# FAST: Font-Agnostic Scene Text Editing

**Alloy Das[1], Prasun Roy[2], Saumik Bhattacharya[3], Subhankar Ghosh[2], Umapada Pal[1], and Michael Blumenstein[2]**

[1]Computer Vision and Pattern Recognition Unit, Indian Statistical Institute, Kolkata, India

[2] Faculty of Engineering and IT, University of Technology Sydney, NSW, Australia

[3] Electronics and Electrical Communication Engg., Indian Institute of Technology Kharagpur, Kharagpur, India

`alloyuit@gmail.com, prasun.roy@student.uts.edu.au, saumik@ece.iitkgp.ac.in`
`subhankar.ghosh@student.uts.edu.au, umapada@isical.ac.in, michael.blumenstein@uts.edu.au`

## ABSTRACT

Scene Text Editing (STE) is a challenging research problem, and it aims to modify existing texts in an image while preserving the background and the font style of the original text of the image. Due to its various real-life applications, researchers have explored several approaches toward STE in recent years. However, most of the existing STE methods show inferior editing performance because of (1) complex image backgrounds, (2) various font styles, and (3) varying word lengths within the text. To address such inferior editing performance issues, in this paper, we propose a novel font-agnostic scene text editing framework, named FAST, for simultaneously generating text in arbitrary styles and locations while preserving a natural and realistic appearance through combined mask generation and style transfer. The proposed approach differs from the existing methods as they directly modify all image pixels. Instead, the proposed method has introduced a filtering mechanism to remove background distractions, allowing the network to focus solely on the text regions where editing is required. Additionally, a text-style transfer module has been designed to mitigate the challenges posed by varying word lengths. Extensive experiments and ablations have been conducted, and the results demonstrate that the proposed method outperforms the existing methods both qualitatively and quantitatively.

## 1 Introduction

In recent years, there has been a growing interest in the field of scene text editing due to its numerous practical applications across various domains. These include text image synthesis [1, 2, 3], styled text transfer [4, 5, 6], and augmented reality translation [7, 8, 9].

The task of modifying or inserting text into an image while preserving its natural and realistic appearance has been extensively researched [10, 11, 12]. Previous approaches to this task have predominantly formulated it as a style transfer problem using generative models like GANs [13, 14, 15, 16, 12, 17]. These methods require a reference image that serves as a template for the desired style, typically a cropped section of the image containing the target text. The approach involves rendering the desired text in the desired spelling to match the reference style and background. However, these methods have limitations in their ability to generate text in largely varying font styles, sizes, and colors with arbitrary geometrical transforms. Additionally, the existing process of cropping, transferring style, and then replacing it often results in unnatural-looking results with discordance between the edited and surrounding areas, featuring distinct boundaries and messy distortions.

| Source Image | Generated Image | Source Image | Generated Image | Source Image | Generated Image |



Figure 1: Examples of scene text editing with FAST on real-world image samples.

To address these limitations, in this paper, we propose a novel GAN-based style transfer module that generates the target mask from the input source mask using an image-to-image transfer module. Our method aims to generate text in arbitrary styles and at arbitrary locations while maintaining a natural and realistic appearance. By generating the mask and transferring the style simultaneously, our method is able to produce more accurate and natural-looking results. In Fig1, we have shown the result on a Real database[15]. The main contributions of this paper are as follows:

- Unlike most of the existing pipelines that perform STE at a character level, the proposed algorithm performs the editing at the word level. This ensures faster editing with less distortion in the background.

- The algorithm is able to generate target texts that may have different character-length compared to the source text.

- The proposed algorithm is largely font agnostic and can be used for real-world scene text editing, irrespective of the font size, color, or orientation.

The rest of the paper is organized as follows. We explore the existing STE methods and prerequisites of STE in Section 2. The proposed approach is described in Section 3. We discuss the datasets, implementation details, evaluation metrics, comparisons with previous methods, and ablation study in Section 4. Section 5 discusses the limitations of the proposed method. Finally, we conclude the paper with an overview of our approach and potential future scopes.

## 2 Related Works

**Text Image Synthesis:** Image synthesis and rendering have been intensively researched in the realm of computer graphics [18]. Text image synthesis is regarded as a data augmentation technique in the context of text identification and detection. For instance, while Gupta et al. [19] design a powerful engine to produce synthetic text images, Jaderberg et al. [20]. employ a word generator to create synthetic word images. The goal of text image synthesis is to insert texts in areas of a background image that have semantically significant content. However, the realism of synthesized text images is affected by a number of parameters, including font size, perspective, and illumination. In order to solve this problem, Zhan et al. [3] integrate semantic coherence, visual attention, and adaptive text presentation. The synthesized images are visually accurate, but they differ dramatically from actual photographs in a number of ways, such as the fact that the text and backdrop images have fewer font options.

To get over these restrictions, recent studies have looked into GAN-based picture synthesis algorithms. In order to achieve synthesis realism in both geometry and appearance spaces, Zhan et al. [21] present a spatial fusion GAN that combines a geometry synthesizer with an appearance synthesizer. The stylistic degree of glyphs can be controlled by Yang et al.'s [22] framework using an adjustable parameter, whereas GA-DAN [23] simultaneously models cross-domain shifts in both geometry and appearance spaces. Additionally, MC-GAN [6] is introduced for font style transfer in the set of letters from A to Z, and Wu et al. [12]. offer an end-to-end trainable style retention network to edit the text in natural photos.

