

# Generate Like Experts: Multi-Stage Font Generation by Incorporating Font Transfer Process into Diffusion Models

## Supplementary Materials

### 6. Proofs and Derivations

In this section, we provide the detailed proofs and derivations for our proposed model. Our derivations mainly follow the derivations of denoising diffusion probabilistic models in Ref. [1, 2], and people can refer to them for more details.

#### 6.1. The Derivation of Eq. (9)

In this paper, our proposed model is constructed by utilizing a font transfer process to connect the traditional diffusion processes (DDPM) of the latent features  $z_t^c$  and  $z_t^g$ , where the latent feature  $z_t^g$  of the target font is gradually replaced by the latent feature  $z_t^c$  of the source font from  $t = t_1$  to  $t = t_2$ . Therefore, in data sub-space, only  $z_t^g$  exists at  $t = t_1$  while only  $z_t^c$  exists at  $t = t_2$ . The corresponding latent features  $\tilde{z}_{t_1}$  and  $\tilde{z}_{t_2}$  can be expressed as

$$\tilde{z}_{t_1} = z_{t_1}^g = \sqrt{\bar{\alpha}_{t_1}} z_0^g + \sqrt{1 - \bar{\alpha}_{t_1}} \epsilon_0, \quad (24)$$

$$\tilde{z}_{t_2} = z_{t_2}^c = \sqrt{\bar{\alpha}_{t_2}} z_0^c + \sqrt{1 - \bar{\alpha}_{t_2}} \epsilon_0. \quad (25)$$

As the latent feature  $z_t^g$  will be gradually replaced by the latent feature  $z_t^c$  in the font transfer region  $t \in (t_1, t_2]$ , we formulate the forward latent feature  $\tilde{z}_t$  in this region as

$$\tilde{z}_t = \sqrt{\alpha_t} \tilde{z}_{t-1} + \sqrt{1 - \alpha_t} \epsilon_{t-1}^* - \Psi(z_0^g - z_0^c), \quad (26)$$

where  $\{\epsilon_t^*, \epsilon_t\} \sim \mathcal{N}(0, I)$ .  $\Psi$  is the coefficient for the term  $z_0^g - z_0^c$ , which will be determined in this section. With the reparameterization trick, the forward latent  $\tilde{z}_t$  in region  $(t_1, t_2]$  can be written as

$$\tilde{z}_t = \sqrt{\alpha_t} \tilde{z}_{t-1} + \sqrt{1 - \alpha_t} \epsilon_{t-1}^* - \Psi(z_0^g - z_0^c) \quad (27)$$

$$= \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}} \tilde{z}_{t-2} + \sqrt{1 - \alpha_{t-1}} \epsilon_{t-2}^* - \Psi(z_0^g - z_0^c) \right) + \sqrt{1 - \alpha_t} \epsilon_{t-1}^* - \Psi(z_0^g - z_0^c) \quad (28)$$

$$= \sqrt{\alpha_t \alpha_{t-1}} \tilde{z}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \epsilon_{t-2} - (1 + \sqrt{\alpha_t}) \Psi(z_0^g - z_0^c) \quad (29)$$

$$= \sqrt{\alpha_t \alpha_{t-1} \cdots \alpha_{t_1+1}} \tilde{z}_{t_1} + \sqrt{1 - \alpha_t \alpha_{t-1} \cdots \alpha_{t_1+1}} \epsilon_{t_1} - \mathcal{H}_t \Psi(z_0^g - z_0^c), \quad (30)$$

and this equation can be further reduced by aligning the forward latent with the traditional diffusion process (DDPM) at  $t = t_1$ . Using Eq. (24), the latent feature  $\tilde{z}_t$  in region  $(t_1, t_2]$  can be expressed as

