

Calligrapher: Freestyle Text Image Customization

YUE MA* and QINGYAN BAI*, Hong Kong University of Science and Technology, China

HAO OUYANG, Ant Group, China

KA LEONG CHENG, Hong Kong University of Science and Technology, China

QIUYU WANG, Ant Group, China

HONGYU LIU and ZICHEN LIU, Hong Kong University of Science and Technology, China

HAOFAN WANG, InstantX, Independent Research Team

JINGYE CHEN, Hong Kong University of Science and Technology, China

YUJUN SHEN†, Ant Group, China

QIFENG CHEN†, Hong Kong University of Science and Technology, China

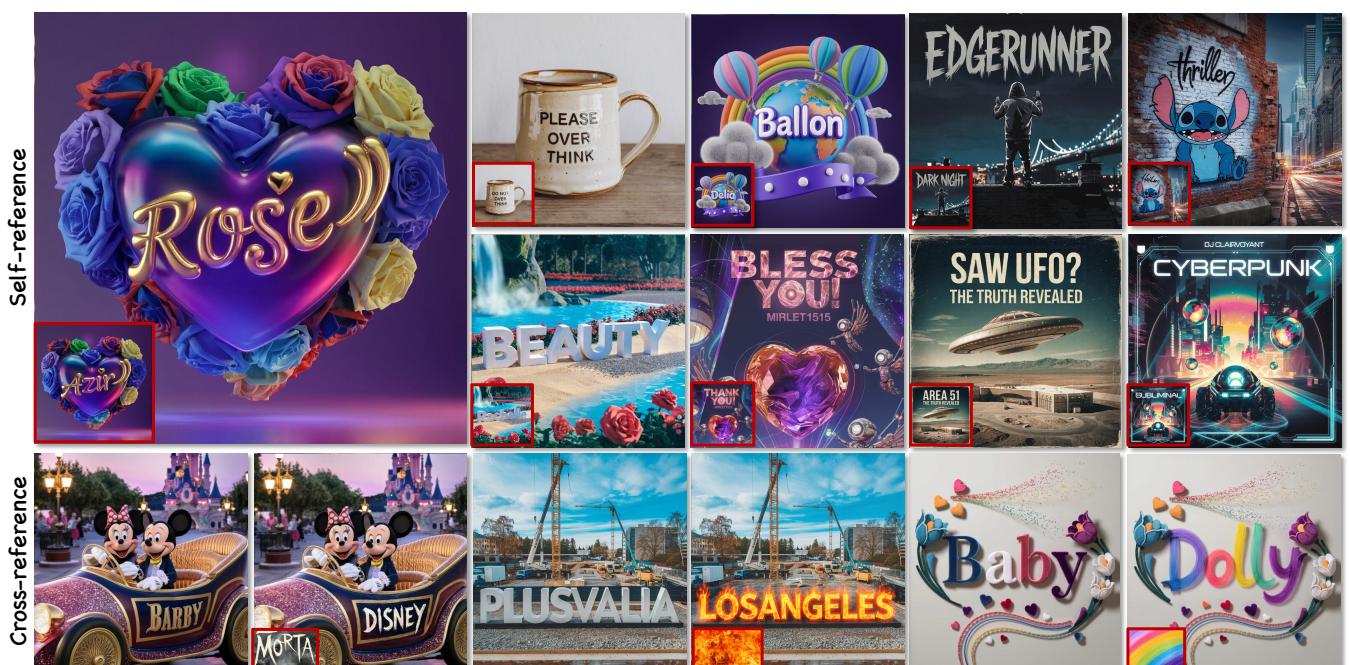


Fig. 1. **Photorealistic text image customization results** produced by our proposed **Calligrapher**, which allows users to perform customization with diverse stylized images and text prompts. The input and reference images are shown in the lower left corner of the generated results, respectively for the setting of self-reference and cross-reference text image customization.

We introduce Calligrapher, a novel diffusion-based framework that innovatively integrates advanced text customization with artistic typography for digital calligraphy and design applications. Addressing the challenges of precise style control and data dependency in typographic customization, our framework incorporates three key technical contributions. First, we

*Both authors contributed equally to this research.

†Corresponding authors.