**Style Transfer:** The difficult task of transferring an image's visual style from a reference image to a target image is known as image style transfer. Many existing techniques embed the input into a subspace using an encoder-decoder architecture, then decode it to produce the desired images. Isola et al. [24] implement a learnable mapping from input to output photos, for instance, using a training set of aligned image pairings. Zhu et al. [25] introduce cycle consistency loss to generalize the mapping relationship to unpaired cross-domain data. Various challenges, like creating photographs from sketches and synthesizing faces from attribute and semantic data, have been approached using similar concepts. In [26], authors proposed an algorithm to transfer text effects using statistical measurements. Here, distance-based properties of text impacts are analyzed, modeled, and used to direct the synthesis process. [27] uses stylization and de-stylization sub-networks to accomplish both the goals of style transfer and removal.

**Scene Text Editing:** The variety of uses for GAN-based scene text editing has piqued researchers' growing curiosity. To alter a single character, for instance, [16] created the Font Adaptive Neural Network (STEFANN). This character-level alteration, however, falls short of replacing words with length alterations, which restricts practical uses. In [12], the authors developed the word-level editing approach employing three sub-networks: background inpainting, text conversion, and fusion to overcome this constraint. Each module can manage a reasonably straightforward task thanks to the divide-and-conquer method. In order to make the text conversion module easier to learn, [17] enhanced SRNet and added the TPS module, which separates the spatial transformation from text styles. Furthermore, [28] suggested the forge-and-recapture procedure to reduce visual artifacts and applied scene text editing to document forging. In order to keep the consistency of all other edited frames, [29] used SRNet for video text editing, altering only the chosen frame as a reference and applying certain photometric modifications. These methods, which are essentially extensions of SRNet, can only be trained on artificial datasets and may not be able to replace text with complicated styles. A stroke-level
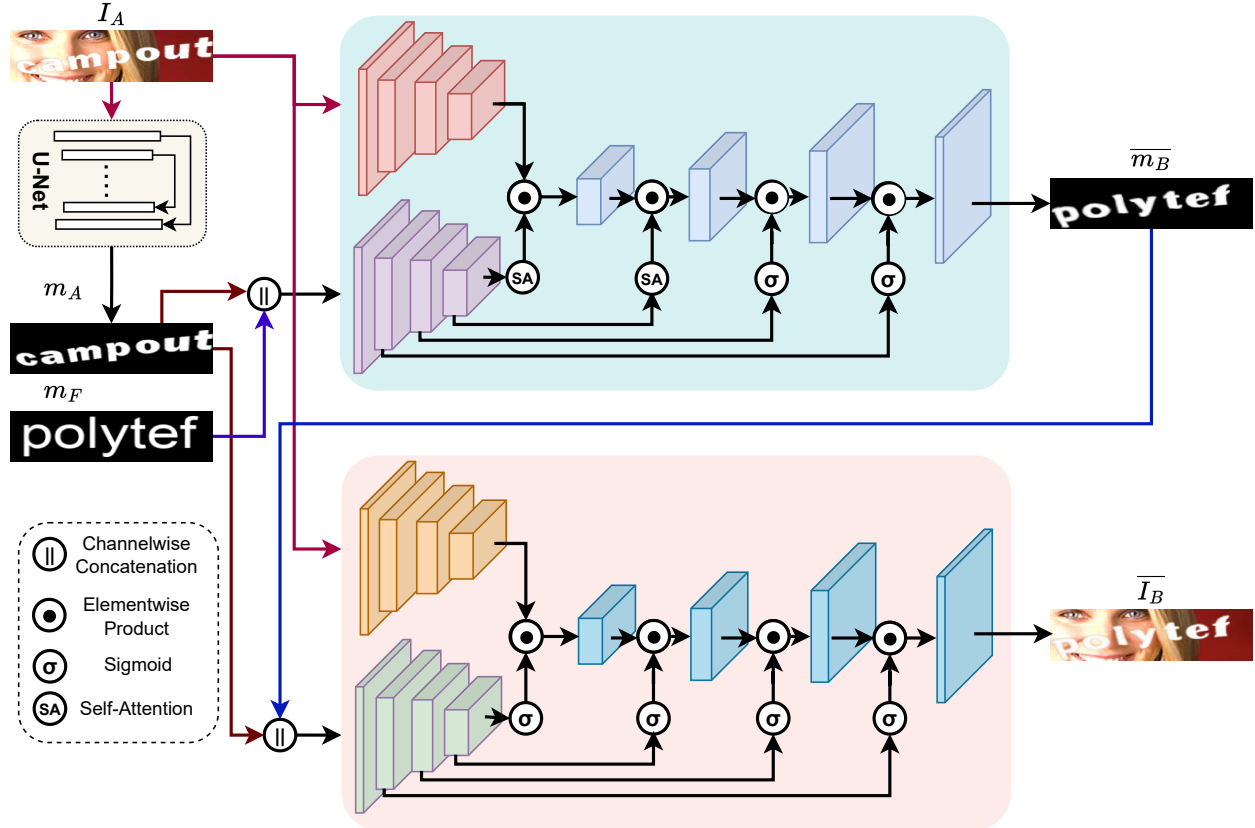
Figure 2: Architecture of the proposed method. In **stage-I**, an approximate target style mask $\overline{m_B}$ is estimated from the source image $I_A$, source style mask $m_A$, and a fixed style mask $m_F$ of the target text. In **stage-II**, the target image $\overline{I_B}$ is generated by transferring image attributes from $I_A$ and conditioning the image translation on the structural guidance $(m_A, \overline{m_B})$.

alteration technique that creates more readable text graphics are suggested by [15]. Their technique can be trained on both labeled synthetic datasets and unpaired scene text images, and it supports the semi-supervised training scheme.

