i, y_k, \theta^*) \quad (2)$$

## Proposed Model

Figure 2 shows the training process of TCN. TCN consists of typeface and content encoders, a discriminator, and a generator. The encoders extract desired feature vectors from an image. The two encoders each return a latent vector that combines different information. The generator, along with the two vectors above, receives character labels corresponding to the input and target. Through this process, the generator makes a target character image that has the same style as the input. The discriminator receives the generated image and determines if it is real. In this process, the discriminator returns the probability that the image corresponds to a certain typeface and character.

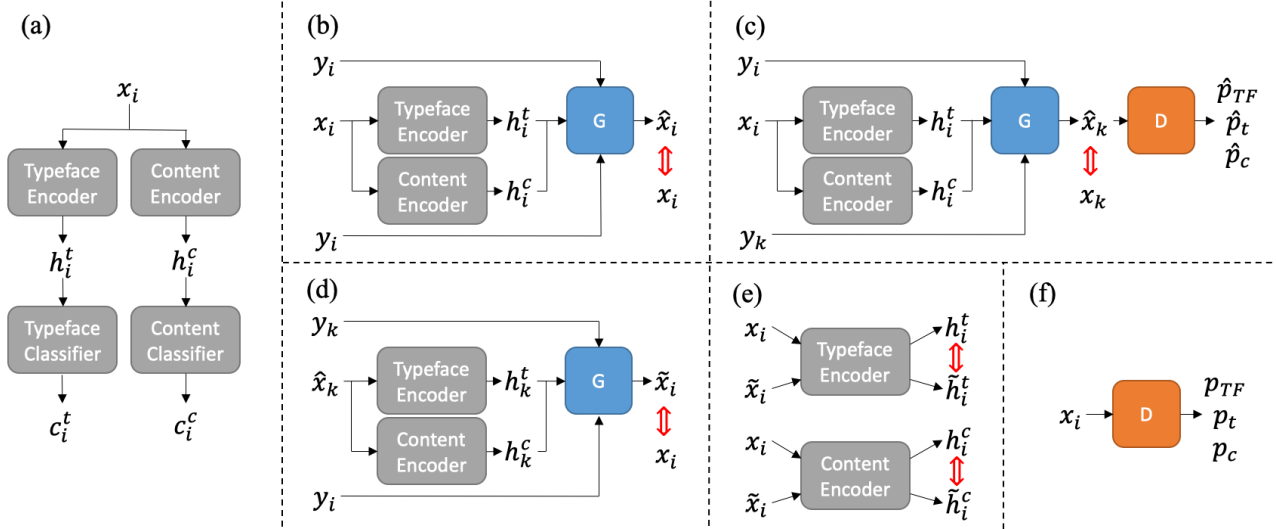


Figure 2: The overall flow chart of the TCN. The red arrow refers to the loss from the difference. Encoders,  $G$ , and  $D$  share the same weights in all the experiments. (a) Encoder pretraining. (b) Identity loss. (c) SSIM loss comparing  $\hat{x}_k$  to  $x_k$ , and adversarial loss leading to the result of the discriminator corresponding to  $x_k$ . (d) Reconstruction loss that transforms  $\hat{x}_k$  to  $\tilde{x}_i$  and compares it to the original  $x_i$ . (e) Perceptual reconstruction loss. (f) Discriminator training.

## Encoders

In our model, the encoder is divided into a content encoder and a typeface encoder. The content encoder extracts the symbolic representation of a character from an image, and the typeface encoder extracts the typographical representation from an image. Each encoder consists of a ResNet, and a 1-layer fully connected (FC) classifier, and returns the output of the ResNet and classifier, respectively.

**Typeface and Content Feature** The encoders receive the image from the main training and return the latent vectors containing the typeface and the content feature, respectively. Each expression is as follows:

$$\begin{aligned} h_i^t &= E_h^t(x_i) \\ h_i^c &= E_h^c(x_i). \end{aligned} \quad (3)$$

where the subscript  $h$  of  $E_h^t$  and  $E_h^c$  represent the ResNets of the encoders, so each  $h_i$  is an encoded latent vector of  $x_i$ .

**Encoder Pretraining** Since the two encoders have the same structure, we don't know what information each encoder extracts unless we guide them. Therefore, we pretrain the encoders so that each encoder extracts disjoint information.

For pretraining, we perform the classification task to distinguish the typeface and character of an image. The output of the classifier creates a cross entropy (CE) loss, encouraging the output of each ResNet to contain corresponding feature. As shown in Figure 2 (a), the losses from the typeface ( $L_{cls}^t$ ) and content ( $L_{cls}^c$ ) encoders are defined as

$$\begin{aligned} L_{cls}^t &= \text{CE}(c_i^t, l_i^t), c_i^t = E_c^t(h_i^t) \\ L_{cls}^c &= \text{CE}(c_i^c, l_i^c), c_i^c = E_c^c(h_i^c) \end{aligned} \quad (4)$$

where the subscript  $c$  of  $E_c^t$  and  $E_c^c$  represent the classifier of encoders, therefore each  $c_i^t$  and  $c_i^c$  is a classification result of the typeface and content of  $x_i$ , and  $l_i^t$  and  $l_i^c$  are the typeface and content labels of  $x_i$ , respectively. Especially,  $l_i^t$  is equal to  $y_i$ .