$$\begin{aligned} \tilde{z}_t &= \sqrt{\alpha_t \alpha_{t-1} \cdots \alpha_{t_1+1}} \left( \sqrt{\bar{\alpha}_{t_1}} z_0^g + \sqrt{1 - \bar{\alpha}_{t_1}} \epsilon_0^* \right) + \sqrt{1 - \alpha_t \alpha_{t-1} \cdots \alpha_{t_1+1}} \epsilon_{t_1} - \mathcal{H}_t \Psi(z_0^g - z_0^c) \\ &= \sqrt{\alpha_t \alpha_{t-1} \cdots \alpha_1} z_0^g + \sqrt{1 - \alpha_t \alpha_{t-1} \cdots \alpha_1} \epsilon_0 - \mathcal{H}_t \Psi(z_0^g - z_0^c) \\ &= \sqrt{\bar{\alpha}_t} z_0^g + \sqrt{1 - \bar{\alpha}_t} \epsilon_0 - \mathcal{H}_t \Psi(z_0^g - z_0^c), \end{aligned} \quad (31)$$

where we have defined  $\mathcal{H}_t = 1 + \sqrt{\alpha_t} + \sqrt{\alpha_t \alpha_{t-1}} + \cdots + \sqrt{\alpha_t \alpha_{t-1} \cdots \alpha_{t_1+2}}$ , and  $\mathcal{H}_t$  has the following recursion formula,

$$\begin{aligned} \mathcal{H}_t &= 1 + \sqrt{\alpha_t} + \sqrt{\alpha_t \alpha_{t-1}} + \cdots + \sqrt{\alpha_t \cdots \alpha_{t_1+2}} \\ &= 1 + \sqrt{\alpha_t} (1 + \sqrt{\alpha_{t-1}} + \sqrt{\alpha_{t-1} \alpha_{t-2}} + \cdots + \sqrt{\alpha_{t-1} \alpha_{t-2} \alpha_{t-3} \cdots \alpha_{t_1+2}}) \\ &= 1 + \sqrt{\alpha_t} \mathcal{H}_{t-1}, \end{aligned} \quad (32)$$

which has been utilized in our implementation for calculating  $\mathcal{H}_t$ . The coefficient  $\Psi$  can be determined by aligning the forward latent with the traditional diffusion process at  $t = t_2$ ,

$$\Psi = \frac{\sqrt{\bar{\alpha}_{t_2}}}{\mathcal{H}_{t_2}} = \frac{\sqrt{\tilde{\alpha}(1, t_2)}}{1 + \sum_{m=t_1+2}^{t_2} \sqrt{\tilde{\alpha}(m, t_2)}}, \quad (33)$$

where we have defined  $\tilde{\alpha}(t_1, t_2) = \prod_{i=t_1}^{t_2} \alpha_i$ . Note that  $\mathcal{H}_t$  is a function with respect to time  $t$ , while  $\Psi$  is a constant once we have defined  $t_1$  and  $t_2$ .

Therefore, with the above derivations, the forward latent  $\tilde{z}_t$  in the region  $(t_1, t_2]$  can be expressed as

$$\tilde{z}_t = \sqrt{\bar{\alpha}_t} z_0^g + \sqrt{1 - \bar{\alpha}_t} \epsilon_0 - (z_0^g - z_0^c) \Psi \mathcal{H}_t, \quad (34)$$

which is the Eq. (9) in the main body.

## 6.2. The Derivation of Eq. (15)

Based on the variational bound on negative log-likelihood (Eq. (5) of DDPM [1]), we learn  $p_\theta(\tilde{z}_{t-1}|\tilde{z}_t)$  as an approximation to the ground-truth denoising transition step  $q(\tilde{z}_{t-1}|\tilde{z}_t, \tilde{z}_0)$ , where  $\tilde{z}_0 = z_0^g$ . Once we obtain the forward latent  $\tilde{z}_t$  of our diffusion model, we can derive  $q(\tilde{z}_{t-1}|\tilde{z}_t, \tilde{z}_0)$  in region  $(t_1, t_2]$  by using Bayes' theorem:

$$q(\tilde{z}_{t-1}|\tilde{z}_t, \tilde{z}_0) = \frac{q(\tilde{z}_t|\tilde{z}_{t-1}, \tilde{z}_0) q(\tilde{z}_{t-1}|\tilde{z}_0)}{q(\tilde{z}_t|\tilde{z}_0)}. \quad (35)$$