## 3  Methodology

Given a scene text image $I_A$, the objective of the proposed STE is to generate an image $I_B$ with a modified text. To enforce a classifier-free image translation, we aim to condition the generative process on the structural guidance $(m_A, m_B)$, where $m_A$ and $m_B$ correspond to the binary masks of text content in $I_A$ and $I_B$, respectively. However, obtaining $m_B$ before generating $I_B$ is unrealistic, making the end-to-end text style transfer difficult. We address this issue by splitting the generative process into two independent stages. In the first stage, we replace the initially unknown mask $m_B$ with another mask $m_F$ having the same textual content but in a fixed font of known style. In this stage, the generator $G_m$ attempts to produce an approximation of the target style mask $\overline{m_B}$. In the second stage, an identical generator $G_i$ synthesizes the approximate target image $\overline{I_B}$ by transferring attributes from $I_A$ and using $(m_A, \overline{m_B})$ as the structural guidance. We use synthetically generated $(I_A, m_A)$ pairs to train both $G_m$ and $G_i$. However, an additional U-Net [30] is separately trained for estimating the mask $\overline{m_A}$ from $I_A$ during inference on real scene text samples. Fig. 2 shows an overview of our proposed architecture.

### 3.1   Stage – I: Target Style Mask Estimation

The stage – I architecture is a generative adversarial network (GAN) containing a target mask generator $G_m$ and a PatchGAN [24] discriminator $D_m$. $G_m$ takes the source image $I_A$ and the channel-wise concatenated masks $m_\theta = (m_A, m_F)$ as inputs and produces an approximate target style mask $\overline{m_B}$ as the output. $D_m$ discriminates

between a real and a fake transformation by taking channel-wise concatenated masks $(m_A, m_B)$ or $(m_A, \overline{m_B})$ and predicting a binary class probability map of ones (real) or zeroes (fake).

$G_m$ comprises two encoding branches for $I_A$ and $m_\theta$. The condition image and the guidance masks are resized to a dimension of $64 \times 256$. At each branch, the encoder first projects the input into a 64-channel feature space using $3 \times 3$ convolution kernels with stride 1, padding 1, and without adding any bias. The final block output is obtained after batch normalization and ReLU activation. The projected feature space is then downscaled four times. At every downsampling block, the spatial feature dimension is downscaled to half while doubling the number of channels. Each downsampling block uses $4 \times 4$ convolution kernels, stride 1, padding 1, and zero bias. Each downscaling convolution is immediately followed by batch normalization, ReLU activation and a basic residual block [31] to produce the final block output.

During decoding, the outputs of both encoding branches are channel-wise concatenated and passed through the decoder. The decoder consists of four upscaling blocks, each doubling the spatial feature dimension while decreasing the number of channels to half. The upscaling is performed with $4 \times 4$ transposed convolution kernels, stride 1, padding 1, and zero bias. Like the encoder, the final block output is generated following batch normalization, ReLU activation and a basic residual block.

In $G_m$, we apply two different attention mechanisms at matching feature resolutions of the encoder and decoder to attend to both coarse and fine image attributes during structural transformation. At the two lowest resolutions, we use *self-attention* similar to [32], and at the other two higher resolutions, we use *sigmoid attention*. Mathematically, at the two lowest resolutions $k = \{1, 2\}$ with self-attention,

$$m_1^{\phi_m} = \phi_{m1}(I_4^{\varphi_m} \odot SA(m_4^{\varphi_m}))$$

$$m_2^{\phi_m} = \phi_{m2}(m_1^{\phi_m} \odot SA(m_3^{\varphi_m}))$$

and for the following decoder blocks at higher resolutions $k = \{3, 4\}$ that use sigmoid attention,

$$m_k^{\phi_m} = \phi_{mk}(m_{k-1}^{\phi_m} \odot \sigma(m_{5-k}^{\varphi_m}))$$

where, $I_k^{\varphi_m}$ is the output of $k$-th image encoder block, $m_k^{\varphi_m}$ is the output of $k$-th mask encoder block, $m_k^{\phi_m}$ is the output of the $k$-th decoder block, $SA$ and $\sigma$ denote *self-attention* and *sigmoid attention*, respectively, and $\odot$ denotes element-wise product.

The output feature maps from the decoder are post-processed through four consecutive basic residual blocks. The resulting feature space is projected to a 3-channel image space of spatial resolution of $64 \times 256$ by a point convolution with $1 \times 1$ kernel, unit stride, zero padding, and without bias. The final normalized output of $G_m$ is obtained following a hyperbolic tangent (**tanh**) activation function.

## 3.2 Stage – II: Text Style Transfer with Structural Guidance

The stage – II architecture is identical to stage – I, with slightly different input specifications and attention mechanisms. In this case, the generator $G_i$ takes the source image $I_A$ and the channel-wise concatenated masks $(m_A, \overline{m_B})$ as inputs and produces an approximate target image $\overline{I_B}$ as the output. The PatchGAN discriminator $D_i$ discerns between a real and a fake transformation by taking channel-wise concatenated images $(I_A, I_B)$ or $(I_A, \overline{I_B})$ and predicting a binary class probability map of ones (real) or zeroes (fake).

In $G_i$, we apply only *sigmoid attention* at every matching feature resolution of the encoder and decoder. Mathematically, at the lowest resolution $k = 1$,

$$I_1^{\phi_i} = \phi_{i1}(I_4^{\varphi_i} \odot \sigma(m_4^{\varphi_i}))$$

and for the following decoder blocks at higher resolutions $k = \{2, 3, 4\}$,

$$I_k^{\phi_i} = \phi_{ik}(I_{k-1}^{\phi_i} \odot \sigma(m_{5-k}^{\varphi_i}))$$

where, $I_k^{\varphi_i}$ is the output of $k$-th image encoder block, $m_k^{\varphi_i}$ is the output of $k$-th mask encoder block, $I_k^{\phi_i}$ is the output of the $k$-th decoder block, $\sigma$ denotes *sigmoid attention*, and $\odot$ denotes element-wise product.

