# RewriteNet: Realistic Scene Text Image Generation via Editing Text in Real-world Image

Junyeop Lee[1*], Yoonsik Kim[2*]

Seonghyeon Kim[2], Moonbin Yim[2], Seung Shin[2], Gayoung Lee[2], Sungrae Park[1†]

[1]Upstage AI Research, Upstage
[2]Clova AI Research, NAVER Corp.
{junyeop.lee, sungrae.park}@upstage.ai
{yoonsik.kim90, kim.seonghyeon, moonbin.yim, seung.shin, gayoung.lee}@navercorp.com

## Abstract

*Scene text editing (STE), which converts a text in a scene image into the desired text while preserving an original style, is a challenging task due to a complex intervention between text and style. To address this challenge, we propose a novel representational learning-based STE model, referred to as RewriteNet that employs textual information as well as visual information. We assume that the scene text image can be decomposed into content and style features where the former represents the text information and style represents scene text characteristics such as font, alignment, and background. Under this assumption, we propose a method to separately encode content and style features of the input image by introducing the scene text recognizer that is trained by text information. Then, a text-edited image is generated by combining the style feature from the original image and the content feature from the target text. Unlike previous works that are only able to use synthetic images in the training phase, we also exploit real-world images by proposing a self-supervised training scheme, which bridges the domain gap between synthetic and real data. Our experiments demonstrate that RewriteNet achieves better quantitative and qualitative performance than other comparisons. Moreover, we validate that the use of text information and the self-supervised training scheme improves text switching performance. The implementation and dataset will be publicly available.*

## 1. Introduction

Scene text editing (STE) is a task of image synthesis that replaces the text in a scene image to the desired text while preserving a style such as a font type, font size, text alignment, and background. As a core technology for virtual reality, STE can be employed to replace the text contents and visualize the result of machine translation (e.g. Figure 1 (a) and (b)). Moreover, as shown in Figure 1 (c), STE also can be applied to augment the training data by converting text into various target texts for scene text recognition (STR) [1, 29, 45] and detection [2, 48, 19]. STE still has not been widely adopted, because it is intricately intertwined with various tasks such as image in-painting, style extraction, character rendering, and localization.

To generate scene text edited images, previous STE methods [32, 42, 44] generally follow a framework with two stages: text deletion and text conversion. Text deletion module generates text erased background, which can be thought of as an image in-painting task specialized for scene text images [27, 46, 43, 38]. Text conversion module renders the desired text where the text-related styles in the original image are transferred, and then, two outputs generated from text deletion and conversion module are harmonically fused. These methods show successful performance, but we think they have two limitations. These methods heavily depend on visual features when distinguishing between the text region and background region and do not fully utilize the text information that could be helpful for understanding and generating scene text images. Moreover, due to the requirement of additional supervisions for the deletion, real-world data cannot be employed as a training data and thus results in the model being biased in synthetic styles.

From these observations, we present a novel repre-

---

(a) Original scene images     (b) Edited scene images
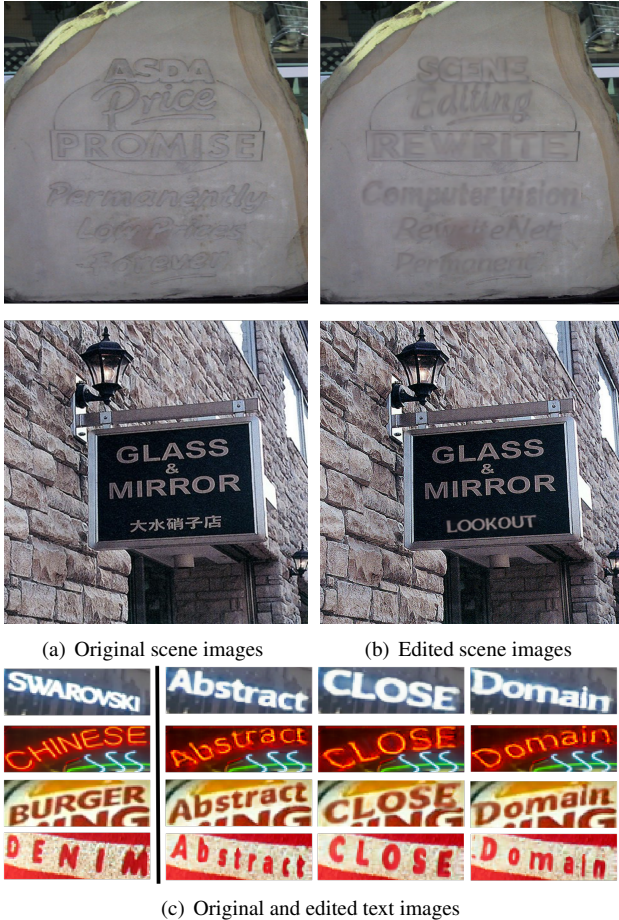
(c) Original and edited text images

Figure 1. (a) and (b) are original scene images and the results of RewriteNet respectively. (c) is the original text images (leftmost) and the text edited images where target texts are "Abstract", "CLOSE", and "Domain", respectively.

sentational learning-based STE framework, referred to as RewriteNet, which employs textual and visual information. The employment of text information is motivated by the human perception that recognizing text in the image can be helpful to identify the text region from the complex background. We assume that the scene text image can be decomposed into content and style features that respectively indicate text information and anything other than text information such as font, text alignment, and background. Under this assumption, we propose a method to separately represent content and style features of the input image by introducing the STR architecture [1] that purposes to recognize text in a scene text image. With separately extracted style and content features, a generator can be trained to synthesize an image with a target text while preserving the style of the style-image. Thus, RewriteNet replaces the text deletion and conversion stages of previous work with a simple encoder in the latent space and the model can be trained in

an end-to-end manner.

We also propose a self-supervised training scheme that does not require additional annotation cost and enables to exploit unlabeled real-world images. The proposed self-supervised training scheme prevents the trained model to be biased in synthetic styles and bridges the domain gap between training and test environments. As shown in Figure 1, our model robustly generates text-edited images where the styles of original images are well preserved.

Our contributions can be summarized as follows:

- We propose a novel representational learning-based STE framework that employs textual information to separately encode content features from style features. To the best of our knowledge, our RewriteNet is the first method that exploits textual information for STE.

- We also propose the self-supervised training scheme for STE model to employ unlabeled real scene images.