Authors' Contact Information: Yue Ma; Qingyan Bai, Hong Kong University of Science and Technology, Hong Kong, China; Hao Ouyang, Ant Group, Hangzhou, China; Ka Leong Cheng, Hong Kong University of Science and Technology, Hong Kong, China; Qiuyu Wang, Ant Group, Hangzhou, China; Hongyu Liu; Zichen Liu, Hong Kong University of Science and Technology, Hong Kong, China; Haofan Wang, InstantX, Independent Research Team; Jingye Chen, Hong Kong University of Science and Technology, Hong Kong, China; Yujun Shen, Ant Group, Hangzhou, China; Qifeng Chen, Hong Kong University of Science and Technology, Hong Kong, China.

develop a self-distillation mechanism that leverages the pre-trained text-to-image generative model itself alongside the large language model to automatically construct a style-centric typography benchmark. Second, we introduce a localized style injection framework via a trainable style encoder, which comprises both Qformer and linear layers, to extract robust style features from reference images. An in-context generation mechanism is also employed to directly embed reference images into the denoising process, further enhancing the refined alignment of target styles. Extensive quantitative and qualitative evaluations across diverse fonts and design contexts confirm Calligrapher's accurate reproduction of intricate stylistic details and precise glyph positioning. By automating high-quality, visually consistent typography, Calligrapher surpasses traditional models, empowering creative practitioners in digital art, branding, and contextual typographic design. The code, model, and data can be found at the [Project Page](#).

Additional Key Words and Phrases: Text image customization, style transfer, diffusion models

1 Introduction

The advertising and promotion industry, encompassing digital media, branding, packaging, and printed materials, relies on vivid and meticulously crafted typography to effectively communicate messages and solidify brand identity. Currently, designers often dedicate substantial time to manually fine-tuning fonts to achieve specific aesthetic objectives. This process is not only labor-intensive but can also introduce inconsistencies. Therefore, an automated method capable of generating text that emulates a reference style while ensuring precise character positioning would significantly streamline the design workflow and enhance overall visual consistency, as demonstrated in Fig. 1.

Modern typography design, as illustrated in Fig. 2, predominantly employs two main categories of methods. The first category centers on the use of standardized font libraries [mdn web docs 1996]. While these libraries offer considerable accessibility, they often present challenges in seamless integration with diverse backgrounds and typically necessitate substantial manual adjustment to achieve specific aesthetic or artistic outcomes. The second branch of methods employs neural generative models [Huang et al. 2023; Mou et al. 2024; Yang et al. 2024; Zhang et al. 2023b] to enable typography generation, editing, automating text modification, and font creation. Although promising, this technique frequently fails to capture the precise nuances of specific font styles or handle styles different from those in the source image, which are difficult to express through textual description. Our work bridges these gaps by enhancing generative techniques to automate the typography customization process while ensuring that the final output closely adheres to the desired visual style.

Specifically, we propose a diffusion-based framework to address data dependency and precise style control through three technical contributions. Firstly, we introduce a self-distillation framework to construct a style-oriented typography training dataset. This framework leverages the pre-trained text-to-image generative model in conjunction with a large language model to synthesize a comprehensive set of text images. These are processed and paired with corresponding reference images, prompts, and masks, thereby creating self-supervised training data that facilitates style learning without requiring manual annotation. Based on the aforementioned data generation pipeline, we propose a style-centric text customization benchmark. This benchmark, including training and test sets, is expected to further boost the development of the typography research community. Secondly, a local style injection mechanism is designed to employ a trainable style encoder, including both Qformer [Li et al. 2023] and linear layers, to extract robust style-related features from references. By replacing cross-attention features in the denoising transformer network with these style embeddings, the method achieves granular typographic control in the latent space. Thirdly, an in-context generation mechanism directly integrates reference images into the denoising process. This integration significantly enhances the fidelity of style alignment between the generated output and the target references. These design elements enable our method to uniquely generate highly desirable text images that accurately reflect the style of reference inputs, even with *arbitrary text or non-text* images.