### 3.3 Learning Objectives

**Stage – I objectives:** The optimization objective of generator $G_m$ consists of four different loss components – (a) pixel-wise $L_2$ loss $\mathcal{L}_{L_2}^{G_m}$, (b) discriminator loss $\mathcal{L}_{GAN}^{G_m}$ estimated by the discriminator $D_m$, (c) perceptual loss $\mathcal{L}_{P_\rho}^{G_m}$ computed using a pre-trained VGG-19 network [33], and (d) multi-scale structural similarity [34] loss $\mathcal{L}_{SSIM}^{G_m}$.

The $L_2$ loss is estimated as the mean squared error between the generated and target masks as follows

$$\mathcal{L}_{L_2}^{G_m} = \|\overline{m_B} - m_B\|_2^2 \tag{1}$$

The discriminator loss is defined as the binary cross-entropy estimated by $D_m$ as follows

$$\mathcal{L}_{GAN}^{G_m} = \mathcal{L}_{BCE}\left(D_m(m_A, \overline{m_B}), 1\right) \tag{2}$$

The perceptual loss is defined as

$$\mathcal{L}_{P_\rho}^{G_m} = \frac{1}{h_\rho w_\rho c_\rho} \sum_{a=1}^{h_\rho} \sum_{b=1}^{w_\rho} \sum_{c=1}^{c_\rho} \|\psi_\rho(\overline{m_B}) - \psi_\rho(m_B)\|_1 \tag{3}$$

where $\psi_\rho$ is the output of dimension $(h_\rho \times w_\rho \times c_\rho)$ from the $\rho$-th layer of a pre-trained VGG-19 network and $\|\cdot\|_1$ denotes the $L_1$ norm (mean absolute error). We include two perceptual loss terms for $\rho = 4$ and $\rho = 9$ in the objective function.

The multi-scale structural similarity loss $\mathcal{L}_{SSIM}^{G_m}$ is estimated as

$$\mathcal{L}_{SSIM}^{G_m} = 1 - MSSSIM\left(\overline{m_B}, m_B\right) \tag{4}$$

where $MSSSIM(a, b)$ denotes the multi-scale structural similarity between $a$ and $b$.

Mathematically, the full objective function of $G_m$ is given by

$$\mathcal{L}_{G_m} = \lambda_1 \cdot \mathcal{L}_{L_2}^{G_m} + \lambda_2 \cdot \mathcal{L}_{GAN}^{G_m} + \lambda_3 \cdot (\mathcal{L}_{P_4}^{G_m} + \mathcal{L}_{P_9}^{G_m}) + \lambda_4 \cdot \mathcal{L}_{SSIM}^{G_m} \tag{5}$$

where $\lambda_1$, $\lambda_2$, $\lambda_3$, and $\lambda_4$ are the weighing parameters for respective loss components.

The optimization objective of the discriminator $D_m$ is given by

$$\mathcal{L}_{D_m} = \frac{1}{2}\left[\mathcal{L}_{BCE}(D_m(m_A, m_B), 1) + \mathcal{L}_{BCE}(D_m(m_A, \overline{m_B}), 0)\right] \tag{6}$$

**Stage – II objectives:** The optimization objective of generator $G_i$ consists of three loss components – (a) pixel-wise $L_1$ loss $\mathcal{L}_{L_1}^{G_i}$, (b) discriminator loss $\mathcal{L}_{GAN}^{G_i}$ estimated by the discriminator $D_i$, and (c) perceptual loss $\mathcal{L}_{P_\rho}^{G_i}$ computed using a pre-trained VGG-19 network. Mathematically,

$$\mathcal{L}_{L_2}^{G_i} = \left\|\overline{I_B} - I_B\right\|_1 \tag{7}$$

$$\mathcal{L}_{GAN}^{G_i} = \mathcal{L}_{BCE}\left(D_i(I_A, \overline{I_B}), 1\right) \tag{8}$$

$$\mathcal{L}_{P_\rho}^{G_i} = \frac{1}{h_\rho w_\rho c_\rho} \sum_{a=1}^{h_\rho} \sum_{b=1}^{w_\rho} \sum_{c=1}^{c_\rho} \left\|\psi_\rho(\overline{I_B}) - \psi_\rho(I_B)\right\|_1 \tag{9}$$

The full objective function of $G_i$ is given by

$$\mathcal{L}_{G_i} = \beta_1 \cdot \mathcal{L}_{L_2}^{G_i} + \beta_2 \cdot \mathcal{L}_{GAN}^{G_i} + \beta_3 \cdot (\mathcal{L}_{P_4}^{G_i} + \mathcal{L}_{P_9}^{G_i}) \tag{10}$$

where $\beta_1$, $\beta_2$, and $\beta_3$ are the weighing parameters for respective loss components.

The optimization objective of the discriminator $D_i$ is given by

$$\mathcal{L}_{D_i} = \frac{1}{2}\left[\mathcal{L}_{BCE}(D_i(I_A, I_B), 1) + \mathcal{L}_{BCE}(D_i(I_A, \overline{I_B}), 0)\right] \tag{11}$$

In our experiments, we use $\lambda_1 = 1$, $\lambda_2 = 5$, $\lambda_3 = 1$, $\lambda_4 = 100$, $\beta_1 = 5$, $\beta_2 = 1$, and $\beta_3 = 5$ as the weighing parameters. The values of these weighing parameters are estimated experimentally through an extensive ablation study.