Although the classification accuracy of each encoder is higher than 80%, redundancy may exist between the two features to some extent because both typeface and content features are generated from the same character image. To solve this problem, we applied triplet loss (Schroff, Kalenichenko, and Philbin 2015) as follows.

$$\begin{aligned} L_{\text{triplet}}^t &= ||h_i^t - h_j^t|| - ||h_i^t - h_k^t|| \\ L_{\text{triplet}}^c &= ||h_i^c - h_k^c|| - ||h_i^c - h_j^c|| \end{aligned} \quad (5)$$

where  $h_i^t$  and  $h_i^c$  are ResNet results of the typeface and content, respectively. the typeface  $x_j$  is the same, the content is different from  $x_i$ .  $x_k$  has same content as  $x_i$  and has different typeface.

Last, we have obtained a reconstruction loss as in auto-encoder. The reconstruction loss ensures that no feature is lost during exclusive extraction. The Decoder used in the reconstruction task has the same structure as our generator.

## Generator

In addition to the outputs of the encoder, the generator receives input and target character labels. The generator consists of two submodules: the feature combination submodule that combines the four inputs, and the image generation submodule that generates the image using the combined inputs.

**Feature Combination** Before generating an image, we combine four inputs. The input of the generator includes the typeface/content feature vectors of the input image ( $h_i^t, h_i^c$ ),

the character label of the input ( $y_i$ ), and the target ( $y_k$ ). Ideally, the typeface encoder should extract only the typeface feature, but due to the structural nature of CNN, the typeface encoder also extracts the content feature. Accordingly, even if we extract the same typeface information from other content, we will obtain a different result. By the combination of the inputs, we want to make the feature vectors the same as the feature vectors obtained from the target image. We thus made the typeface transfer task into an auto-encoder task.

Input character labels are inserted for the multi-domain task. In the multi-domain image-to-image task with various domains of input, it is helpful for the model to know the input label with the target label rather than just the target label.

We define the combination function  $f$  by the following equation:

$$u_{i \rightarrow k} = f(h_i^t, h_i^c, y_i, y_k). \quad (6)$$

where  $f$  is a 1x1 convolutional network to establish a correlation between each channel. The 1x1 convolutional network better captures correlations than concatenating vectors and requires fewer parameters than a FC layer.

**Image Generation** Next, the generator creates the image after receiving the result of the feature combination. The image generation model is composed of deconvolutional model. We define the image generation function  $g$  by the following equation:

$$\hat{x}_k = g(u_{i \rightarrow k}) \quad (7)$$

where  $g$  can be seen as a generator of a vanilla GAN that takes a latent vector and generates an image.

We concatenate the two functions of the generator and define it as  $G$ , and it can be expressed as follows:

$$\hat{x}_k = G(h_i^t, h_i^c, y_i, y_k) = g \cdot f(h_i^t, h_i^c, y_i, y_k). \quad (8)$$

We do not distinguish between  $f$  and  $g$  in the future, but we only use  $G$ .

## Discriminator

The discriminator takes an image and determines whether the image is a real image or a fake image generated by the generator. This is used as a loss so that the image created by the generator will appear real enough to fool the discriminator.

The discriminator consists of ResNet, as in the encoders. The difference between the discriminator and the encoder is the classification part. For the encoder, there is a separate ResNet for each typeface/content classifier to distinguish the typeface and content. On the other hand, the discriminator uses one ResNet and three classifiers. One returns the T/F probability of whether the image is real, just like the discriminator of the basic GAN. The others determine which typeface and content the input has, as in (Mirza and Osindero 2014; Odena 2016; Odena, Olah, and Shlens 2016;

Perarnau et al. 2016). Our discriminator did not use a separate ResNet for each classifier and thus uses fewer parameters and normalizes losses for the three tasks. Another difference is that our discriminator does not return the output of ResNet because it is not necessary.

$$\begin{aligned} p_{TF}, p_t, p_c &= D(x_i) \\ \hat{p}_{TF}, \hat{p}_t, \hat{p}_c &= D(\hat{x}_i) \end{aligned} \quad (9)$$

We define the discriminator as  $D$  and express the results that correspond to the real image and fake image. The  $\hat{\cdot}$  denotes to be associated with a fake image by the generator.