- We confirm that RewriteNet outperforms previous image-to-image translation methods and recent STE methods on qualitative and quantitative comparisons.

- We broadly demonstrate the validations of RewriteNet such as feature decomposition, model design, robustness, and extensibility on data augmentation for STR.

## 2. Related Works

### 2.1. Scene Text Editing

As the growth of generation model [49, 50], STE has been actively studied for its various applications. Previous STE methods mainly have proposed multiple sub-modules to extract a background and spatial text alignment, and a single fusion module to generate a text-edited image with the identified information. Specifically, initial STE work [32] segments binary mask for each character and switches it into the desired character. Although it has shown that their character correction method can be applied in real-world images, it cannot deal with different lengths between the text in the original image and the desired text. Moreover, its simple rule-based segmentation module could critically affect the generation performance.

Recently, Wu *et al*. [42] and Yang *et al*. [44] have proposed word-level STE methods using text conversion, background in-painting and fusion modules. These methods attempt to train the model to separate the text region and the background region using the text-erased image. They could successfully conduct word-level STE, however, sometimes they fail to edit the text of the complex style images. We think that it is because these methods solely rely on visual features to distinguish the text area of the input image. To deal with this problem, we introduce the novel architecture which can understand text information in the input image
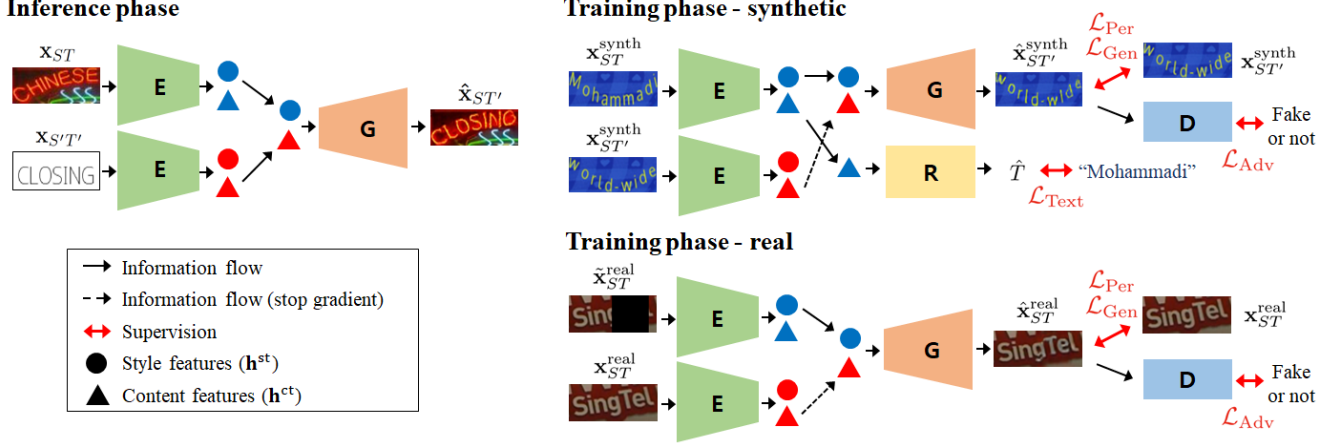
Figure 2. Overview of RewriteNet during inference and training process. The RewriteNet is composed of Style-content encoder ($\mathbf{E}$), Generator ($\mathbf{G}$), Text recognizer ($\mathbf{R}$), and Discriminator ($\mathbf{D}$). In each phase, we use the same encoder to extract style and content features from two different images. The output image is generated by combining the style feature extracted from the top image (style-image) and the content feature extracted from the bottom image (content-image).

and separate content and style features from the image. Furthermore, we propose the self-supervised training scheme that enables to use the unlabeled real-world images. The use of real-world images helps the model to cope with complex and various scene images where the synthetic rendering system cannot represent.

## 2.2. Image-to-Image Translation

Image-to-image translation methods have been widely researched due to their practical usages. Isola *et al.* [13] have proposed paired image-to-image translation with conditional GAN to learn a mapping from the input image to the output image. To address unpaired image-to-image translation, UNIT [20] and MUNIT [11] assume fully and partially shared latent space, respectively. Ma *et al.* [23] have proposed an exemplar guided image translation method with semantic feature mask that does not require additional labels for feature mask. Motivated by previous works, we introduce partially shared latent space assumption and the self-supervised training scheme with STE specialized proposals.

## 3. RewriteNet

RewriteNet consists of an encoder that extracts decomposed style and content features, and a generator that generates a text image with the identified features. This section first describes how encoder and generator are utilized to generate an image of desired text, and explains how the modules are trained on synthetic and real datasets. Followed by it, architectural details are provided.

### 3.1. Inference Process

Let $\mathbf{x}_{ST}$ be the style-image with text $T$ and style $S$. When a target text $T'$ is given, our model aims to generate $\mathbf{x}_{ST'}$ whose text is switched into the target text $T'$ from $\mathbf{x}_{ST}$ while holding its style $S$. To achieve the goal, RewriteNet assumes two disentangled latent features, $\mathbf{h}_S^{\mathrm{st}}$ for the style $S$ and $\mathbf{h}_T^{\mathrm{ct}}$ for the content $T$, and conducts the content switch.

Following the encoder and generator framework, the inference model consists of the below two modules.

- **Style-content encoder** ($\mathbf{E}^{\mathrm{st}}: \mathbf{x}_{ST} \to \mathbf{h}_S^{\mathrm{st}}$, $\mathbf{E}^{\mathrm{ct}}: \mathbf{x}_{ST} \to \mathbf{h}_T^{\mathrm{ct}}$) extracts latent style feature $\mathbf{h}_S^{\mathrm{st}}$ and content feature $\mathbf{h}_T^{\mathrm{ct}}$ from an input image $\mathbf{x}_{ST}$.

- **Generator** ($\mathbf{G}: \mathbf{h}_S^{\mathrm{st}}, \mathbf{h}_T^{\mathrm{ct}} \to \hat{\mathbf{x}}_{ST}$) generates an output image of text $T$ under the style $S$.