Table 1: Quantitative comparison of different methods.

| Method | MSE ↓ | PSNR ↑ | SSIM ↑ | LPIPS (SqzNet) ↓ |
|--------|-------|--------|--------|------------------|
| pix2pix [24] | 0.0732 | 12.01 | 0.349 | - |
| SRNet [12] | 0.0193 | 18.66 | 0.610 | 0.2076 |
| SwapText [17] | 0.0174 | 19.43 | 0.652 | - |
| MOSTEL [15] | **0.0121** | **20.75** | 0.707 | 0.1770 |
| FAST (Proposed) | 0.0135 | 20.20 | **0.776** | **0.1556** |

## 4 Experiments

### 4.1 Datasets

**Synthetic Data:** For the supervised training of our pipeline, we utilized a dataset of 150,000 labeled images as mentioned in the work of Qu et al. [15]. For evaluation, we use the Tamper-Syn2k [15] dataset, which consists of 2,000 paired images specifically designed for evaluation purposes. These paired images were created by rendering different texts with consistent styles, including font, size, color, spatial transformation, and background image. To generate these paired images, a collection of 300 fonts and 12,000 background images are used. These background images were subjected to random rotation, curve, and perspective transformations to introduce diversity. It is worth noting that we also generated an additional set of 100,000 labeled images using 2,500 fonts. This expanded dataset was employed specifically for training purposes of the proposed method.

**Real Data:** In our research on enhancing natural screen text editing on real scene images, we utilized a dataset consisting of real scene images obtained from the MOSTEL [15]. This dataset is generated by the authors using random cropped images from ICDAR 2013 [35], MLT-2017 [36] MLT-2019 [37] datasets.

### 4.2 Implementation Details

The pre-transformation stage includes resizing the input images to $64 \times 256$. We utilize the Adam optimizer [38] with $\beta_1 = 0.5$ and $\beta_2 = 0.999$, and set the learning rate to $1 \times 10^{-3}$ for both the Stages. We train FAST for a total of 100k iterations using a batch size of 40. Our pipeline is implemented using PyTorch [39] and it is trained on a single NVIDIA 4080 OC GPU. Also we have trained it using the input image and its binary mask as ground truth. As a dataset we have used the same dataset as mentioned in the synthetic dataset. With batch size 1, for 20 epochs, with Adam optimizer and learning rate 0.001.

### 4.3 Evaluation Metrics

To evaluate the effectiveness of FAST in generating edited images, we employ the following commonly used metrics on paired synthetic images: (1) Mean Squared Error (MSE), which measures the L-2 distance; (2) Peak Signal-to-Noise Ratio (PSNR), which measures the ratio of peak signal to noise, and (3) Structural Similarity Index (SSIM) [34], which measures the mean structural similarity, and (4) Learned Perceptual Image Patch Similarity (LPIPS) [40], which measures the similarities of activation of two image patches of squzeenet [41]. Higher PSNR and SSIM, and lower MSE and LPIPS scores indicate better performance.

### 4.4 Comparisons with Previous Methods

In Table 1, we present a detailed comparison of our proposed method, FAST, with the recent existing approaches. We trained FAST on the same dataset used in the MOSTEL paper [15] and our synthetically generated data. We compared the performance based on the reported values in MOSTEL. Additionally, we calculated the LPIPS (SqzNet) scores for both MOSTEL [15] and SRNet [12]. The results show that our propose method (FAST) outperforms MOSTEL in terms of SSIM (Structural Similarity Index) and LPIPS (SqzNet) by approximately 0.05 and 0.02, respectively. This indicates that our approach generates images with higher similarity to the ground truth and exhibits better perceptual quality. In Figure 3, we provide visual comparisons of the generated images from different previous methods and FAST. Notably, the images generated by SRNet demonstrate difficulties in accurately capturing the desired features. While MOSTEL's generated images are of good quality, when compared directly with FAST, they are not as impressive. It is important to note that although MOSTEL shows better results in some cases compared to FAST, the visual quality and fidelity of the generated images produced by FAST surpass MOSTEL's performance.

In any generative task, human evaluation is important to understand the quality of the generation. Thus, we perform an
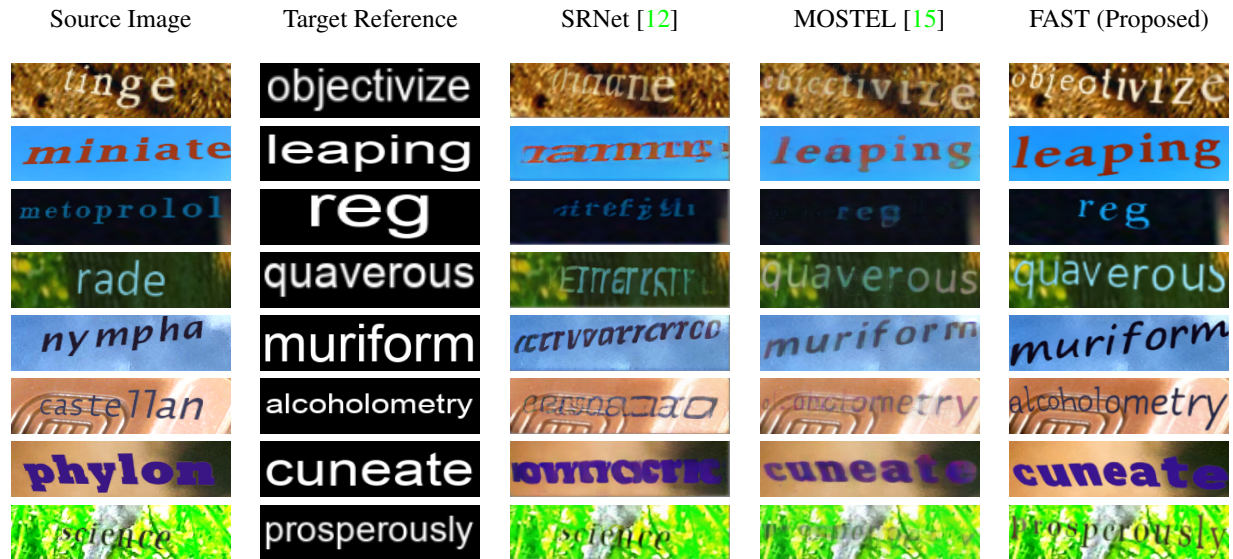
Figure 3: Qualitative comparison of FAST with the existing methods on synthetic datasets.