With Eq. (27) and Eq. (34), this equation can be further reduced to

$$\begin{aligned} q(\tilde{z}_{t-1}|\tilde{z}_t, \tilde{z}_0) &= \frac{\mathcal{N}(\tilde{z}_t; \sqrt{\alpha_t} \tilde{z}_{t-1} - \Psi(z_0^g - z_0^c), (1 - \alpha_t)I) \mathcal{N}(\tilde{z}_{t-1}; \sqrt{\bar{\alpha}_{t-1}} z_0^g - \mathcal{H}_{t-1} \Psi(z_0^g - z_0^c), (1 - \bar{\alpha}_{t-1})I)}{\mathcal{N}(\tilde{z}_t; \sqrt{\bar{\alpha}_t} z_0^g - \mathcal{H}_t \Psi(z_0^g - z_0^c), (1 - \bar{\alpha}_t)I)} \\ &\propto \exp \left\{ -\frac{1}{2} \left[ \frac{(\tilde{z}_t - (\sqrt{\alpha_t} \tilde{z}_{t-1} - \Psi(z_0^g - z_0^c)))^2}{1 - \alpha_t} + \frac{(\tilde{z}_{t-1} - (\sqrt{\bar{\alpha}_{t-1}} z_0^g - \mathcal{H}_{t-1} \Psi(z_0^g - z_0^c)))^2}{1 - \bar{\alpha}_{t-1}} \right. \right. \\ &\quad \left. \left. - \frac{(\tilde{z}_t - (\sqrt{\bar{\alpha}_t} z_0^g - \mathcal{H}_t \Psi(z_0^g - z_0^c)))^2}{1 - \bar{\alpha}_t} \right] \right\} \\ &\propto \exp \left\{ -\frac{1}{2} \left[ \left( \frac{\alpha_t}{1 - \alpha_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) \tilde{z}_{t-1}^2 \right. \right. \\ &\quad \left. \left. - 2 \left( \frac{\sqrt{\alpha_t} \tilde{z}_t + \sqrt{\alpha_t} \Psi(z_0^g - z_0^c)}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}} z_0^g - \mathcal{H}_{t-1} \Psi(z_0^g - z_0^c)}{1 - \bar{\alpha}_{t-1}} \right) \tilde{z}_{t-1} \right] \right\} \\ &= \exp \left\{ -\frac{1}{2} \left[ \left( \frac{\alpha_t (1 - \bar{\alpha}_{t-1}) + 1 - \alpha_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \right) \tilde{z}_{t-1}^2 \right. \right. \\ &\quad \left. \left. - 2 \left( \frac{\sqrt{\alpha_t} \tilde{z}_t + \sqrt{\alpha_t} \Psi(z_0^g - z_0^c)}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}} z_0^g - \mathcal{H}_{t-1} \Psi(z_0^g - z_0^c)}{1 - \bar{\alpha}_{t-1}} \right) \tilde{z}_{t-1} \right] \right\} \\ &= \exp \left\{ -\frac{1}{2} \left( \frac{1 - \bar{\alpha}_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \right) \right. \\ &\quad \left. \left[ \tilde{z}_{t-1}^2 - 2 \frac{\left( \frac{\sqrt{\alpha_t} \tilde{z}_t + \sqrt{\alpha_t} \Psi(z_0^g - z_0^c)}{1 - \alpha_t} + \frac{\sqrt{\bar{\alpha}_{t-1}} z_0^g - \mathcal{H}_{t-1} \Psi(z_0^g - z_0^c)}{1 - \bar{\alpha}_{t-1}} \right) (1 - \alpha_t) (1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \tilde{z}_{t-1} \right] \right\} \\ &= \exp \left\{ -\frac{1}{2} \left( \frac{1 - \bar{\alpha}_t}{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})} \right) \right. \\ &\quad \left. \left[ \tilde{z}_{t-1}^2 - 2 \frac{\sqrt{\alpha_t} (\tilde{z}_t + \Psi(z_0^g - z_0^c)) (1 - \bar{\alpha}_{t-1}) + (\sqrt{\bar{\alpha}_{t-1}} z_0^g - \mathcal{H}_{t-1} \Psi(z_0^g - z_0^c)) (1 - \alpha_t)}{1 - \bar{\alpha}_t} \tilde{z}_{t-1} \right] \right\} \\ &\propto \mathcal{N}(\tilde{z}_{t-1}; \mu_q(\tilde{z}_t, z_0^g, t), \Sigma_q(t)), \end{aligned} \quad (36)$$