By switching off the latent content features, the model becomes enabled to generate a text-switched image $\hat{\mathbf{x}}_{ST'}$ as follows:

$$\hat{\mathbf{x}}_{ST'} = \mathbf{G}(\mathbf{E}^{\mathrm{st}}(\mathbf{x}_{ST}), \mathbf{E}^{\mathrm{ct}}(\mathbf{x}_{S'T'})). \tag{1}$$

The left of Figure 2 explains the inference process in a view of the information flow. A content-image $\mathbf{x}_{S'T'}^{\mathrm{synth}}$ is synthetically rendered with simple style $S'$ and target text $T'$.

### 3.2. Training Process

The right of Figure 2 shows two training processes of our model. One is for paired synthetic images, and the other is for unpaired real-world images.

#### 3.2.1 Modules Utilized in Training Process

Here, we introduce two modules only used in the training process to encourage the content-switched image genera-

tion.

- **Text recognizer** ($\mathbf{R}$: $\mathbf{h}_T^{\text{ct}} \rightarrow T$) identifies text label from the latent content feature. By learning content features $\mathbf{h}^{\text{ct}}$ to predict text label, the content feature can represent the text upon on the input image and is used as a content condition of $\mathbf{G}$. We should note that content feature $\mathbf{h}^{\text{ct}}$ is only trained with text label in the whole training process.

- **Style-content discriminator** ($\mathbf{D}$: $\hat{\mathbf{x}}_{ST'}, \mathbf{x}_{ST}, \mathbf{h}_{T'}^{\text{ct}} \rightarrow [0, 1]$) determines whether an input image $\hat{\mathbf{x}}_{ST'}$ is synthetically generated with a style reference $\mathbf{x}_{ST}$ and a content feature $\mathbf{h}_{T'}^{\text{ct}}$, where $\hat{x}_{ST'}$ is an output of $\mathbf{G}$. As a competitor of $\mathbf{G}$, its adversarial training improves generation quality.

By utilizing these two modules, the $\mathbf{E}$ enables to identify the latent content and the $\mathbf{G}$ enables to generate high-quality images.

### 3.2.2 Learning from Synthetic Data

We train the modules to decompose style and content features by using synthetic image pairs. As shown in the top right of Figure 2, synthetic image pairs share the same style but have different text contents. The content feature is learned to capture text information in an image by utilizing $\mathbf{E}^{\text{ct}}$ and $\mathbf{R}$. The encoded content feature is fed into the recognizer, and to let the recognizer predict correct labels, the encoder is trained to produce favorable content features. The style feature is learned to represent style information by allowing $\mathbf{E}^{\text{st}}$ and $\mathbf{G}$ to maintain style consistency after content switched generation.

We can obtain paired images $\{\mathbf{x}_{ST}^{\text{synth}}, \mathbf{x}_{ST'}^{\text{synth}}\}$ by feeding different texts to synthesizing engine with same rendering parameters such as background, font style, alignment, and so on. Then, a single training set becomes $\{\mathbf{x}_{ST}^{\text{synth}}, \mathbf{x}_{ST'}^{\text{synth}}, T\}$ where $T$ is a text label. Therefore, $\mathbf{E}$, $\mathbf{G}$ and $\mathbf{R}$ can be trained with reconstruction loss and recognition loss as follows:

$$\mathcal{L}_{\text{Gen}}^{\text{synth}} = \|\mathbf{G}(\mathbf{E}^{\text{st}}(\mathbf{x}_{ST}^{\text{synth}}), \bar{\mathbf{E}}^{\text{ct}}(\mathbf{x}_{ST'}^{\text{synth}})) - \mathbf{x}_{ST'}^{\text{synth}}\|_1, \quad (2)$$

$$\mathcal{L}_{\text{Text}}^{\text{synth}} = \sum_i \text{CrossEntropy}(\mathbf{R}(\mathbf{E}^{\text{ct}}(\mathbf{x}_{ST}^{\text{synth}}))_i, T_i), \quad (3)$$

where $\bar{\mathbf{E}}^{\text{ct}}$ indicates a frozen encoder that does not get any back-propagation flow and $T_i$ represents $i$-th character of the ground truth text label. If we do not freeze $\mathbf{E}^{\text{ct}}$ at reconstruction loss, $\mathbf{E}$ and $\mathbf{G}$ will quickly fall into a local minimum by simply copying content-image. Thus, we freeze the $\mathbf{E}^{\text{ct}}$ to prevent the content feature from being affected by the reconstruction loss and train $\mathbf{E}^{\text{ct}}$ only with the recognition loss. These losses guide the model to learn the content switch, but the trained model might fail to address

real-world images caused by the limitation of the synthetic styles.

### 3.2.3 Learning from Real-world Data

Synthetic data provides proper guidance for content switching, but it does not fully represent a style of real-world images. To compensate for the limitation, we propose a self-supervised training process for RewriteNet utilizing real-world data as shown in the right bottom of Figure 2.

In the case of real-world images, there are no paired images that have different texts with the same style. Moreover, it is expensive to obtain text labels of real-world images. Therefore, we introduce conditioned denoising autoencoder loss to allow the model to learn style and content representations for unpaired real-world images. Specifically, we cut out a region randomly selected in the width direction with length $w$ to lose some characters [5], and then the noisy image is used as a style-image to extract the style feature from the left regions. By combining the content feature extracted from the original image, $\mathbf{G}$ fills the blank by referring to the style of the surrounding area of the blank region. The proposed self-supervised scheme will forbid the model to trivially autoencode style-image by using the corrupted image as style-image and enforce model to learn separated representations. The denoising autoencoder loss is defined as:

$$\mathcal{L}_{\text{Gen}}^{\text{real}} = \|\mathbf{G}(\mathbf{E}^{\text{st}}(\tilde{\mathbf{x}}_{ST}^{\text{real}}), \bar{\mathbf{E}}^{\text{ct}}(\mathbf{x}_{ST}^{\text{real}})) - \mathbf{x}_{ST}^{\text{real}}\|_1, \quad (4)$$

where $\tilde{\mathbf{x}}_{ST}^{\text{real}}$ indicates a noisy image corrupted from $\mathbf{x}_{ST}^{\text{real}}$. Here, we should note that the proposed self-supervised method does not require any text labels and paired images.