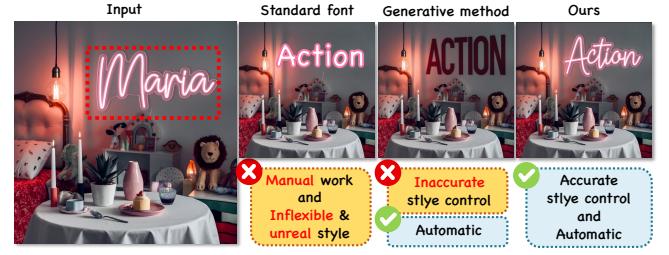


Fig. 2. Motivation and technical differentiation of our approach. Existing typography design methods face critical limitations: (1) Standard font libraries prioritize accessibility but require extensive manual adjustments for integration into diverse backgrounds, resulting in inflexible and unrealistic outputs. (2) Neural generative models automate typography but often fail to capture precise font style nuances, especially when relying on textual descriptions. In contrast, the proposed method addresses these challenges by enabling fully automated typography generation with precise style control and various kinds of references, including non-text images.

Overall, we construct and propose a style-centric text customization benchmark based on the self-distillation strategy, specifically addressing the critical need for the model learning and standardized evaluation in this field. The proposed model learned on the training set of this benchmark achieves success in text customization with various kinds of references, and has been extensively evaluated using both qualitative and quantitative methods, including user studies, demonstrating superior performance across multiple metrics compared to existing approaches. Further results also suggest the model also could be applied to the task of reference-based generation without tuning. Our work represents a significant step toward automated, efficient, and artistically driven typography design, with substantial potential applications in both design and branding processes, potentially revolutionizing workflows in creative industries by reducing manual labor while maintaining artistic integrity.

2 Related work

Visual text rendering. Visual text rendering has been a longstanding research problem in the era of generative AI. Traditional image generation models such as Stable Diffusion [Rombach et al. 2022a], Imagen [Saharia et al. 2022], and DALL-E [Ramesh et al. 2021] have fallen short in accurately rendering text. Consequently, some researchers have resorted to incorporating additional conditions into the generation process [Chen et al. 2023b,a, 2024, 2025; Ji et al. 2023; Jiang et al. 2025; Koo et al. 2025; Liu et al. 2022, 2024a,b; Ma et al. 2024, 2023; Tuo et al. 2024, 2023; Wang et al. 2025; Zhao and Lian 2023; Zhao et al. 2024]. For instance, GlyphDraw [Ma et al. 2023] introduces two diffusion branches: one for determining the text location and another for the actual text generation. While this approach somewhat alleviates the issues associated with text rendering, it remains limited to single-line text generation. TextDiffuser series [Chen et al. 2023a, 2024] seeks to address this limitation by using Transformers and Large Language Models (LLMs) to handle text positioning tasks, thereby extending capabilities to multi-line text rendering. AnyText [Tuo et al. 2024, 2023] introduces the generation of multi-language text images. Brush Your Text [Zhang et al. 2024] utilizes the canny map of a text template as a condition, achieving a



Fig. 3. **Self-distillation pipeline** for style-oriented typography dataset construction and model training. We emulate natural language processing practices by leveraging pre-trained text-to-image generative models and large language models to synthesize stylized text images, paired with reference prompts and masks. This generates self-supervised training pairs for robust style learning without manual annotation.

higher accuracy in multi-language rendering, though it falls short in terms of diversity. Overall, while existing efforts have focused on improving the accuracy of text rendering, the challenge of rendering controllable text remains substantial. Overall, although existing work has focused on the accuracy of text rendering, there is still a significant need to create more visually appealing and controllable text. This is the main focus of our research.

Text attributes customization. Early work has focused on font attribute customization [Gal et al. 2022; Hayashi et al. 2019; He et al. 2024a,a,c, 2023; Kondo et al. 2024; Wang et al. 2020; Yang et al. 2024]. For example, Attribute2Font [Wang et al. 2020] can automatically generate font styles by synthesizing visually pleasing glyph images based on user-specified attributes with corresponding values. However, compared to font stylizing, the task of customizing scene text attributes [He et al. 2024b; Paliwal et al. 2024; Su et al. 2023; Tuo et al. 2024] is more challenging due to complex factors such as perspective distortions and unique textures in scenes. For instance, MetaDesigner [He et al. 2024b] is a system that uses LLMs to facilitate the creation of customized artistic typography. It employs a multi-agent framework enhanced by a feedback loop from multimodal models and user evaluations, to produce aesthetically pleasing and contextually relevant WordArt that adapts to user preferences. AnyText2 [Tuo et al. 2024] explicitly designs a font encoder and a color encoder, providing additional style-related guidance during the rendering process. However, it is typically limited to generating simple fonts and struggles in producing artistic typography nor following the given styles, and occasionally generates blurred results. We believe that accommodating free-style fonts can make the visuals more dynamic and engaging.