## Training Process

**Identity Loss** Identity loss is similar to the loss in the auto-encoder in that it helps an output to be equal to an input (Fig. 2 (b)). The generator uses the character label of an input image as an input label and a target label. This experiment prevents possible loss during feature compression.

$$L_{id} = \|x_i - \hat{x}_i\|_1 \quad (10)$$

where  $\hat{x}_i$  is the image generated so that has same typeface and content with  $x_i$ .

**SSIM Loss** Structural SIMilarity index (SSIM) is used to measure the structural similarity between two images. We use SSIM index as an evaluation metric for the performance, and we also use it as a loss (Fig. 2 (c) red arrow). Using SSIM index as a loss was proposed in (Zhao et al. 2017; Snell et al. 2017). In our experiments, we applied l1-loss along with the SSIM index, in the way that showed the best performance at (Zhao et al. 2017).

$$L_{ssim} = \|x_k - \hat{x}_k\|_1 - SSIM(x_k, \hat{x}_k) \quad (11)$$

The SSIM function returns the SSIM index between two inputs. Detailed formula is in session 5.2.1.

**Adversarial Losses** Adversarial losses that help outputs to look real and deceive the discriminator are the ones that are same as those of vanilla GAN. Additionally, there is also a typeface/content classification loss between the true typeface/content label and the output that the discriminator returns. We can show these losses at the end of Figure 2 (c).

$$L_{gan} = CE(\hat{p}_{TF}, l_k^{TF}) \quad (12)$$

$$L_{cls} = CE(p_t, l_k^t) + CE(p_c, l_k^c) \quad (13)$$

In generator training, input image is always fake (Fig. 2 (c)), but since it should deceive discriminator,  $l_i^{TF}$  set to 1.

**Reconstruction Loss** Reconstruction loss, proposed by CycleGAN(Zhu et al. 2017), is a loss between the original image and the reconstruction image that is translated back to the original image from the typeface-changed image (Fig. 2 (d)). To calculate the reconstruction loss, we use the following loss function:

$$L_{rec} = \|x_i - \tilde{x}_i\|_1 \quad (14)$$

where  $\tilde{x}_i$  is the reconstruction image which has the same typeface and content as  $x_i$ .

**Perceptual Reconstruction Loss** Reconstruction loss was proposed for a pixel by pixel comparison between images, but we also apply perceptual loss to this concept. A perceptual loss was first proposed by (Johnson, Alahi, and Fei-Fei 2016) in the style transfer field. This loss compares high-dimensional semantic information in the feature vector space. Since a character image is an image composed of strokes rather than pixel units, it is appropriate to apply the perceptual loss for reconstruction image, as shown Fig. 2 (e). The equation is as follows:

$$L_{\text{per}} = ||h_i^t - \tilde{h}_i^t||_2^2 + ||h_i^c - \tilde{h}_i^c||_2^2 \quad (15)$$

where  $\tilde{h}_i^t$  and  $\tilde{h}_i^c$  is the output of the  $E_z^t$  and  $E_z^c$  for  $\hat{x}_i$ .

Perceptual reconstruction loss is the difference between the outputs of the encoder. We compare the input image and the image translated twice, not once. The typeface of the input image and that of the image translated once are the same. However, applying perceptual loss to the two images is not effective because these two images have different content features. Hence, we compare the input image and the twice-translated image with the same typeface/content as the input image.

The final loss of the generator is as follows:

$$L_g = L_{\text{gan}} + \lambda_{\text{cls}} L_{\text{cls}} + \lambda_{\text{ssim}} L_{\text{ssim}} + \lambda_{\text{rec}} (L_{\text{rec}} + L_{\text{per}} + L_{\text{id}}) \quad (16)$$

where  $\lambda_{\text{cls}}$ ,  $\lambda_{\text{ssim}}$  and  $\lambda_{\text{rec}}$  are hyper-parameters that control the importance of each loss.

**Discriminator Loss** The learning method of the discriminator is similar to that of the existing GAN. The discriminator receives two types of input: one is a real image and the other is a fake image generated by the generator. For real images, the model computes the classification loss using T/F, typeface, and content output. For fake images, only the classification loss of the T/F output is calculated because the discriminator does not need to take a loss for poor images that the generator makes.

$$L_d = \text{CE}(\hat{p}_{\text{TF}}, l_i^{\text{TF}}) + \text{CE}(p_{\text{TF}}, l_i^{\text{TF}}) + \text{CE}(p_t, l_i^t) + \text{CE}(p_c, l_i^c) \quad (17)$$

where  $l_i^{\text{TF}}$  is 1 if the input image is real and otherwise 0.  $l_i^{\text{TF}}$  will be 0 at Figure 2 (c) because the discriminator receives a fake image, and will be 1 at Figure 2 (f) because  $x_i$  is real.