### 3.2.4 Adversarial Training

Generally, text image in the wild has high-frequency regions like complex background, diverse textures, and high contrast regions. However, pixel-wise reconstruction loss, referred to as $\mathcal{L}_{\text{Gen}}$, has a limitation to address the high-frequency and tends to capture the low-frequencies [13]. To encourage high-frequency crispness, we apply the generative adversarial network (GAN) framework to generate realistic text images [18, 13, 25, 24]. Specifically, we design the $\mathbf{D}$ to represent a fake or real probability of the input image under the given conditions of its style-image and latent content. We denote $\mathbf{D}(X|X^{\text{st}}, \mathbf{h}^{\text{ct}})$ for the probability $p(X \text{ is not fake}|X^{\text{st}}, \mathbf{h}^{\text{ct}})$, where $X$ and $X^{\text{st}}$ indicate the input image and the style-image respectively. The adversarial losses are calculated as follows:

$$\mathcal{L}_{\text{Adv}}^{\text{synth}} = \log \mathbf{D}(\mathbf{x}_{ST'}^{\text{synth}}|\mathbf{x}_{ST}^{\text{synth}}, \bar{\mathbf{E}}^{\text{ct}}(\mathbf{x}_{ST'}^{\text{synth}}))$$
$$+ \log\left(1 - \mathbf{D}(\hat{\mathbf{x}}_{ST'}^{\text{synth}}|\mathbf{x}_{ST}^{\text{synth}}, \bar{\mathbf{E}}^{\text{ct}}(\mathbf{x}_{ST'}^{\text{synth}}))\right), \quad (5)$$

$$\mathcal{L}_{\text{Adv}}^{\text{real}} = \log \mathbf{D}(\mathbf{x}_{ST}^{\text{real}}|\tilde{\mathbf{x}}_{ST}^{\text{real}}, \bar{\mathbf{E}}^{\text{ct}}(\mathbf{x}_{ST}^{\text{real}}))$$
$$+ \log\left(1 - \mathbf{D}(\hat{\mathbf{x}}_{ST}^{\text{real}}|\tilde{\mathbf{x}}_{ST}^{\text{real}}, \bar{\mathbf{E}}^{\text{ct}}(\mathbf{x}_{ST}^{\text{real}}))\right), \quad (6)$$

where $\hat{\mathbf{x}}_{ST'}^{\text{synth}}$ and $\hat{\mathbf{x}}_{ST}^{\text{real}}$ denote generated images from synthetic and real-world style-images, respectively (See Figure 2). Here, it should be noted that the latent contents used as the conditions are frozen to block back-propagation flow to the $\mathbf{E}$ from the adversarial loss.

We also employ feature matching loss that stabilizes the training of various GAN models [33, 41, 21]. Specifically, we extract intermediate feature maps of the $\mathbf{D}$ and minimize the distance between generated and target samples:

$$\mathcal{L}_{\text{Per}}^{\text{synth}} = \sum_l \frac{1}{M_l} \|\phi_l(\mathbf{x}_{ST'}^{\text{synth}}) - \phi_l(\hat{\mathbf{x}}_{ST'}^{\text{synth}})\|_1, \qquad (7)$$

$$\mathcal{L}_{\text{Per}}^{\text{real}} = \sum_l \frac{1}{M_l} \|\phi_l(\mathbf{x}_{ST}^{\text{real}}) - \phi_l(\hat{\mathbf{x}}_{ST}^{\text{real}})\|_1, \qquad (8)$$

where $\phi_l$ and $M_l$ are the output feature map and its size of the $l$-th layer. For each loss, the same conditions are used to calculate the activation maps $\{\mathbf{x}_{ST}^{\text{synth}}, \bar{\mathbf{E}}^{\text{ct}}(\mathbf{x}_{ST'}^{\text{synth}})\}$ for $\mathcal{L}_{\text{Per}}^{\text{synth}}$ and $\{\tilde{\mathbf{x}}_{ST}^{\text{real}}, \bar{\mathbf{E}}^{\text{ct}}(\mathbf{x}_{ST}^{\text{real}})\}$ for $\mathcal{L}_{\text{Per}}^{\text{real}}$. Feature matching losses could facilitate $\mathbf{G}$ to match multi-scale statistics with target samples [41], thus beneficial for overall sample qualities.

#### 3.2.5 Final Loss Term

The final losses are formalized as follows:

$$\mathcal{L} = \underset{\mathbf{E}, \mathbf{G}, \mathbf{R}}{\operatorname{argmin}} \Big( \mathcal{L}_{\text{Gen}}^{\text{synth}} + \alpha \mathcal{L}_{\text{Text}}^{\text{synth}} + \gamma \mathcal{L}_{\text{Per}}^{\text{synth}}$$
$$+ \lambda \left( \mathcal{L}_{\text{Gen}}^{\text{real}} + \gamma \mathcal{L}_{\text{Per}}^{\text{real}} \right) \qquad (9)$$
$$+ \beta \underset{\mathbf{D}}{\operatorname{argmax}} \left( \mathcal{L}_{\text{Adv}}^{\text{synth}} + \lambda \mathcal{L}_{\text{Adv}}^{\text{real}} \right) \Big),$$

where $\alpha$, $\beta$, $\gamma$, and $\lambda$ are hyper-parameters balancing the losses.

### 3.3. Architectural Details

**Style-content Encoder** The style-content encoder follows *partially shared latent space assumption* as in MUNIT [12], where an image $\mathbf{x}_{ST}$ is composed of its latent style feature $\mathbf{h}_S^{\text{st}}$ and content feature $\mathbf{h}_T^{\text{ct}}$. The network is based on a ResNet [7] similar to the feature extractor used in [3]. In addition, we apply bidirectional LSTM [9] layers upon the content features to alleviate spatial dependencies from the input image.

**Text recognizer** Text Recognizer estimates a sequence of characters in an image and it has an important role to distinguish contents from styles. It consists of an LSTM decoder with an attention mechanism [1] from the identified content features. Since the text labels are required to train this module, we only train the module with the synthetic dataset.

**Generator** Given the latent style and content features as an input, the generator outputs an image with a given style and content. The generator network is similar to the decoder used in the Unet [31] architecture. The style features in multiple $\mathbf{E}^{\text{st}}$ layers are fed into the generator using short-connections. The network design is inspired by the Pix2Pix [13] model.

**Style-content discriminator** The style-content discriminator determines whether an image is fake or not. The network architecture follows PatchGAN [13, 36].

## 4. Experiments

This section describes the details about datasets and implementation of RewriteNet and validates the generation performance of RewriteNet.