Image style transfer. Image style transfer has been a longstanding research problem [Bai et al. 2024; Chen et al. 2017, 2021; Chung et al. 2024; Frenkel et al. 2024; Gal et al. 2022; Hertz et al. 2024; Isola et al. 2017; Karras et al. 2019, 2020; Kotovenko et al. 2019; Kwon and Ye 2022; Ouyang et al. 2025; Patashnik et al. 2021; Shah et al. 2024; Sohn et al. 2023; Wang et al. 2021, 2023; Zhang et al. 2023a, 2022; Zhu et al. 2017]. Within this domain, various methods have been proposed to tackle the challenge of transferring the style from one image to another. Pix2Pix [Isola et al. 2017] models style transfer as a low-level CNN prediction task, treating it as an image-to-image translation within a conditional GAN framework, where the model is trained on paired images to directly predict the transformation from content to style at the pixel level. On the other hand, CycleGAN [Zhu et al. 2017] employs a cycle consistency loss to enable style transfer in scenarios with unpaired images, using a GAN-based loss that

encourages the model to generate images that are indistinguishable from real images in the target style domain. InstructPix2Pix [Brooks et al. 2023] is built on top of the Stable Diffusion [Rombach et al. 2022a], enabling the usage of text prompt to conduct style transfer trained on curated paired images. Our main focus is to transfer the style of text based on a reference image and additional guidance such as color palette. Generally, text areas are relatively small, and there are also some minor stroke details that require precise rendering, which makes this task particularly challenging.

3 Methodology

The data generation and training pipeline of our method are shown in Fig. 3 and Fig. 4. Given the input image with mask, reference style image, and prompt, the purpose of our approach is to generate the text following the font style and customize it to the input source image, even for reference fonts of uncommon styles (i.e., cartoon, handwriting, and 3D style). In this section, we first discuss the motivation in Section 3.1, followed by the three carefully designed design components: The self-distillation learning strategy to cope with data scarcity is introduced in Section 3.2. Then, we describe the localized style injection mechanism in Section 3.3. Finally, we demonstrate the design of In-context inference for finer style consistency in Section 3.4.

3.1 Motivation

In this subsection, we identify several key limitations in current state-of-the-art approaches for real-world typography design and present corresponding motivations and solutions.

Scarcity of artistic typography data. A significant challenge in this domain is the limited availability of large-scale datasets dedicated to artistic typography. Our observations indicate that current diffusion models [Black-Forest-Labs 2024a] are capable of synthesizing high-quality stylized text when paired with robust post-processing and careful selection. We propose incorporating the model to generate a synthesized artistic typography benchmark and employ a self-distillation training strategy that leverages a high-quality synthesized dataset to effectively transfer artistic styles, as elaborated in Section 3.2.

Failure to capture subtle font details. Existing methods often rely on global stylization techniques that are insufficient for capturing shapes and textures, focusing only on the task of self-reference inpainting. To address this limitation, we introduce a novel training pipeline that emphasizes localized style injection. This pipeline concentrates on fine-grained detail refinement, and yields a more