where  $\mu_q(\tilde{z}_t, z_0^g, t)$  and  $\Sigma_q(t)$  are defined as

$$\mu_q(\tilde{z}_t, z_0^g, t) = \frac{\sqrt{\alpha_t} (\tilde{z}_t + \Psi(z_0^g - z_0^c)) (1 - \bar{\alpha}_{t-1}) + (\sqrt{\bar{\alpha}_{t-1}} z_0^g - \mathcal{H}_{t-1} \Psi(z_0^g - z_0^c)) (1 - \alpha_t)}{1 - \bar{\alpha}_t}, \quad (37)$$

$$\Sigma_q(t) = \sigma_t^2 I = \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} I. \quad (38)$$

Note that the variance remains the same as DDPM, since we only incorporate font transfer in the data sub-space and keep the noise sub-space unchanged. We further rewrite the mean function  $\mu_q(\tilde{z}_t, z_0^g, t)$  as

$$\begin{aligned}\mu_q(\tilde{z}_t, z_0^g, t) &= \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}\tilde{z}_t + \frac{\mathcal{H}_{t-1}(1-\alpha_t)-\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}\Psi z_0^c \\ &+ \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\Psi+(\sqrt{\alpha_{t-1}}-\mathcal{H}_{t-1}\Psi)(1-\alpha_t)}{1-\bar{\alpha}_t}z_0^g.\end{aligned}\quad (39)$$

In order to match  $p_\theta(\tilde{z}_{t-1}|\tilde{z}_t)$  to the ground-truth denoising transition step  $q(\tilde{z}_{t-1}|\tilde{z}_t, \tilde{z}_0)$ , we model  $p_\theta(\tilde{z}_{t-1}|\tilde{z}_t)$  as a Gaussian with  $p_\theta(\tilde{z}_{t-1}|\tilde{z}_t) = \mathcal{N}(\tilde{z}_{t-1}|\mu_\theta(\tilde{z}_t, t), \Sigma_\theta(\tilde{z}_t, t))$ . We first set the variances of  $p_\theta(\tilde{z}_{t-1}|\tilde{z}_t)$  and  $q(\tilde{z}_{t-1}|\tilde{z}_t, \tilde{z}_0)$  to match exactly, thus

$$\Sigma_\theta(\tilde{z}_t, t) = \sigma_t^2 I = \frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} I. \quad (40)$$

Once the variances of the two Gaussian match exactly, Ref. [2] further proves that optimizing the KL-divergence between  $q(\tilde{z}_{t-1}|\tilde{z}_t, \tilde{z}_0)$  and  $p_\theta(\tilde{z}_{t-1}|\tilde{z}_t)$  reduces to minimize the difference between the means of these distributions. Therefore, following the common practice [1, 2], we utilize the prediction network  $\tilde{z}_\theta^{(g)}(\tilde{z}_t, t)$  to predict  $z_0^g$ , and  $\mu_\theta(\tilde{z}_t, t)$  can be expressed as

$$\begin{aligned}\mu_\theta(\tilde{z}_t, t) &= \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}\tilde{z}_t + \frac{\mathcal{H}_{t-1}(1-\alpha_t)-\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}\Psi z_0^c \\ &+ \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\Psi+(\sqrt{\alpha_{t-1}}-\mathcal{H}_{t-1}\Psi)(1-\alpha_t)}{1-\bar{\alpha}_t}\tilde{z}_\theta^{(g)}(\tilde{z}_t, t),\end{aligned}\quad (41)$$

which is the Eq. (15) in the main body.

## 7. Implement Details

In this section, we provide more implementation details about our multi-stage font generation model.

### 7.1. Multi-Stage Font Generative Process

We provide the pseudo-code of multi-stage font generative process in Algorithm 1.

---

#### Algorithm 1: Multi-Stage Font Generative Process

---

**Input:** The coefficients  $\alpha_t, \bar{\alpha}_t, \sigma_t, \mathcal{H}_t$ , and  $\Psi$ . The latent feature  $z_0^c$  of the source image. The conditions  $\mathbf{y}_1$  and  $\mathbf{y}_2$  for font transfer stage and font refinement stage, respectively.