### 4.1. Datasets and Implementation Details

#### 4.1.1 Synthetic Data for Training

Since RewriteNet requires paired synthetic datasets for supervised-learning, we generate 8M synthetic images with our synthesizing engine[1] that is based on MJSynth [14] and SynthText [6]. Specifically, we compose paired data $\{\mathbf{x}_{ST}^{\text{synth}}, \mathbf{x}_{ST'}^{\text{synth}}\}$ with same rendering parameters $S$ such as font styles, background textures, shape for the text alignments (rotation, perspective, curve), and artificial blur noises except only for input texts $(T, T')$. The employed texts are the union of MJSynth [14] and SynthText [6] corpus and the paired synthetic dataset will be publicly available.

#### 4.1.2 Real Data for Training and Evaluation

We combine multiple benchmark training datasets such as IIIT [26], IC13 [16], IC15 [15], and COCO [39]. The total number of training images is 59,856. Although these datasets contain ground-truth text labels, our model does not employ the text labels that are expensive in a practical scenario. For evaluation, we use the test a split of each public dataset such as IIIT [26], IC03 [22], IC13 [16], IC15 [15], SVT [40], SVTP [28] and CT80 [30] where the total number of test images is 8,536.

#### 4.1.3 Implementation Details

We rescale the input image to $32 \times 128$ and empirically set $w$ as 42, which is the proper length to cut out some characters and capture style information for the self-supervised training. To balance the multiple loss terms, we empirically set $\alpha = 20$, $\beta = 10$, $\gamma = 1$, $\delta = 0.1$, and $\lambda = 1$. The

---

[1]Our engine can generate various rendered images by controlling multiple parameters.

Figure 3. Visual comparisons on generation performance.

model is optimized by Adam optimizer [17] with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. A cyclic learning rate [37] is applied with an initial learning rate of 1e-4 and an maximum iteration number of 300K. The batch size is 192 including 144 for synthetic data and 48 for real-world data. The total training takes 7 days using two Tesla V100s. At the inference phase, $\mathbf{x}_{S'T'}^{\text{synth}}$ is generated by $ImageDraw$ function from $PIL$ package.

### 4.2. Comparison on Generation Performance

We compare our model to four models: MUNIT [12], EGSC [23], ETIW [42], and STEFANN [32]. Although MUNIT and EGSC are not specifically designed for STE task, we train the model targeted for STE task and make a comparison with our model to validate the STE performance of the representative image translation models[23]. ETIW is the exact comparison method for RewriteNet and its results are achieved from re-implementation[4]. STEFANN is designed for the character-wise correction method that requires manual text region selection, so test environments are different from other methods. We try to achieve high-quality results for STEFANN by testing multiple times with its official code[5].

In the quantitative comparison, we employ three measurements: (1) recognition accuracy on generated images utilizing a pre-trained STR model[6], (2) learned perceptual image patch similarity (LPIPS) [47], and (3) Fréchet Inception Distance (FID) [8]. The recognition accuracy measures whether the generated images contain switched contents or not. The LPIPS and FID represent style consistency between a style-image and a generated image. Here, we

---

[2]https://github.com/NVlabs/MUNIT

[3]https://github.com/charliememory/EGSC-IT

[4]https://github.com/youdao-ai/SRNet

[5]https://github.com/prasunroy/stefann

[6]https://github.com/clovaai/deep-text-recognition-benchmark

| Models | Accuracy ($\uparrow$) | LPIPS ($\downarrow$) | FID ($\downarrow$) |
|---|---|---|---|
| MUNIT [12] | 80.75 | 0.33 | 65.7 |
| EGSC [23] | 0.04 | 0.19 | 43.3 |
| ETIW [42] | 31.16 | **0.12** | **13.7** |
| Ours w/o real | <u>89.00</u> | <u>0.18</u> | 18.7 |
| Ours | **90.30** | 0.19 | <u>16.7</u> |

Table 1. Quantitative comparison between scene text generation methods: "Accuracy" represents the content-switch performance (higher is better) and "LPIPS" and "FID" show style consistency (lower is better). The bold and underline indicate the best method and the second-best method, respectively.

would note that measurements between text switching performance and style preserving performance have a trade-off. It is because the best performance on LPIPS and FID is achieved when the output is the same as the input. Thus, balanced performance and visual results should be considered to compare the model performance. Table 1 presents the quantitative comparison results. Quantitative performance on STEFANN is not evaluated, because it requires manual region selection for each image and it considerably takes a long time.

The naive application of MUNIT tends not to maintain original styles, which can be confirmed in its high LPIPS and FID scores. We observe that the naive application of EGSC would inappropriate for scene text editing. ETIW shows the best performance on LPIPS and FID, however, it achieves comparably lower accuracy. These results indicate ETIW often fails to convert the content and simply copies style-image where the examples are shown in Figure 3 (2nd and 6th rows). Since other methods could not employ unlabeled real images as a training set, we report the performance of two of our models: (1) only trained with synthetic images and (2) trained with synthetic and unlabeled real im-
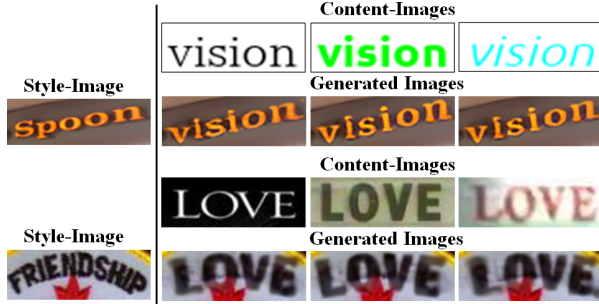
Figure 4. Generated images with diverse content-images that have different styles with same content.



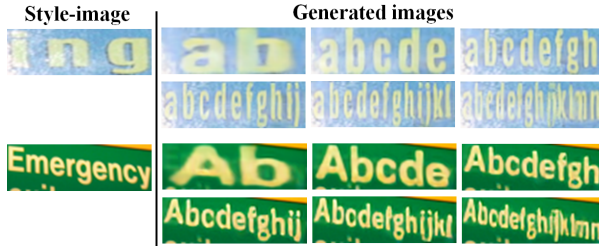Figure 5. Generated images from our RewriteNet when the lengths of desired texts are different.



Figure 6. Generated images from our RewriteNet trained with ablated training processes. Target texts are "exposes", "golf", "changed", "cottage" and "chocolate", respectively.