experiment to compare the mean opinion score of the proposed method with the SOTA technique MOSTEL as both the methods have close quantitative values. In our experiment, we randomly selected five images from each of the methods, MOSTEL and FAST (proposed). Additionally, we included 5 real images for comparison. These images were presented to a group of 47 individuals, and we collected a total of 705 responses. Each user is independently asked to mark the images as real (original) or fake (generated). Finally, we calculate how many images generated by a particular method has been identified as real. A higher value in this experiment typically indicates better generative performance. For the MOSTEL method, the rate was found to be 26%, indicating that 26% of the images generated by MOSTEL were perceived as real. On the other hand, for the FAST method, the rate was 45%, suggesting that 45% of the images generated by FAST were perceived as real. Overall, our method, FAST, showcases superior results both quantitatively and qualitatively when compared to the previous methods considered in the evaluation. In Figure 4, and Figure 5 we have shown the editing results in real data and synthetic data respectively.
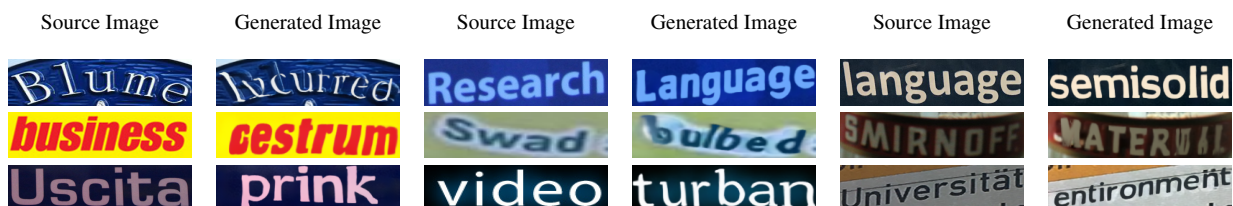


Figure 4: Typical examples of image editing using FAST on the Real dataset.



Figure 5: Typical examples of image editing using FAST on the Synthetic dataset.

## 4.5 Ablation Study

**Effectiveness of Self-Attention:**

We conducted a series of experiments in our FAST pipeline to investigate the impact of different attention mechanisms on image upsampling. In the first experiment, we employed four sigmoid attention blocks within the upsampling block.

Table 2: Ablation study – Ablation of different attention mechanisms.

| Block-wise attention combination | | | | | | |
|---|---|---|---|---|---|---|
| Block A | Block B | Block C | Block D | MSE ↓ | PSNR ↑ | SSIM ↑ |
| $\sigma$ | $\sigma$ | $\sigma$ | $\sigma$ | 0.0216 | 18.22 | 0.690 |
| SA | SA | SA | SA | 0.0213 | 18.00 | 0.672 |
| SA | SA | $\sigma$ | $\sigma$ | **0.0171** | **19.04** | **0.707** |

Table 3: Ablation study – Effectiveness of Data Mixing.

| Dataset | MSE ↓ | PSNR ↑ | SSIM ↑ |
|---|---|---|---|
| MOSTEL [15] | 0.0171 | 19.04 | 0.707 |
| Mixed | **0.0162** | **19.19** | **0.718** |

The objective was to assess the performance of sigmoid attention in isolation. For the second experiment, we integrated four self-attention blocks within the upsampling block. In the third experiment, we combined two sigmoid attention blocks and two self-attention blocks in the upsampling block. Notably, we adopted a strategy of utilizing self-attention at lower levels and sigmoid attention at higher levels. This arrangement was found to yield superior metrics compared to the previous experiments as shown in Table 2. Please note that these results are obtained by training the proposed model using the dataset mentioned in [15] only. Here we consider $I_A$ as the binary source mask of the textual part, and the channel-wise concatenated masks $m_\theta = (m_A, m_F)$ as another input.

**Effectiveness of Data Mixing:**

In the dataset ablation study, we employed two datasets. Initially, we utilized the MOSTEL dataset for training. Subsequently, we created a mixed dataset by combining the MOSTEL and SRNET data. We conducted two experiments and compared their results, as shown in Table 3. The findings indicate that the mixed data model outperforms the single dataset model. As mentioned above, here also we consider $I_A$ as the binary source mask of the textual part, and the channel-wise concatenated masks $m_\theta = (m_A, m_F)$ as another input.

**Effectiveness of Inputs:**

Our proposed pipeline was first guided by the concatenation of the input mask $m_A$ and a fixed mask $m_F$. Then, we ran an experiment to see what would happen if we only sent the fixed mask $m_F$ without concatenation. We also evaluated the efficacy of producing results from both binary source mask and image $I_A$. Table 4 illustrates our strategy's effectiveness.

## 5 Limitations

One limitation of the proposed STE described in the passage is the reliance on mask guidance maps to explicitly indicate editing regions. While this approach filters out background distractions and guides the network to focus on editing rules of text regions, it may encounter challenges when dealing with complex or irregular text layouts. The structure guidance maps assume that the text regions have well-defined masks or boundaries that can be easily segmented. However, in real-world images, text can appear in various forms, such as handwritten text, distorted text, or text with overlapping elements. In such cases, accurately generating mask guidance maps can be difficult, leading to potential errors or inconsistencies in the editing process. In Figure 6 we have shown some failure cases of our method.

Table 4: Ablation study – Effectiveness of Inputs.

| Input Combination | | | | |
|---|---|---|---|---|
| Input Image Type | Mask ($m_a$) Concatenation | MSE ↓ | PSNR ↑ | SSIM ↑ |
| Mask | Concat | 0.0162 | 19.19 | 0.718 |
| Mask | w/o Concat | 0.0178 | 18.86 | 0.701 |
| Image | Concat | **0.0135** | **20.20** | **0.776** |
| Image | w/o Concat | 0.0156 | 19.55 | 0.759 |

Figure 6: Some failure cases of our model.