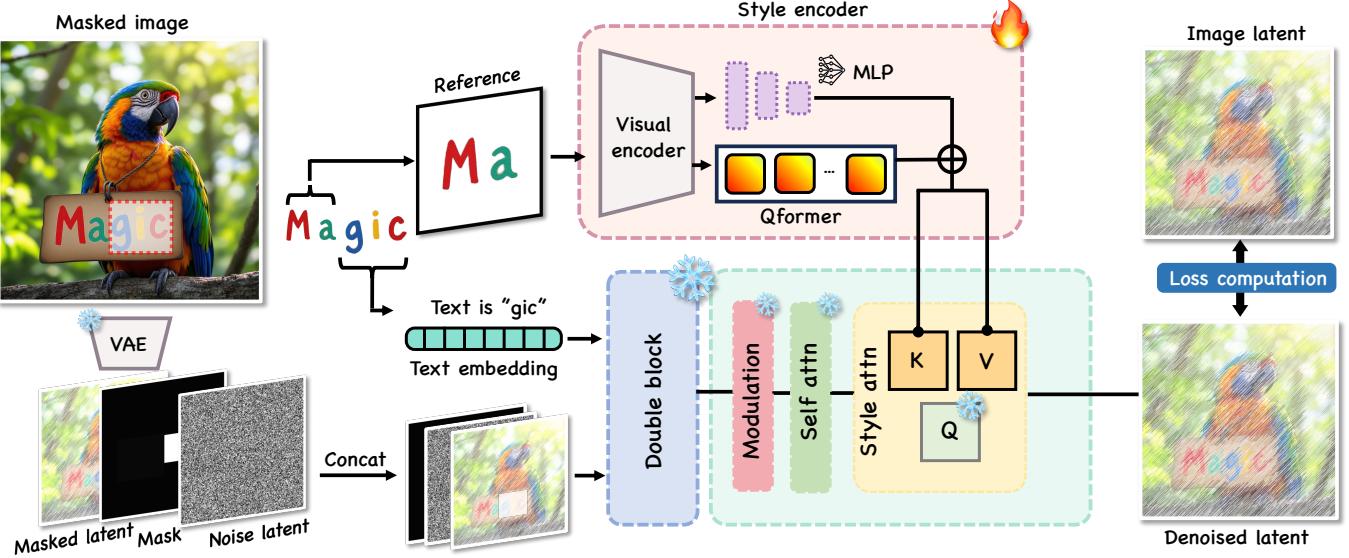


Fig. 4. Training framework of Calligrapher, demonstrating the integration of localized style injection and diffusion-based learning. The framework processes masked images through a Variational Auto-Encoder (VAE) to obtain latent representations, concatenated with mask and noise latents. A style encoder comprising a visual encoder, Qformer, and linear layers is designed to extract style-related features from the reference style image, while text embeddings (e.g., “gic” in the case) modulate the denoising transformer. In the denoising block, style attention predicted from the style features replaces the original cross-attention, injecting style embeddings (K_E, V_E) with the denoiser’s query Q to enable granular typographic control in the latent space. The model is optimized under the flow-matching learning objective with the self-distillation typography dataset.

faithful reproduction with the in-context generation techniques, as outlined in Section 3.3 and Section 3.4.

3.2 Self distillation & stylized typography benchmark

Unlike image translation tasks in ControlNet [Zhang et al. 2023b], acquiring high-quality supervised training data for text style transfer remains challenging [Black-Forest-Labs 2024a,b; Ye et al. 2023] due to the prohibitive cost and effort required to manually curate large-scale datasets of text pairs, which exhibit identical semantic content but distinct stylistic attributes. Furthermore, such datasets require diverse and sufficiently rich stylistic variations to enable models to robustly capture nuanced style features and adapt to complex style transfer scenarios. With the finding that modern generative models [Black-Forest-Labs 2024a] could produce text images with desirable quality, we draw inspiration from recent advances in self-training paradigms within large language model (LLM) research [Huang et al. 2022], where a robust generative model is adopted to yield data to train itself. As in Fig. 3, our proposed framework introduces a novel methodology where the pretrained generative model is employed to: (1) synthesize stylistically consistent training data through controlled generation, and (2) refine the style transfer model using the self-generated corpus. This approach establishes a learning system that effectively leverages the internal knowledge representation of generative models while circumventing the dependency on human-annotated paired examples.

Specifically, as in Fig. 3, we first leverage large language models (LLMs) to generate a diverse set of semantic-coherent prompts p annotated with explicit typographic style descriptors (e.g., “3D metallic

text,” “watercolor calligraphy”). These style-conditioned prompts are subsequently fed into the flow-matching diffusion model \mathcal{G}_θ [Black-Forest-Labs 2024a], to synthesize high-fidelity stylized text images through iterative denoising processes. To construct training pairs from the synthesized corpus, we first adopt the neural text understanding method [AI 2023] to detect the text locations and employ a strategic cropping mechanism that preserves typographic consistency while enabling effective self-supervision. For each generated image, we randomly crop a local region containing stylized characters as the reference style exemplar, while maintaining the remaining text region as the target for style transfer learning. Based on the aforementioned data generation pipeline, we establish and propose a style-centric text customization benchmark to benefit the development of the community. The details of this stylized typography benchmark can be found in the data webpage.