**Output:** The generated font image  $I_G$ .

```
// step 1: Structure Construction Stage
1  $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$ ;
2  $\tilde{z}_{t_2} = \sqrt{\alpha_{t_2}}z_0^c + \sqrt{1-\bar{\alpha}_{t_2}}\epsilon$ ;
   // step 2: Font Transer Stage
3 for  $t = t_2$  to  $t = t_1 + 1$  do
4    $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$ ;
5    $\mu_\theta = \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}\tilde{z}_t + \left(\frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\Psi+(\sqrt{\alpha_{t-1}}-\mathcal{H}_{t-1}\Psi)(1-\alpha_t)}{1-\bar{\alpha}_t}\right)\tilde{z}_{\theta_1}^{(g,1)}(\tilde{z}_t, t, \mathbf{y}_1) + \frac{[\mathcal{H}_{t-1}(1-\alpha_t)-\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})]}{1-\bar{\alpha}_t}\Psi z_0^c$ ;
6    $\tilde{z}_{t-1} = \mu_\theta + \sigma_t\epsilon$ ;
   // step 3: Font Refinement Stage
7 for  $t = t_1$  to  $t = 1$  do
8    $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$  if  $t > 1$ , else  $\epsilon = \mathbf{0}$ ;
9    $\mu_\theta = \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}\tilde{z}_t + \frac{\sqrt{\alpha_{t-1}}(1-\alpha_t)}{1-\bar{\alpha}_t}\tilde{z}_{\theta_2}^{(g,2)}(\tilde{z}_t, t, \mathbf{y}_2)$ ;
10   $\tilde{z}_{t-1} = \mu_\theta + \sigma_t\epsilon$ ;
11  $z_0^g = \tilde{z}_0$ ;
12 return  $I_G = \mathcal{D}(z_0^g)$ 
```

---

## 7.2. Optimization Process

We use a two-stage optimization process to train our proposed model. The pseudo-code of the optimization process is provided in Algorithm 2.