Figure 7. Generated images from our RewriteNet trained with validations of design choices. Target texts are "board", "world", "paramount", and "GAS", respectively.

| Models | Accuracy (↑) | LPIPS (↓) | FID (↓) |
|---|---|---|---|
| w/o $\mathbf{R}$ | 0.54 | 0.29 | 114.9 |
| w/o real | 89.00 | 0.18 | 18.7 |
| w/ real w/o noise | 82.97 | **0.17** | 20.9 |
| Proposed | **90.30** | 0.19 | **16.7** |

Table 2. Ablation study about the training process on absence of recognizer and unpaired real-world data.

| Models | Accuracy (↑) | LPIPS (↓) | FID (↓) |
|---|---|---|---|
| w/o Stop Gradient | **97.22** | 0.36 | 89.9 |
| w/$\mathcal{L}_{\text{con-text}}$ | 94.34 | **0.19** | **16.7** |
| Proposed | 90.30 | **0.19** | **16.7** |

Table 3. Validations of stop gradient and consistency loss. "Proposed" indicates w/ Stop Gradient and w/o $\mathcal{L}_{\text{con-text}}$.

ages. The proposed model, trained only with synthetic data, achieves the best accuracy among the models trained with synthetic data and also shows comparable performance on LPIPS and FID. The use of a real dataset improves the text editing performance, but it represents different tendencies on the performance of style preservation.

The visual comparisons are presented in Figure 3. MU-NIT looks failed to preserve the style of style-image and ETIW tends to simply copy style-image for challenging style without content switching. STEFANN also cannot robustly edit texts when the lengths of texts are different and backgrounds are complex. In contrast, the proposed methods show promising results on multiple examples compared to other methods. The use of real image finely upgrade visual quality. More visual results are presented in Figure 1 where multi texts also can be edited by employing text region detector [2].

## 4.3. Discussions

### 4.3.1 Decomposition between Content and Style

To validate whether our model successfully separates the content feature from the style feature, we show that the style of content-image does not affect the style of the generated image. We feed various images for content-images that have different styles with the same content, and observe whether the generated results are affected. As shown in Figure 4, the generated results are quite stable to the change of content-images. This result suggests that our model well separates the content feature from the style feature of input images.

### 4.3.2 Robustness for Different Text Lengths

To show the robustness of text editing for different text lengths between desired text and style-image, we present more generation examples when the length of the desired text are extremely different with the text of style-image. As can be seen in Figure 5, RewriteNet can edit different length texts robustly. Specifically, RewriteNet generates great quality outputs (1st example) when converting the 3 characters (ing) to 14 characters (abcdefghijklmn). Interestingly, we observe that the model can properly adjust the height, width, and spacing of characters as the number of characters changed.

| Model | Train Data | Regular | | | | Irregular | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | IC13 | SVT | IIIT | Average | IC15 | CT80 | SVTP | Average |
| BEST [1] | Synth | 92.8 | **88.7** | 91.9 | 91.1 | 78.0 | 76.7 | 79.5 | 78.1 |
| BEST [1] | Synth+MUNIT | 89.2 | 76.2 | 88.8 | 84.7 (-6.4) | 62.0 | 57.3 | 66.2 | 61.8 (-16.3) |
| BEST [1] | Synth+ETIW | 87.2 | 77.1 | 84.6 | 83.0 (-8.1) | 64.0 | 62.8 | 61.1 | 62.6 (-15.5) |
| BEST [1] | Synth+Ours | **93.2** | **88.7** | **92.8** | 91.6 (+0.5) | **79.6** | **84.4** | **81.6** | 81.8 (+3.7) |

Table 4. STR accuracies over six benchmark test datasets depending on the training data. "Synth" indicates font-rendered data from MJSynth [14] and SynthText [6]. "MUNIT", "ETIW", and "Ours" represent fully generated data from unlabeled real images using MUNIT, ETIW, and RewriteNet, respectively.

### 4.3.3 Ablation Studies on Main Proposals

We describe the effectiveness of employing the recognizer that is designed for feature decomposition and unpaired real-world images for training sets. To reveal that, we train RewriteNet with ablated training processes: a model without the recognizer (w/o $\mathbf{R}$), a model without the training branch utilizing real-world data (w/o real), and a model feeding an original real-world image into the training branch instead of its noisy variant (w/ real w/o input noise). Table 2 and Figure 6 show the comparison results. We observe that RewriteNet cannot be trained without $\mathbf{R}$ where the performance of switching text is dramatically degraded. The visual results also present the necessity of $\mathbf{R}$. The performance of "w/o real" achieves lower accuracy than the proposed method and training with real images present better visual results when the texts are irregularly shaped and background are complex. The model trained without noises shows the lower accuracy but with the best scores on the LPIPS. These quantitative performance and its visual results suggest that it has strong tendencies to simply autoencoding style-images, rather than edit text with given contents when the style is challenging.

### 4.3.4 Validation on Design Choices

In the training of RewriteNet, we have frozen $\mathbf{E}^{ct}$ at the generation process to prevent simple copying of content-image. We conduct the ablation study when the $\mathbf{E}^{ct}$ is not frozen. We also validate the employment of consistency loss, because it is widely adopted in generation tasks [49, 11, 10, 4] and improves performance by regularizing the generator. Following the previous works, we train RewriteNet with additional consistency loss as follows:

$$\mathcal{L}_{\text{con-text}}^{\text{synth}} = \sum_i \text{CrossEntropy}(\mathbf{R}(\mathbf{E}^{ct}(\hat{\mathbf{x}}_{ST}^{\text{synth}}))_i, T_i), \quad (10)$$

where $\hat{x}_{ST}^{\text{synth}}$ denotes the generated image with the style $S$ and content $T$, and $T_i$ represents $i$-th character of the ground truth text label. This loss re-enforces the generated image to have desired text and it can only be applied on the synthetic data due to the requirement of text label.

The quantitative performance of "w/o Stop Gradient" (without frozen) and "w/ $\mathcal{L}_{\text{con-text}}$" are reported in Table 3. "w/o Stop Gradient" achieves the higher accuracy than the proposed, however, the performances of LPIPS and FID were much worse than the proposed. On the other hand, "w/ $\mathcal{L}_{\text{con-text}}$" improves accuracy without significant performance degradation of LPIPS and FIDS. We also present the generated images in Figure 7. "w/o Stop Gradient" simply writes the desired texts without preserving styles. As can be seen in second row of Figure 7, "w/ $\mathcal{L}_{\text{con-text}}$" sometimes erases text that is out of interest in the background (part of characters are erased in the top right region), for enhancing the recognition accuracy on the generated image. Moreover, it fails to preserve font information when the text shapes and textures are complex.