## Test Process

In training, as shown in Figure 2, we induced the losses through several steps, but in the test, we carry out one step, with only encoders and generator, not using discriminator. The typeface completion task in the test is expressed as follows:

$$x_k \approx \hat{x}_k = G(h_i^t, h_i^c, y_i, y_k) \quad (18)$$

By repeating this equation  $N - 1$  times according to  $k$ , we can complete one typeface consisting of  $N$  characters.

Table 1: Dataset composition

	Chinese		English	
	#Typeface	#Image	#Typeface	#Image
<i>Train</i>	105	96,426	635	16,495
<i>Validation</i>	15	13,776	90	2,357
<i>Test</i>	30	27,637	181	4,733

## Evaluation

### Datasets

**Chinese Character** Since there are more than 50K characters in Chinese, we chose the top 1,000 most used characters<sup>2</sup>. Chinese images were collected from true-type format (TTF) and open-type format (OTF) files obtained from the Web<sup>34</sup>. A total of 150 files were manually selected. Since all files do not contain all of the 1,000 characters, we have a dataset with a total of 137,839 character images. All character image sizes are 128x128, and are gray-scale 1-channel images.

**English Character** We also build an English dataset for comparison. We used a total of 907 typographies and 26 uppercase characters. As a result of using the same selection process as the Chinese dataset, we obtained 23,583 images in total. The detailed composition is shown in Table 1.

**CelebA** We performed a style transition experiment on the CelebA(Liu et al. 2015) dataset to measure the performance of TCN. We used 202,599 images and resized them all to 128x128, as was done with the other dataset. We used three features: black, blond, brown hair colors. The data composition and other settings are the same as those of the baseline (Choi et al. 2017).

### Metrics

**SSIM** Unlike general images, a dataset of character images can be used to evaluate the output using an objective metric because character data has all the input-target pairs. We used the Structural Similarity (SSIM) index to objectively evaluate the performance on the character data. SSIM is a metric that measures the quality of images using structural information, and is defined by the following equation:

$$\text{SSIM}(a, b) = \frac{(2\mu_a\mu_b + c_1)(2\sigma_{ab} + c_2)}{(\mu_a^2 + \mu_b^2 + c_1)(\sigma_a^2 + \sigma_b^2 + c_2)} \quad (19)$$

where  $\mu_a$  is the average value of the  $a$  which denotes the brightness of the image.  $\sigma_a$  is the distribution of the  $a$  which denotes the contrast ratio of the image.  $\sigma_{ab}$  is the covariance of  $a$  and  $b$ , which denotes the correlation of the two images.  $c_1$  and  $c_2$  are small constants that prevent the denominator from being zero. The closer the score is to 1, the more similar the image is to the original image.

<sup>2</sup><http://www.qqxiuzi.cn/zh/xiandaihanyu-changyongzi.php>

<sup>3</sup><https://chinesefontdesign.com>

<sup>4</sup><http://www.sozi.cn>

Table 2: Accuracy of classifier

	Typeface	Content
<i>Chinese</i>	99.4	99.7
<i>English</i>	100.0	100.0

Table 3: Baselines

	Multi-Domain Available	Rep. of Domain Index
CycleGAN	X	None
MUNIT	X	Real-valued Distributed
StarGAN	O	One-hot
TCN	O	Real-valued Distributed

**L1 distance** In general, the  $L1$  distance is used in computer vision area (Wu, Xu, and Hall 2017). The  $L1$  distance is pixel-wise difference between the generated image and the target image. We can measure the pixel-wise performance with a intuitive and easily implemented way.

**Classification Accuracy** We introduce the classification accuracy as metric used in (Chang et al. 2018). At test dataset, We train the typeface and the content classifier, which are the ResNet. Because these classifiers show high performance accuracy over each dataset, it is reliable to be used for the metric. Each accuracy of each dataset can see at Table 2.

### Implementation Details

We selected the learning strategy and hyper-parameters of the models for the experiment.

The encoder, discriminator, and generator are all trained using the Adam optimizer, with a learning rate of 0.0001,  $\beta_1 = 0.5$ ,  $\beta_2 = 0.999$ . The learning rate gradually decreases to zero as the number of epochs is increased. The dimensions of  $h_i$  and  $u_{i \rightarrow k} \in \mathbb{R}^{256}$ .  $\lambda_{cls}$ ,  $\lambda_{sim}$  is 5,  $\lambda_{rec}$  is 10. In fact,  $\lambda_{texttcls}$  can be assigned differently for typeface and content. This ratio is a coefficient of trade-off between the typeface accuracy and the content accuracy of the resulting image. This ratio acts as a trad-off coefficient between typeface accuracy and content accuracy. The source code, implemented with Pytorch(Paszke et al. 2017), is also available at <https://github.com/yongqyu/TCN>.