**Algorithm 2:** Optimization Process

---

```

// Stage 1-1: Optimize  $E_c^1$ ,  $E_s^1$ , and  $\tilde{z}_{\theta_1}^{(g,1)}(\tilde{z}_t, t, \mathbf{y}_1)$  in Region  $(t_1, t_2]$ 
1 while True do
2   Loading training pairs  $(I_s, I_c, I_g)$  by picking  $I_s$  and  $I_c$  based on  $I_g$ ;
3   Calculating condition  $y_1 = \mathcal{C}(\mathcal{F}(e_c^1), \mathcal{F}(e_s^1))$  by flattening and concatening  $e_c^1 = E_c^1(I_c)$  and  $e_s^1 = E_s^1(I_s)$ ;
4   Calculating  $z_0^g = \mathcal{E}(I_g)$  and  $z_0^c = \mathcal{E}(I_c)$ ;
5   Sampling  $t \in (t_1, t_2]$  and calculating  $\tilde{z}_t$  by Eq. (9);
6   Using conditional prediction network  $\tilde{z}_{\theta_1}^{(g,1)}(\tilde{z}_t, t, y_1)$  to predict  $z_0^g$ ;
7   Using loss function Eq. (22) to train  $E_c^1, E_s^1, \tilde{z}_{\theta_1}^{(g,1)}(\tilde{z}_t, t, \mathbf{y}_1)$  in an end-to-end manner;
8   Until Converged;

// Stage 1-2: Optimize  $E_c^2$ ,  $E_s^2$ , and  $\tilde{z}_{\theta_2}^{(g,2)}(\tilde{z}_t, t, \mathbf{y}_2)$  in Region  $(0, t_1]$ 
9 while True do
10  Loading training pairs  $(I_s, I_c, I_g)$  by picking  $I_s$  and  $I_c$  based on  $I_g$ ;
11  Calculating condition  $y_2 = \mathcal{C}(\mathcal{F}(e_c^2), \mathcal{F}(e_s^2))$  by flattening and concatening  $e_c^2 = E_c^2(I_c)$  and  $e_s^2 = E_s^2(I_s)$ ;
12  Calculating  $z_0^g = \mathcal{E}(I_g)$  and  $z_0^c = \mathcal{E}(I_c)$ ;
13  Sampling  $t \in (0, t_1]$  and calculating  $\tilde{z}_t$  by Eq. (8);
14  Using conditional prediction network  $\tilde{z}_{\theta_2}^{(g,2)}(\tilde{z}_t, t, y_2)$  to predict  $z_0^g$ ;
15  Using loss function Eq. (23) to train  $E_c^2, E_s^2, \tilde{z}_{\theta_2}^{(g,2)}(\tilde{z}_t, t, \mathbf{y}_2)$  in an end-to-end manner;
16  Until Converged;

// Stage 2: Using  $\tilde{z}_{\theta_1}^{(g,1)}(\tilde{z}_t, t, \mathbf{y}_1)$  to Fine-tune  $E_c^2$ ,  $E_s^2$ , and  $\tilde{z}_{\theta_2}^{(g,2)}(\tilde{z}_t, t, \mathbf{y}_2)$  in Region
// (0, t_1]
17 while True do
18  Loading training pairs  $(I_s, I_c, I_g)$  by picking  $I_s$  and  $I_c$  based on  $I_g$ ;
19  Calculating condition  $y_1 = \mathcal{C}(\mathcal{F}(e_c^1), \mathcal{F}(e_s^1))$  by flattening and concatening  $e_c^1 = E_c^1(I_c)$  and  $e_s^1 = E_s^1(I_s)$ ;
20  Calculating  $z_0^g = \mathcal{E}(I_g)$  and  $z_0^c = \mathcal{E}(I_c)$ ;
21  Sampling  $t' \in (t_1, t_2]$  and calculating  $\tilde{z}_{t'}$  by Eq. (9);
22  Using conditional prediction network  $\tilde{z}_{\theta_1}^{(g,1)}(\tilde{z}_{t'}, t', y_1)$  to predict  $z_0^g$ ;
23  Sampling  $t'' \in (0, t_1]$  and using the predicted  $z_0^g$  to calculate  $\tilde{z}_{t''}$  by Eq. (8);
24  Calculating condition  $y_2 = \mathcal{C}(\mathcal{F}(e_c^2), \mathcal{F}(e_s^2))$  by flattening and concatening  $e_c^2 = E_c^2(I_c)$  and  $e_s^2 = E_s^2(I_s)$ ;
25  Using conditional prediction network  $\tilde{z}_{\theta_2}^{(g,2)}(\tilde{z}_{t''}, t'', y_2)$  to predict  $z_0^g$ ;
26  Using loss function Eq. (23) to fine-tune  $E_c^2, E_s^2, \tilde{z}_{\theta_2}^{(g,2)}(\tilde{z}_{t''}, t'', \mathbf{y}_2)$  in an end-to-end manner;
27  Until Converged;

```

---

## 7.3. Additional Details of Implementation and Optimization

We utilize the widely-used Stable Diffusion platform [3] to construct our proposed diffusion process in latent space. In the training stage, we obtain the latent feature  $z_0^c \in \mathcal{R}^{16 \times 16 \times 4}$  and  $z_0^g \in \mathcal{R}^{16 \times 16 \times 4}$  by using the VAE encoder  $\mathcal{E}$  to project the source image  $I_c \in \mathcal{R}^{128 \times 128 \times 3}$  and the target image  $I_g \in \mathcal{R}^{128 \times 128 \times 3}$  into latent space, respectively. In the inference stage, the prediction network  $\tilde{z}_{\theta_i}^{(g,i)}(\tilde{z}_t, t, y_i)$  will generate  $z_0^g \in \mathcal{R}^{16 \times 16 \times 4}$  by gradually denoising  $\tilde{z}_t$ , and the generated font image  $I_G$  can be obtained by using the VAE decoder  $\mathcal{D}$  to project  $z_0^g$  back to the image space. Following common practice in FFG, we construct the content encoder  $E_c^i$  and the style encoder  $E_s^i$  to generate the condition  $y_i$ . The style encoder  $E_s^i$  and the content encoder  $E_c^i$  have the same architecture, including 16 convolution layers and 3 down-sampling layers. We utilize the content encoder  $E_c^i$  and the style encoder  $E_s^i$  to extract the content feature  $e_c^i \in \mathcal{R}^{16 \times 16 \times 1024}$  from the source image  $I_c$  and the style feature  $e_s^i \in \mathcal{R}^{16 \times 16 \times 1024}$  from the reference images  $I_s$ , respectively. Finally, the content and style features are flattened and concatenated to build the condition  $y_i \in \mathcal{R}^{512 \times 1024}$ . Since the dimension of the condition  $y_i$  is 1024, we

slightly modify the cross-attention module in  $\tilde{z}_{\theta_i}^{(g,i)}(\tilde{z}_t, t, y_i)$  to perform attention operation on the dimension 1024.