### 4.3.5 Scene Text Editing for Scene Text Recognition

Since RewriteNet can provide reliable performance on text-editing, we extend RewriteNet to generate training samples for STR models. To analyze the effect of the dataset generated by our model and other comparisons on scene text recognition, we train BEST [1] with the generated and synthetic images. Then, we investigate the improvements in the recognition performance. For the generation, the unified real-world data (59,856 in total), combining four benchmark training datasets such as IIIT [26], IC13 [16], IC15 [15], and COCO [39], is used as the style-images. We generate 18 text images from a single style image. As a result, the total amount of generated images is 1M (59,856 × 18). Here, we would note that STE methods do not require additional text labels for generating samples. Following the evaluation protocol of STR [1], all trained models are evaluated on the six real-word benchmarks where total number of images is 7,406; 3,000 from IIIT [26], 647 from SVT [40], 1,015 from IC13 [16], 1,811 from IC15 [15], 645 from SVTP [28], and 288 from CT80 [30].

In Table 4, we observe other comparison methods including STE method are harmful to train STR model. These performance degradation might result from noise labels where models cannot reliably edit texts and make a mismatch between images and labels. On the other hand, the proposed

RewriteNet contributes to the performance improvement on both regular and irregular benchmarks. Interestingly, irregular benchmarks are significantly improved and these improvements represent that our generated data mainly complements insufficient irregularly shaped data by providing realistic images.

## 5. Conclusions

In this paper, we have proposed RewriteNet which could edit text content in a scene text image. The novel feature decomposition method that makes use of an STR network successfully extracts content and style features and show that unlabeled real-world images can be further utilized by adopting cutout. Compared to previously proposed STE and image translation methods, the outputs generated by RewriteNet achieve better generation quality. Based on our experiments, we expect that the proposed RewriteNet can be applied to practical usages.

## References

[1] Jeonghun Baek, Geewook Kim, Junyeop Lee, Sungrae Park, Dongyoon Han, Sangdoo Yun, Seong Joon Oh, and Hwalsuk Lee. What is wrong with scene text recognition model comparisons? dataset and model analysis. In *International Conference on Computer Vision (ICCV)*, 2019. 1, 2, 5, 8, 12

[2] Youngmin Baek, Bado Lee, Dongyoon Han, Sangdoo Yun, and Hwalsuk Lee. Character region awareness for text detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9365–9374, 2019. 1, 7

[3] Zhanzhan Cheng, Fan Bai, Yunlu Xu, Gang Zheng, Shiliang Pu, and Shuigeng Zhou. Focusing attention: Towards accurate text recognition in natural images. In *ICCV*, pages 5086–5094, 2017. 5

[4] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8789–8797, 2018. 8

[5] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017. 4

[6] Ankush Gupta, Andrea Vedaldi, and Andrew Zisserman. Synthetic data for text localisation in natural images. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 5, 8, 12

[7] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 5

[8] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 6629–6640, 2017. 6

[9] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997. 5

[10] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *International conference on machine learning*, pages 1989–1998. PMLR, 2018. 8

[11] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 172–189, 2018. 3, 8

[12] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 172–189, 2018. 5, 6

[13] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. 3, 4, 5

[14] Max Jaderberg, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Reading text in the wild with convolutional neural networks. *International Journal of Computer Vision*, 116(1):1–20, 2016. 5, 8, 12

[15] Dimosthenis Karatzas, Lluis Gomez-Bigorda, Anguelos Nicolaou, Suman Ghosh, Andrew Bagdanov, Masakazu Iwamura, Jiri Matas, Lukas Neumann, Vijay Ramaseshan Chandrasekhar, Shijian Lu, et al. Icdar 2015 competition on robust reading. In *ICDAR*, pages 1156–1160, 2015. 5, 8, 12

[16] Dimosthenis Karatzas, Faisal Shafait, Seiichi Uchida, Masakazu Iwamura, Lluis Gomez i Bigorda, Sergi Robles Mestre, Joan Mas, David Fernandez Mota, Jon Almazan Almazan, and Lluis Pere De Las Heras. Icdar 2013 robust reading competition. In *ICDAR*, pages 1484–1493, 2013. 5, 8, 12

[17] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. 6

[18] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, page 1558–1566, 2016. 4

[19] Minghui Liao, Zhaoyi Wan, Cong Yao, Kai Chen, and Xiang Bai. Real-time scene text detection with differentiable binarization. In *Proc. AAAI*, 2020. 1

[20] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *Advances in neural information processing systems*, pages 700–708, 2017. 3

[21] Ming-Yu Liu, Xun Huang, Arun Mallya, Tero Karras, Timo Aila, Jaakko Lehtinen, and Jan Kautz. Few-shot unsupervised image-to-image translation. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. 5

[22] Simon M Lucas, Alex Panaretos, Luis Sosa, Anthony Tang, Shirley Wong, and Robert Young. Icdar 2003 robust reading competitions. In *ICDAR*, pages 682–687, 2003. 5

[23] Liqian Ma, Xu Jia, Stamatios Georgoulis, Tinne Tuytelaars, and Luc Van Gool. Exemplar guided unsupervised image-to-image translation with semantic consistency. In *International Conference on Learning Representations*, 2019. 3, 6

[24] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2794–2802, 2017. 4

[25] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014. 4

[26] Anand Mishra, Karteek Alahari, and CV Jawahar. Scene text recognition using higher order language priors. In *BMVC*, 2012. 5, 8, 12

[27] Toshiki Nakamura, Anna Zhu, Keiji Yanai, and Seiichi Uchida. Scene text eraser. In *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, volume 1, pages 832–837. IEEE, 2017. 1

[28] Trung Quy Phan, Palaiahnakote Shivakumara, Shangxuan Tian, and Chew Lim Tan. Recognizing text with perspective distortion in natural scenes. In *ICCV*, pages 569–576, 2013. 5, 8, 12

[29] Siyang Qin, Alessandro Bissacco, Michalis Raptis, Yasuhisa Fujii, and Ying Xiao. Towards unconstrained end-to-end text spotting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 1