## 6   Conclusion

This work proposes a font-agnostic scene text editing framework for simultaneously generating text in arbitrary styles and locations, while preserving a natural and realistic appearance through combined mask generation and style transfer. Extensive experiments conducted in the study demonstrate that the proposed method outperforms the existing methods in both qualitatively and quantitatively. By preserving the background and font style of the original text while modifying the text itself, the method holds promise for a wide range of real-world applications where text editing is required. However, it is important to acknowledge the limitations and challenges discussed earlier, such as the difficulty in accurately generating mask guidance maps for complex text layouts and the potential limitations in generalizing to all real-world scenarios. Overall, the proposed method represents a significant advancement in the field of STE, providing a more effective and robust approach to modifying scene text while preserving important visual aspects of the original image. Further research and refinement of these techniques can contribute to improving the accuracy and versatility of text editing in various practical applications.

## References

[1] Levent Karacan, Zeynep Akata, Aykut Erdem, and Erkut Erdem. Learning to generate images of outdoor scenes from attributes and semantic layouts. *arXiv preprint arXiv:1612.00215*, 2016.

[2] Xugong Qin, Yu Zhou, Youhui Guo, Dayan Wu, Zhihong Tian, Ning Jiang, Hongbin Wang, and Weiping Wang. Mask is all you need: Rethinking mask r-cnn for dense and arbitrary-shaped scene text detection. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 414–423, 2021.

[3] Fangneng Zhan, Shijian Lu, and Chuhui Xue. Verisimilar image synthesis for accurate detection and recognition of texts in scenes. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 249–266, 2018.

[4] Yexun Zhang, Ya Zhang, and Wenbin Cai. Separating style and content for generalized style transfer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8447–8455, 2018.

[5] Gantugs Atarsaikhan, Brian Kenji Iwana, Atsushi Narusawa, Keiji Yanai, and Seiichi Uchida. Neural font style transfer. In *2017 14th IAPR international conference on document analysis and recognition (ICDAR)*, volume 5, pages 51–56. IEEE, 2017.

[6] Samaneh Azadi, Matthew Fisher, Vladimir G Kim, Zhaowen Wang, Eli Shechtman, and Trevor Darrell. Multi-content gan for few-shot font style transfer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7564–7573, 2018.

[7] Jacky Cao, Kit-Yung Lam, Lik-Hang Lee, Xiaoli Liu, Pan Hui, and Xiang Su. Mobile augmented reality: User interfaces, frameworks, and intelligence. *ACM Computing Surveys*, 55(9):1–36, 2023.

[8] Victor Fragoso, Steffen Gauglitz, Shane Zamora, Jim Kleban, and Matthew Turk. Translatar: A mobile augmented reality translator. In *2011 IEEE workshop on applications of computer vision (WACV)*, pages 497–502. IEEE, 2011.

[9] Jun Du, Qiang Huo, Lei Sun, and Jian Sun. Snap and translate using windows phone. In *2011 International Conference on Document Analysis and Recognition*, pages 809–813. IEEE, 2011.

[10] Qirui Huang, Bin Fu, Yu Qiao, et al. Gentext: Unsupervised artistic text generation via decoupled font and texture manipulation. *arXiv preprint arXiv:2207.09649*, 2022.

[11] Yangming Shi, Haisong Ding, Kai Chen, and Qiang Huo. Aprnet: Attention-based pixel-wise rendering network for photo-realistic text image generation. *arXiv preprint arXiv:2203.07705*, 2022.

[12] 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*, pages 1500–1508, 2019.

[13] Junyeop Lee, Yoonsik Kim, Seonghyeon Kim, Moonbin Yim, Seung Shin, Gayoung Lee, and Sungrae Park. Rewritenet: Reliable scene text editing with implicit decomposition of text contents and styles. *arXiv preprint arXiv:2107.11041*, 2021.

[14] Song Park, Sanghyuk Chun, Junbum Cha, Bado Lee, and Hyunjung Shim. Multiple heads are better than one: Few-shot font generation with multiple localized experts. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13900–13909, 2021.

[15] Yadong Qu, Qingfeng Tan, Hongtao Xie, Jianjun Xu, Yuxin Wang, and Yongdong Zhang. Exploring stroke-level modifications for scene text editing. *arXiv preprint arXiv:2212.01982*, 2022.

[16] 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*, pages 13228–13237, 2020.

[17] 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*, pages 14700–14709, 2020.

[18] Paul Debevec. Rendering synthetic objects into real scenes: Bridging traditional and image-based graphics with global illumination and high dynamic range photography. In *Acm siggraph 2008 classes*, pages 1–10, 2008.

[19] Ankush Gupta, Andrea Vedaldi, and Andrew Zisserman. Synthetic data for text localisation in natural images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2315–2324, 2016.

[20] Max Jaderberg, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Synthetic data and artificial neural networks for natural scene text recognition. *arXiv preprint arXiv:1406.2227*, 2014.

[21] Fangneng Zhan, Hongyuan Zhu, and Shijian Lu. Spatial fusion gan for image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3653–3662, 2019.

[22] Shuai Yang, Zhangyang Wang, Zhaowen Wang, Ning Xu, Jiaying Liu, and Zongming Guo. Controllable artistic text style transfer via shape-matching gan. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4442–4451, 2019.

[23] Fangneng Zhan, Chuhui Xue, and Shijian Lu. Ga-dan: Geometry-aware domain adaptation network for scene text detection and recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9105–9115, 2019.

[24] 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.

[25] 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.

[26] Shuai Yang, Jiaying Liu, Zhouhui Lian, and Zongming Guo. Awesome typography: Statistics-based text effects transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7464–7473, 2017.

[27] Shuai Yang, Jiaying Liu, Wenjing Wang, and Zongming Guo. Tet-gan: Text effects transfer via stylization and destylization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 1238–1245, 2019.