In this paper, we keep the VAE encoder  $\mathcal{E}$  and decoder  $\mathcal{D}$  frozen and train the remaining networks. Specifically, we utilize the loss function  $L_g^i$  (in Eq. (22) and (23)) to optimize the prediction network  $\tilde{z}_{\theta_i}^{(g,i)}(\tilde{z}_t, t, y_i)$  with the corresponding content encoder  $E_c^i$  and style encoder  $E_s^i$ . The prediction network  $\tilde{z}_{\theta_i}^{(g,i)}(\tilde{z}_t, t, y_i)$  is initialized by the pre-trained Stable Diffusion models, while the style encoder  $E_s^i$  and content encoder  $E_c^i$  are initialized by Kaiming initialization. The AdamW optimizer with the learning rate  $lr = 3.2 \times 10^{-6}$  and batch size 16 is used to train our model.

## 8. Additional Experimental Results

### 8.1. Ablation Study on the Multi-Stage Font Generative Process

The selections of the region for different stages in the multi-stage generative process are hyper-parameters in our model, which significantly impact image quality. Therefore, we conduct experiments on UFUC dataset to study the influence on model performance by selecting different  $t_1$  and  $t_2$ . As shown in Tab. 3, the configuration  $t_1 = 200$  and  $t_2 = 800$  achieves the best record in terms of RMSE, PSNR and SSIM. Moreover, comparing the results in the first two rows (No. 1 and No. 2), the performance of  $t_1 = 100$  is slightly lower than  $t_1 = 200$ , since the font refinement stage of  $t_1 = 100$  is too short to modify the image appearances and local details. Furthermore, for the experiments with the fixed  $t_1$  and different  $t_2$  (e.g., No. 2, 4, 5), we find that the larger font transfer stage tends to generate better results, as the network is easy to model the slow transformation process. However, when  $t_2$  is too large (e.g.,  $t_2 = 900$  in No. 6), the structure construction stage cannot provide enough global structure information for the generative process.

Table 3. Ablation study for the regions of different stages on UFUC dataset. The bold number indicates the best.

No.	$t_1$	$t_2$	RMSE $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
1	100	800	0.2488	12.42	0.7131	0.1553
2	200	800	<b>0.2475</b>	<b>12.45</b>	<b>0.7189</b>	0.1527
3	300	800	0.2493	12.39	0.7155	0.1530
4	200	600	0.2602	12.03	0.7035	0.1708
5	200	700	0.2521	12.29	0.7138	0.1544
6	200	900	0.2496	12.40	0.7175	<b>0.1524</b>

Therefore, the hyper-parameters  $t_1$  and  $t_2$  are trade-off selections, and we select the configuration  $t_1 = 200$  and  $t_2 = 800$  in the following experiments.

### 8.2. Additional Qualitative Results

In this section, we provide more qualitative results of our proposed model.

Content	강 갑 맘 바 뱠 네 꼬 당 담 남 달 당 떼 암 훠 낭 뉘 놔 항 화
MX-Font	<b>강 갑 맘 바 뱠</b> 네 꼬 담 담 남 달 달 떼 암 훠 낭 뉘 놔 항 화
NTF	강 갑 맘 바 뱠 네 꼬 담 담 남 달 달 떼 암 훠 낭 뉘 놔 항 화
Ours	<b>강 갑 맘 바 뱠</b> 네 꼬 당 담 남 달 달 떼 암 훠 낭 뉘 놔 항 화
Reference Image	刻 刻 刻 刻 刻 刻 刻 刻 刻 刻 刻 刻 刻 刻 刻

Figure 6. Qualitative comparisons of our proposed model with other state-of-the-art methods on cross-lingual (Chinese to Korean) generation. Note that we use eight reference images in inference, and we only display one reference image in this figure.