[30] Anhar Risnumawan, Palaiahankote Shivakumara, Chee Seng Chan, and Chew Lim Tan. A robust arbitrary text detection system for natural scene images. In *ESWA*, volume 41, pages 8027–8048, 2014. 5, 8, 12

[31] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 5

[32] Prasun Roy, Saumik Bhattacharya, Subhankar Ghosh, and Umapada Pal. Stefann: Scene text editor using font adaptive neural network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 1, 2, 6

[33] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, Xi Chen, and Xi Chen. Improved techniques for training gans. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 2234–2242. Curran Associates, Inc., 2016. 5

[34] Baoguang Shi, Xiang Bai, and Cong Yao. An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. In *TPAMI*, volume 39, pages 2298–2304, 2017. 12

[35] Baoguang Shi, Xinggang Wang, Pengyuan Lyu, Cong Yao, and Xiang Bai. Robust scene text recognition with automatic rectification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4168–4176, 2016. 12

[36] Ashish Shrivastava, Tomas Pfister, Oncel Tuzel, Joshua Susskind, Wenda Wang, and Russell Webb. Learning from simulated and unsupervised images through adversarial training. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2107–2116, 2017. 5

[37] Leslie N Smith. Cyclical learning rates for training neural networks. In *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 464–472. IEEE, 2017. 6

[38] Osman Tursun, Rui Zeng, Simon Denman, Sabesan Sivapalan, Sridha Sridharan, and Clinton Fookes. Mtrnet: A generic scene text eraser. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, pages 39–44. IEEE, 2019. 1

[39] Andreas Veit, Tomas Matera, Lukas Neumann, Jiri Matas, and Serge Belongie. Coco-text: Dataset and benchmark for text detection and recognition in natural images. In *arXiv:1601.07140*, 2016. 5, 8

[40] Kai Wang, Boris Babenko, and Serge Belongie. End-to-end scene text recognition. In *ICCV*, pages 1457–1464, 2011. 5, 8, 12

[41] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 5

[42] Liang Wu, Chengquan Zhang, Jiaming Liu, Junyu Han, Jingtuo Liu, Errui Ding, and Xiang Bai. Editing text in the wild. In *Proceedings of the 27th ACM International Conference on Multimedia*, page 1500–1508, 2019. 1, 2, 6

[43] Wei Xiong, Jiahui Yu, Zhe Lin, Jimei Yang, Xin Lu, Connelly Barnes, and Jiebo Luo. Foreground-aware image inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5840–5848, 2019. 1

[44] Qiangpeng Yang, Jun Huang, and Wei Lin. Swaptext: Image based texts transfer in scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 1, 2

[45] Deli Yu, Xuan Li, Chengquan Zhang, J. Han, Jingtuo Liu, and E. Ding. Towards accurate scene text recognition with semantic reasoning networks. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12110–12119, 2020. 1

[46] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Generative image inpainting with contextual attention. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5505–5514, 2018. 1

[47] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 586–595, 2018. 6

[48] X. Zhou, C. Yao, H. Wen, Y. Wang, S. Zhou, W. He, and J. Liang. East: An efficient and accurate scene text detector. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2642–2651, 2017. 1

[49] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017. 2, 8

[50] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. In *Advances in neural information processing systems*, pages 465–476, 2017. 2

# Appendices

## A. Experiments on Feature Decomposition

We have validated feature decomposition between style and content in Figure 4 (main manuscript) where the generated images are quite stable to the change of content images. We investigate which component mainly affects the style of the generated image by feeding various fonts and colors of content-image. As can be seen in Figure 8, the generated images are stable according to the change of font but slightly different from each other. On the other hand, the generated images are invariant according to the change of color as shown in Figure 9. We also measure distances between the generated images and variance of the generated images with three metrics:

- PSNR (Peak Signal-to-Noise Ratio): pixel-wise MSE (Mean Squared Error) based distance. Since we cannot achieve the original (reference) image, we measure the distance between generated images.

- SSIM (Structural Similarity): perceptual quality based distance. Since we cannot achieve the original (reference) image, we measure the distance between generated images.

- Variance: the averaged variance of the generated images.

As can be seen in Table 5, the stability of the generated images is more affected by the change of font.

| Changes | PSNR ($\uparrow$) | SSIM ($\uparrow$) | Var ($\downarrow$) |
|---------|------|------|-----|
| Font | 18.51 | 0.6242 | $7.9 \times 10^{-3}$ |
| Color | $Inf$ | 0.9942 | $0.08 \times 10^{-3}$ |

Table 5. PSNR (dB) and SSIM between generated images and the variance of generated images. $Inf$ denotes infinity that indicates some of the generated images are exactly identical.

## B. Validations on Training Data for STR Models

We have presented in Table 4 of the main manuscript that data obtained from the output of RewriteNet improves the scene text recognizer (STR) performance. In this subsection, we validate the effects of our generated data on different STR models such as CRNN [34], RARE [35] and BEST [1]. As presented in Table 6, our generated data improves STR performances on all baselines. In particular, irregular test sets are more improved than regular test sets where averaged irregular and regular improvements are 0.6pp and 3.9pp respectively. These improvements show that our generated data mainly complements insufficient irregularly shaped data by providing realistic images.

| Model | Train Data | Regular | Irregular |
|-------|-----------|---------|-----------|
| CRNN [34] | Synth | 87.0 | 68.9 |
| CRNN [34] | Synth+Ours | **87.3** (+0.3) | **73.3** (+4.4) |
| RARE [35] | Synth | 89.0 | 74.9 |
| RARE [35] | Synth+Ours | **90.0** (+1.0) | **78.4** (+3.5) |
| BEST [1] | Synth | 91.1 | 78.1 |
| BEST [1] | Synth+Ours | **91.6** (+0.5) | **81.8** (+3.7) |

Table 6. Average STR accuracy of regular [16, 40, 26] and irregular [15, 30, 28] shaped benchmarks on three STR models. "Synth" indicates font-rendered data from MJSynth [14] and SynthText [6]. "Ours" represents fully generated data from unlabeled real images using RewriteNet.

Figure 8. The generated image from RewriteNet with changing the font of content-image. The black rectangular image represents content-image and its corresponding output is listed below.



Figure 9. The generated image from RewriteNet with changing the color of content-image. The black rectangular image represents content-image and its corresponding output is listed below.