[28] Lin Zhao, Changsheng Chen, and Jiwu Huang. Deep learning-based forgery attack on document images. *IEEE Transactions on Image Processing*, 30:7964–7979, 2021.

[29] Jeyasri Subramanian, Varnith Chordia, Eugene Bart, Shaobo Fang, Kelly Guan, Raja Bala, et al. Strive: Scene text replacement in videos. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14549–14558, 2021.

[30] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015.

[31] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[32] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *International conference on machine learning*, pages 7354–7363. PMLR, 2019.

[33] K Simonyan and A Zisserman. Very deep convolutional networks for large-scale image recognition. 3rd int conf learn represent iclr 2015-conf track proc. september 2014: 1–14, 2014.

[34] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.

[35] 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 *2013 12th international conference on document analysis and recognition*, pages 1484–1493. IEEE, 2013.

[36] Nibal Nayef, Fei Yin, Imen Bizid, Hyunsoo Choi, Yuan Feng, Dimosthenis Karatzas, Zhenbo Luo, Umapada Pal, Christophe Rigaud, Joseph Chazalon, et al. Icdar2017 robust reading challenge on multi-lingual scene text detection and script identification-rrc-mlt. In *2017 14th IAPR international conference on document analysis and recognition (ICDAR)*, volume 1, pages 1454–1459. IEEE, 2017.

[37] Nibal Nayef, Yash Patel, Michal Busta, Pinaki Nath Chowdhury, Dimosthenis Karatzas, Wafa Khlif, Jiri Matas, Umapada Pal, Jean-Christophe Burie, Cheng-lin Liu, et al. Icdar2019 robust reading challenge on multi-lingual scene text detection and recognition—rrc-mlt-2019. In *2019 International conference on document analysis and recognition (ICDAR)*, pages 1582–1587. IEEE, 2019.

[38] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[39] Sagar Imambi, Kolla Bhanu Prakash, and GR Kanagachidambaresan. Pytorch. *Programming with TensorFlow: Solution for Edge Computing Applications*, pages 87–104, 2021.

[40] 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.

[41] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and< 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.

## Supplementary Material

This is the supplementary material for the main paper. Please read the main version for the model formulation, its performance analysis, and key ablation experiments. All our pre-trained models and code will be made available at https://anonymous.4open.science/r/FAST-B5AE.

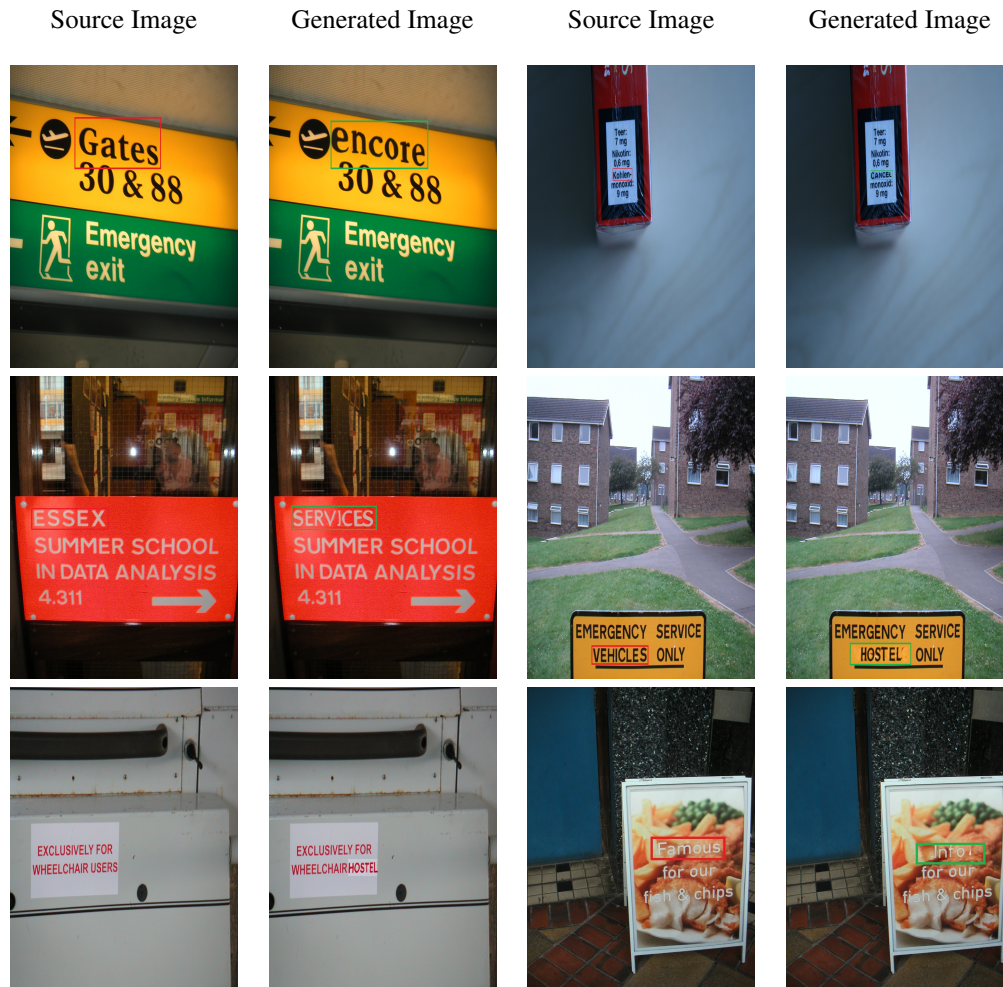| Source Image | Generated Image | Source Image | Generated Image |
| --- | --- | --- | --- |



Figure 7: Some additional samples of genrations from FAST.

Figure 8: Some real-scene editing illustration of our proposed method FAST.