To formalize the task, let x represent the main inputs of the text customization task that include the image latent, mask, and noise latent, while y stands for the reference image. The proposed data generation strategy allows the model to efficiently learn to capture localized stylistic patterns and generate target text images from Gaussian noise $\varepsilon \sim \mathcal{N}(0, I)$, via the flow matching objective [Esser et al. 2024]:

$$\min_{\theta} \mathbb{E}_{x_0 \sim p(x_0), \varepsilon \sim \mathcal{N}(0, I)} [\lambda(t) \|D(x_t, t, p, y) - x_0\|_2^2], \quad (1)$$

where $D(x_t, t, p, y) = x_t - t \cdot \mathcal{G}_\theta(x_t, t, p, y)$ as in [Esser et al. 2024], t stands for the timestep, x_t denotes noisy inputs at t , and $\lambda(t)$ indicates the loss weighting.

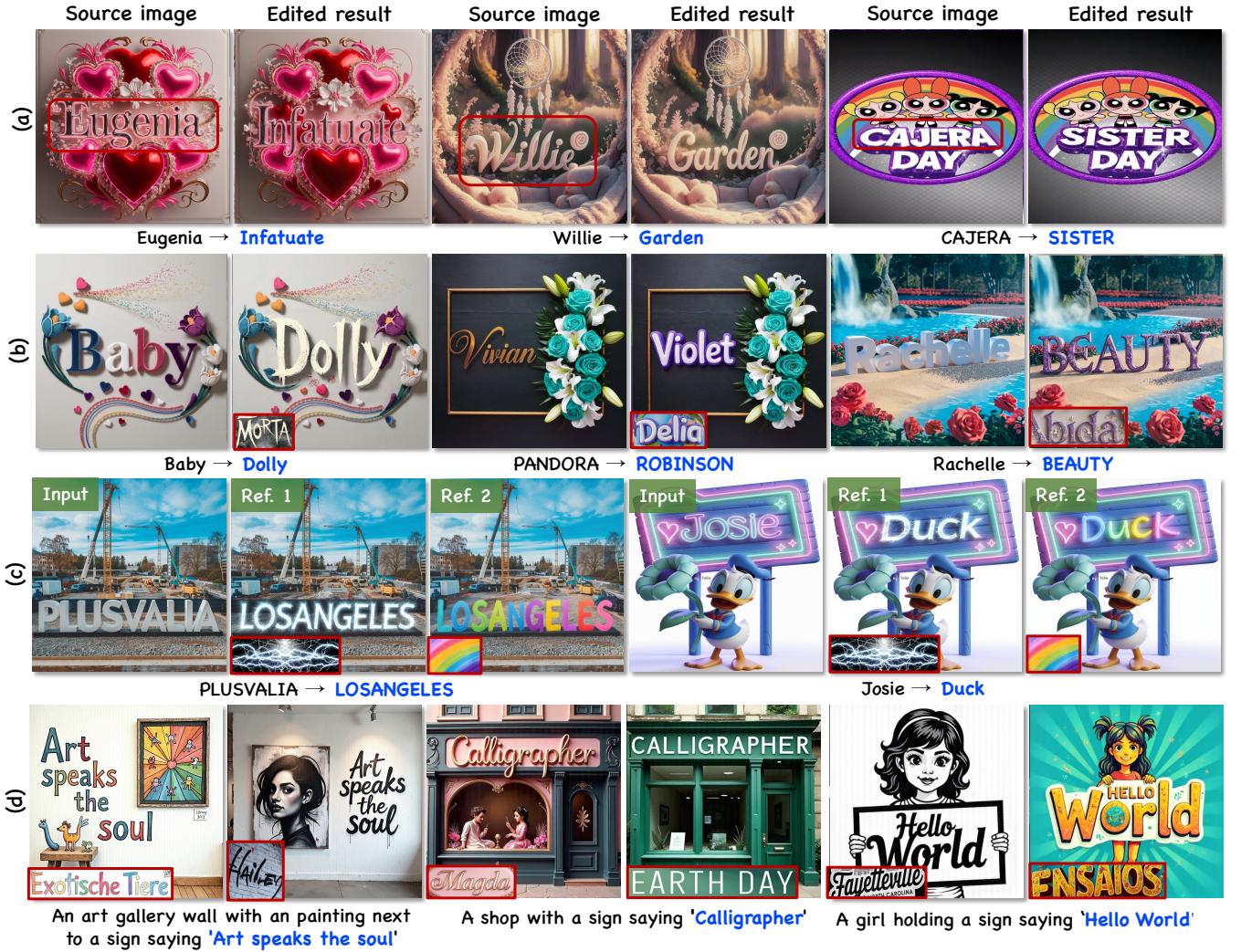


Fig. 5. **Qualitative results of Calligrapher under various settings.** We demonstrate text customization results respectively under settings of (a) self-reference, (b) cross-reference, and (c) non-text reference. Reference-based image generation results are also incorporated in (d).