### Baselines

**CycleGAN** In recent years, CycleGAN(Zhu et al. 2017) has obtained outstanding performance in the image-to-image translation task. CycleGAN was the first to use cycle consistency which makes the image that was once converted back to the original domain equal to the original image, as in Equation 14. We also use vector-wise cycle consistency and pixel-wise cycle consistency at the image level.

**MUNIT** MUNIT(Huang et al. 2018) extracts the typeface and content features as a form of a latent vector using each encoder. By switching these vectors, an image with the desired features can be obtained. MUNIT uses a latent vector,

and the reconstruction loss that uses the latent vector is similar to our perceptual reconstruction loss. However, there is a difference in the concept of reconstruction: we translate twice so that the reconstructed image is the same as the original, but MUNIT translates only once. Another difference is that MUNIT needs two images to generate every image. This is not only dependent on content input image, but also inefficient if the target label is fixed.

**StarGAN** StarGAN(Choi et al. 2017) passes an image with the desired domain label to the generator, like cGAN. To this end, the discriminator returns the true likelihood of the image, along with the domain to which the image corresponds. As a result, StarGAN can generate all characters in one model, like our model. However, unlike our model, StarGAN uses a one-hot vector to represent a content vector. The comparison of the above baselines and TCN is summarized in Table 3.

### Experiment

The single-domain transfer models cannot generate the entire character set using one model. For a fair comparison with the single-domain transfer model, we used two experimental conditions. First, we used sample pairs of characters. In English, Y-G and Q-G pairs were selected to represent the most different and similar pairs, respectively. In Chinese, index number 598-268, and 598-370 pairs were selected. After the sample pair experiments, we compared the performance of our model with that of a multi-domain model on translating all pairs. In these two experimental conditions, we performed the following subtasks: Typeface Completion and Character Reconstruction.

**Typeface Completion** In the image-to-image task, our model takes one character image and learns to complete the rest of the character set while maintaining its typeface. We used a single character image as an input for a fair comparison with the other models. We trained our model on Chinese and English character sets, which we mentioned above.

The image-to-image translation experiment allows us to evaluate the performance of the two encoders in extracting the disjoint features. If the typeface encoder extracts the content feature and the typeface feature, the generated image will have the same content of the typeface input. This also applies to content encoders. In training, the model processes every content of a character set. And in the test, every content of the character set can be generated using the extracted typeface feature, even if the input typeface is new to the typeface encoder. Since the character image set has the target pair and the input, we can objectively evaluate the result based on its score.

**Character Reconstruction** Reconstruction is the process of regenerating an input image using the typeface and content features extracted from the input image. By the reconstruction, we can check if there is any feature missing when the encoder extracts features. It is also possible to check whether the decoder can effectively combine the two types of features. However, in this task, it is not possible to verify whether each of the features is disjointed or overlapped.

Table 4: Ablation Study Results

Input	Target	TCN	$(-)L_{ssim}$	$(-)L_{id}$	$(-)L_{ssim} + L_{id}$	$(-)L_{rec}$	$(-)L_{per}$	$(-)L_{rec} + L_{per}$	$(-)y_{input}$
M	R	R	R	R	R	R	R	R	R
U	K	K	K	K	K	K	K	K	K
H	R	R	R	R	R	R	R	R	R
B	C	C	C	C	C	C	C	C	C
<i>SSIM</i>		0.793	0.765	0.784	0.780	0.771	0.773	0.781	0.784

Table 5: All columns have input as the first column, From left to right: TCN, eliminate SSIM loss, eliminate Identity loss, eliminate SSIM and Identity loss, eliminate reconstruction loss, eliminate perceptual reconstruction loss, eliminate reconstruction and perceptual reconstruction loss. And lastly, generate without target label.

**Ablation Study** We conducted an ablation study to check for redundancy among the various losses in our model. By comparing the performance of eliminating each loss and the performance of the entire model, we can check the influence of each loss. In addition to loss, we also conducted an ablation study on the sub-modules to check the influence of the sub-modules. The ablation study were conducted in Chinese and English datasets, and results were quantitatively and qualitatively evaluated.

**Face Generation** As our model is not limited to character images, we experimented with facial images used for existing image-to-image models. In the facial image experiment, the differences from the character images are that there are no content labels and no target images. Therefore we proceeded the experiment after removed the associated losses and sub-module of TCN. We take a face image and perform a image-to-image translation experiment that changes the style feature label. We also conducted weighted image-to-image translation experiments on weighted style feature labels. Since we cannot quantitatively evaluate in the unpaired dataset, we conducted a quantitative evaluation through comparison only.

## Analysis

The reconstruction performance and the typeface completion performance of single-domain image-to-image models (CycleGAN, MUNIT) vary (Table 3) due to insufficient information of the features. When extracting features from character images, style and content features are duplicated or lost, not being disjoint, which is demonstrated by the translation results of these models. The output of the single-domain models is dependent on the input image, so the models achieve high performance in the reconstruction task where the target is the input. On the other hand, in the typeface completion task, the result appears to be a simple combination of inputs rather than a image-to-image translation.