路米棒城持抽垂典否富根供海和厚呼皇街理留  
期群伸始事束徒武星衣英用育枝知植柱宗族作  
推味引用招枝植州柱族最破泉群弱善舌社失盛  
社伸保倍成持抽告供海何荷呼皇街禁利米取群  
作取若社盛素堂庭徒星序崖宇章枝柱字宗最昨  
秒伸安棒倍本朝成答低副富根供利料流路米去  
徒推味序穴洋役用枝植宗族保倍唱朝成抽堂投  
皇理去舌社伸盛始事授寺推宇在章枝知潮赤抽  
若伸始世束寺推招枝植柱宗族保倍兵唱持抽垂  
群事安岸保倍朝城持放峰服荷街禁精路秒破期  
利流期泉肉舌伸始室辛星要意英真宗昨荷湖皇  
街破期保倍唱朝持抽答服富根供固和湖加禁甲  
持若伸私堂徒要永育章真置柱族作保倍常潮城  
取似授投棒潮城抽垂答度服根供禁精留路破期  
始事室安岸材常城抽打峰富根荷加精理情社伸  
路若性字在柱宗朝典否富供固合荷呼街理利秒  
字落努情特投徒味先役意永由早招真枝知置自  
推味倍博唱城持垂典放服根供呼寄美破情舌伸  
伸世倍常朝抽垂度放服富根荷街理流路秒破社  
情洋安岸保倍曾潮抽峰根供荷湖街精理路美破

Figure 7. Additional qualitative results of our proposed model on UFUC dataset.

Content	萤讳凭菌钻春草惑带闯卦爪阴具功进丢描堡盲吧权香丢刃
LF-Font	萤讳凭菌钻春草惑带闯卦爪阴具功进丢描堡盲吧权香丢刃
DG-Font	萤讳凭菌钻春草惑带闯卦爪阴具功进丢描堡盲吧权香丢刃
MX-Font	萤讳凭菌钻春草惑带闯卦爪阴具功进丢描堡盲吧权香丢刃
FsFont	萤讳凭菌钻春草惑带闯卦爪阴具功进丢描堡盲吧权香丢刃
NTF	萤讳凭菌钻春草惑带闯卦爪阴具功进丢描堡盲吧权香丢刃
Ours	萤讳凭菌钻春草惑带闯卦爪阴具功进丢描堡盲吧权香丢刃
Ground Truth	萤讳凭菌钻春草惑带闯卦爪阴具功进丢描堡盲吧权香丢刃

Figure 8. Qualitative comparisons of our proposed model with other state-of-the-art methods on UFSC dataset.

使梧虾现详醒杏朽鞍掺春铛妇卦惯柜角俱食割  
热艳堡虫带袋搞柜酣贺疾沮络蝶媚钮凭瓶欺权  
胎望胃梧杏喧硬鞍芭虫割卦哼嫁接届竞眼砌窍  
届俏珊猾徙廂醒渲艳渔缘彰睁处纲柜讳郊接俱  
钠枪俏琼权冰春挫搞柜酣虎艰接届俱距溃盟描  
砌珊鞍蚕齿处春妇纲搞惯柜届进俱菌络盟描钮

Figure 9. Additional qualitative results of our proposed model on UFSC dataset.

澡演壁蜻脉蝠漂糕模褂模踪障塞踩僧辗摘躺镇  
碳熬避演藻蜘蛛踪暴臂缤漱播潮撑滴蝠缴鲸漂蜻  
蝇榜壁播缠潮瞅憧醋摧瞪滴噩蝠犒衡境磊潘慎  
漂影癌幢障踩蔡操犒噩稿横塞蝴蝶精颗蜡磷潘稳

Figure 10. Additional qualitative results of our proposed model for handwritten fonts in UFUC setting.

## References

- [1] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in neural information processing systems, 33:6840–6851, 2020. 3, 5, 12.
- [2] Calvin Luo. Understanding diffusion models: A unified perspective. arXiv preprint arXiv:2208.11970, 2022. 13.
- [3] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 10684–10695, 2022. 3, 4, 5, 6, 14