StarGAN obtained good performance on the general images of the image-to-image translation task, but not on the

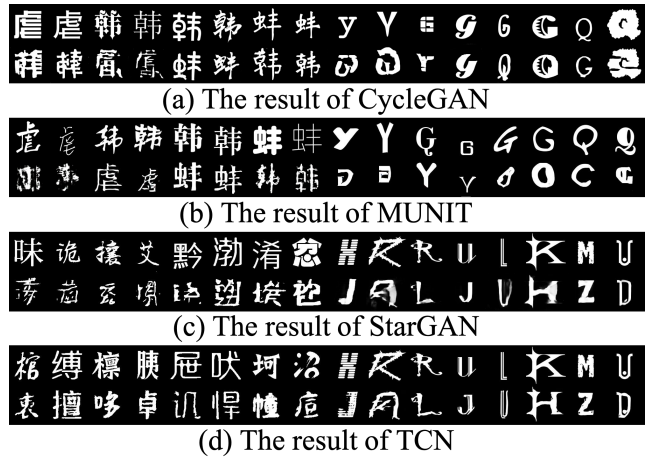


Figure 3: The results of baselines and our model. The first row contains the typeface source image, the second row contains the output.

character dataset. For the StarGAN, content accuracy is similar to our model, but typeface accuracy is not. Because StarGAN uses the reconstruction loss to maintain typeface information, but it places a greater weight on content translation. The resulting image is fairly clean at the character level, but there is a limit to maintaining the typeface information of the input. (Figure 3(c)).

Another difference is that, when calculating a loss in the one-hot vector, the cosine similarity of each vector is either one or zero. Therefore, we can only determine whether two values are matched. To address this issue, we use a latent vector as a domain label which has continuous values for similarity scores between vectors. From the values, the vector determines the similarity and difference of the two vectors, which can help the classifier to learn.

Another difference between TCN and the other models is the use of input labels. Adding input labels for the model results in output images more similar to the real images. Due

Table 6: Sample Pair Experimental Results

	#Parameter		SSIM									
			Typeface Completion				Reconstruction					
	Ch	Eng	598-370	598-268	G-Q	G-Y	598	370	268	G	Q	Y
<i>CycleGAN</i>	2.8e+7 x $N^2$		0.5616	0.5442	0.6915	0.5295	0.8740	0.8467	0.8670	0.9108	0.8941	0.8147
<i>MUNIT</i>	4.6e+7 x $N^2$		0.5438	0.5229	0.6268	0.7028	<b>0.9941</b>	<b>0.9956</b>	<b>0.9937</b>	<b>0.9962</b>	<b>0.9936</b>	<b>0.9960</b>
<i>StarGAN</i>	6.4e+7	5.3e+7	0.5506	0.5443	0.6985	0.7274	0.5413	0.5813	0.5624	0.7463	0.7026	0.7944
<i>TCN</i>	5.6e+7 (7.9e+7)	2.4e+7 (4.5e+7)	<b>0.6673</b>	<b>0.6573</b>	<b>0.7959</b>	<b>0.8264</b>	0.6609	0.7011	0.6862	0.8116	0.8015	0.8604

Table 7: Total Experimental Results

	SSIM				L1				Style Accuracy				Content Accuracy			
	TC		Reconst		TC		Reconst		TC		Reconst		TC		Reconst	
	Ch	Eng	Ch	Eng	Ch	Eng	Ch	Eng	Ch	Eng	Ch	Eng	Ch	Eng	Ch	Eng
<i>StarGAN</i>	0.568	0.749	0.575	0.857	0.212	0.094	0.191	0.046	0.009	0.737	0.009	0.851	0.402	<b>0.820</b>	0.523	0.951
<i>TCN</i>	<b>0.653</b>	<b>0.794</b>	<b>0.676</b>	<b>0.949</b>	<b>0.163</b>	<b>0.088</b>	<b>0.060</b>	<b>0.014</b>	<b>0.645</b>	<b>0.891</b>	<b>0.802</b>	<b>0.998</b>	<b>0.531</b>	0.819	<b>0.972</b>	<b>0.991</b>

to the differences described above, our model outperformed StarGAN by 10% on typeface completion and 12% on reconstruction at SSIM index, as shown in Table 5. And as shown in Figure 4, generated images have consistent typeface of input image. Even the results are unseen typefaces in the training process. Nonetheless, TCN generates an image similar to the target.

We also combined the one-hot label with the encoded latent vector rather than the original image. This improves the parameter efficiency of the model. In Table 4, the number of parameters of the TCN is the smallest. The number in parentheses include the parameters of the classifier used in the pre-train. This module is not used for main-train and inference, hence it is indicated separately. The number of parameters, except for this, is 12.5% and 54% decrease in Chinese and English, respectively, compared with StarGAN.

In ablation study, The full model showed the best performance quantitatively and qualitatively. SSIM loss had the greatest influence on SSIM index in the test set. However, the absence of SSIM loss has little affected another index beyond the SSIM index. The elimination of identity loss did not affect typeface completion performance, but it was the most influential in the reconstruction task. These two losses are available because the dataset is paired. In the model, which eliminated both SSIM loss and identity loss, which did not utilize the target image, we obtained the intermediate result of the previous two studies. Because the two losses are complementary, the overall result is worse when there is one loss only. The reconstruction loss and perceptual reconstruction loss associated with cycle consistency had similar effects on performance. In the absence of cycle consistency, we saw that typeface accuracy is increasing and content accuracy is significantly lowered. Lastly, the absence of input content labels affected overall performance downgrades. As we can show from the qualitative comparison, it is involved in the detailed result without any significant difference from the result of the full model.

Our model can also be applied to unpaired general images. Since we do not have content information, We experimented after removing the content encoder. As shown in

Figure 5, we made a fairly plausible outcome. And when we compared the StarGAN, we found that our model maintains better out-of-style information that has not changed. We also were able to generate the image to maintain a specific ratio between the input and the target, like Figure 6. This suggests that our model can be used for various applications.

## Conclusion

In this paper, we proposed Typeface Completion Network (TCN) which generates an entire set of characters given only one characters while maintaining the typeface of the input characters. TCN utilizes the typeface and content encoders to effectively leverage the information of numerous classes. As a result, TCN learns multi-domain image-to-image translation using a single model, and produces more accurate outputs than existing baseline models. As illustrated in the qualitative analysis, we found that TCN successfully completes the character sets, which could reduce the costs of designing a new typeface. We also tested TCN on the CelebA dataset to demonstrate its applicability. In future work, we are planning a model that takes character subset as input and completes character set more finely.

## References

- [Arjovsky, Chintala, and Bottou 2017] Arjovsky, M.; Chintala, S.; and Bottou, L. 2017. Wasserstein generative adversarial networks. In *International Conference on Machine Learning*.
- [Azadi et al. 2018] Azadi, S.; Fisher, M.; Kim, V.; Wang, Z.; Shechtman, E.; and Darrell, T. 2018. Multi-content gan for few-shot font style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 11, 13.
- [Baluja 2016] Baluja, S. 2016. Learning typographic style. *CoRR*.
- [Bhunia et al. 2018] Bhunia, A. K.; Bhunia, A. K.; Banerjee, P.; Konwer, A.; Bhowmick, A.; Roy, P. P.; and Pal, U. 2018. Word level font-to-font image translation using convolutional recurrent generative adversarial networks. *CoRR*.



Figure 4: The figure consists of sections of which each two row represents the work on same typeface, respectively. Top left image is input image and second row are generated images in each section. The first row of each section are label images corresponding to second row.

- [Campbell and Kautz 2014] Campbell, N. D. F., and Kautz, J. 2014. Learning a manifold of fonts. *ACM Trans. Graph.*
- [Chang and Gu 2017] Chang, J., and Gu, Y. 2017. Chinese typography transfer. *arXiv preprint arXiv:1707.04904*.
- [Chang et al. 2018] Chang, B.; Zhang, Q.; Pan, S.; and Meng, L. 2018. Generating handwritten chinese characters using cyclegan. *arXiv preprint arXiv:1801.08624*.
- [Chen and Koltun 2017] Chen, Q., and Koltun, V. 2017. Photographic image synthesis with cascaded refinement networks. In *The IEEE International Conference on Computer Vision (ICCV)*.
- [Choi et al. 2017] Choi, Y.; Choi, M.; Kim, M.; Ha, J.-W.; Kim, S.; and Choo, J. 2017. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. *arXiv preprint arXiv:1711.09020*.
- [Gatys, Ecker, and Bethge 2016] Gatys, L. A.; Ecker, A. S.; and Bethge, M. 2016. Image style transfer using convolutional neural networks. In *Computer Vision and Pattern Recognition (CVPR), 2016 IEEE Conference on*.
- [Goodfellow et al. 2014] Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative adversarial nets. In *Advances in neural information processing systems*.



Figure 5: The results of the style transfer for CelebA. The four columns on the left and the right are the results of StarGAN and TCN. The first column of each result is the input, and the remaining columns are the result of changing the domain label to black / blond / brown hair, respectively.

- [Goodfellow 2017] Goodfellow, I. J. 2017. NIPS 2016 tutorial: Generative adversarial networks. *CoRR*.
- [Huang and Belongie 2017] Huang, X., and Belongie, S. 2017. Arbitrary style transfer in real-time with adaptive instance normalization. *CoRR, abs/1703.06868*.
- [Huang et al. 2018] Huang, X.; Liu, M.; Belongie, S. J.; and Kautz, J. 2018. Multimodal unsupervised image-to-image translation. *CoRR*.
- [Johnson, Alahi, and Fei-Fei 2016] Johnson, J.; Alahi, A.; and Fei-Fei, L. 2016. Perceptual losses for real-time style transfer and super-resolution. In *Computer Vision – ECCV 2016*. Springer International Publishing.
- [Ledig et al. 2017] Ledig, C.; Theis, L.; Huszar, F.; Caballero, J.; Cunningham, A.; Acosta, A.; Aitken, A.; Tejani, A.; Totz, J.; Wang, Z.; and Shi, W. 2017. Photo-realistic single image super-resolution using a generative adversarial network. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [Li et al. 2017] Li, Y.; Fang, C.; Yang, J.; Wang, Z.; Lu, X.; and Yang, M.-H. 2017. Universal style transfer via feature transforms. In *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc.
- [Li et al. 2018] Li, Y.; Liu, M.; Li, X.; Yang, M.; and Kautz, J. 2018. A closed-form solution to photorealistic image stylization. *CoRR*.
- [Liu et al. 2015] Liu, Z.; Luo, P.; Wang, X.; and Tang, X. 2015. Deep learning face attributes in the wild. In *Proceedings of the IEEE International Conference on Computer Vision*.
- [Lyu et al. 2017] Lyu, P.; Bai, X.; Yao, C.; Zhu, Z.; Huang, T.; and Liu, W. 2017. Auto-encoder guided GAN for chinese calligraphy synthesis. *CoRR*.
- [Mirza and Osindero 2014] Mirza, M., and Osindero, S. 2014. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*.
- [Odena, Olah, and Shlens 2016] Odena, A.; Olah, C.; and



Figure 6: The result of the weighted style transfer from blond to black hair. The four columns on the left and right are the results of StarGAN and TCN, respectively. First column is the original image, and the weight of black color is assigned to 0.3, 0.6, and 0.9 for the rest. The hair color is somewhat reddish in StarGAN when the black hair weight is low, but TCN expresses the color between the blond and black hair based on the weight.

Shlens, J. 2016. Conditional image synthesis with auxiliary classifier gans. *arXiv preprint arXiv:1610.09585*.

[Odena 2016] Odena, A. 2016. Semi-supervised learning with generative adversarial networks. *arXiv preprint arXiv:1606.01583*.

[Paszke et al. 2017] Paszke, A.; Gross, S.; Chintala, S.; and Chanan, G. 2017. Pytorch.

[Perarnau et al. 2016] Perarnau, G.; van de Weijer, J.; Raducanu, B.; and Álvarez, J. M. 2016. Invertible conditional gans for image editing. *arXiv preprint arXiv:1611.06355*.

[Phan, Fu, and Chan 2015] Phan, Q. H.; Fu, H.; and Chan, A. B. 2015. Flexyfont: Learning transferring rules for flexible typeface synthesis. *Comput. Graph. Forum*.

[Radford, Metz, and Chintala 2015] Radford, A.; Metz, L.; and Chintala, S. 2015. Unsupervised representation learning with deep convolutional generative adversarial networks. *CoRR*.

[Schroff, Kalenichenko, and Philbin 2015] Schroff, F.; Kalenichenko, D.; and Philbin, J. 2015. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.

[Snell et al. 2017] Snell, J.; Ridgeway, K.; Liao, R.; Roads, B. D.; Mozer, M. C.; and Zemel, R. S. 2017. Learning to generate images with perceptual similarity metrics. In *Image Processing (ICIP), 2017 IEEE International Conference on*.

[Suveeranont and Igarashi 2010] Suveeranont, R., and Igarashi, T. 2010. Example-based automatic font generation. In *Smart Graphics*.

[Tenenbaum and Freeman 1997] Tenenbaum, J. B., and Freeman, W. T. 1997. Separating style and content. In *Advances in neural information processing systems*.

[Upchurch, Snaveley, and Bala 2016] Upchurch, P.; Snaveley,

N.; and Bala, K. 2016. From A to Z: supervised transfer of style and content using deep neural network generators. *CoRR*.

[Wu, Xu, and Hall 2017] Wu, X.; Xu, K.; and Hall, P. 2017. A survey of image synthesis and editing with generative adversarial networks. *Tsinghua Science and Technology* 22(6):660–674.

[Zhao et al. 2017] Zhao, H.; Gallo, O.; Frosio, I.; and Kautz, J. 2017. Loss functions for image restoration with neural networks. *IEEE Transactions on Computational Imaging*.

[Zhu et al. 2017] Zhu, J.-Y.; Park, T.; Isola, P.; and Efros, A. A. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. *arXiv preprint arXiv:1703.10593*.