3.3 Localized style injection

In order to achieve text customization, we follow ControlNet [Zhang et al. 2023b] and IP-Adapter [Ye et al. 2023] to learn another controllable branch (namely the style encoder \mathcal{E}) to encode the conditional control signals while the original denoiser serves as the main branch, making the denoising formulation as follows:

$$\mathcal{G}(\mathbf{x}_t, t, p, \mathbf{y}) = \mathbf{F}(\mathcal{D}(\mathbf{x}_t, t, p), \mathcal{E}(\mathbf{y})), \quad (2)$$

where \mathbf{F} indicates the fusion function for the features of the style encoder \mathcal{E} and the main denoising network \mathcal{D} . To extract initial features from the reference, we instantiate the style encoder with a pre-trained multi-modal visual encoder [Zhai et al. 2023] and another encoder composed of linear layers and Qformer [Li et al. 2023] with learnable query parameters to transform these features into the key and value matrices. The fusion function \mathbf{F} is instantiated as feature replacement and cross-attention. The key and value matrices predicted from the style encoder are then injected into the main

branch \mathcal{D} , by replacing the original Key and Value matrices in the style attention module of the single block, as in Fig. 4:

$$\text{StyleAttention}(Q_{\mathcal{D}}, K_{\mathcal{E}}, V_{\mathcal{E}}) = \text{softmax}\left(\frac{Q_{\mathcal{D}}K_{\mathcal{E}}^T}{\sqrt{d_K}}\right)V_{\mathcal{E}}, \quad (3)$$

where d_K indicates the tensor dimension. These features from style attention would be added to the original attention activations for modulation. We follow prior art [Rombach et al. 2022b] to perform training and inference in the Variational Auto-Encoder (VAE) latent space for efficiency.

3.4 In-context generation

Motivated by recent works [Huang et al. 2024; Zhang et al. 2025] demonstrating the strong contextual capabilities of diffusion-based

Table 1. **Quantitative comparisons with SOTA baselines** including TextDiffuser-2 [Chen et al. 2024], AnyText [Tuo et al. 2023] and FLUX-Fill [Black-Forest-Labs 2024b]. Our method demonstrates the best or comparable performance across multiple metrics. The metrics for the best-performing method are highlighted in bold.

Method	Metrics				User Study			
	FID ↓	CLIP ↑	DINO ↑	OCR _{Acc} ↑	Style Sync ↑	Text Matching ↑	Aesthetic ↑	Overall ↑
TextDiffuser-2	66.68	0.7097	0.8914	0.81	2.42	2.40	2.37	0.10
AnyText	69.72	0.7041	0.8821	0.45	1.98	1.68	1.95	0.04
FLUX-Fill	67.79	0.7090	0.8984	0.61	2.20	2.52	2.35	0.14
Ours	38.09	0.7401	0.9474	0.84	3.40	3.40	3.32	0.72

generative models [Black-Forest-Labs 2024a], we explore if reference-based text customization enables to be improved by in-context inference. Specifically, our approach explicitly embeds contextual information - serving as a style reference - into the denoising trajectory by means of spatial concatenation at the pixel level. This composite image is then encoded through the shared VAE to yield a unified and contextualized latent representation. This latent feature, together with a correspondingly constructed binary mask that zeroes out the region occupied by the reference, is then sent to DiT to condition the denoising of Gaussian noise. The resulting context-aware latent encapsulates both the semantic content to be edited and the stylistic cues from the reference, forming a holistic conditioning signal for the subsequent diffusion process. As a result, this design enables fine-grained style coherence while preserving structural fidelity in the generated text.

4 Experiment

4.1 Implementation Details

In our experiment, we adopt FLUX-Fill [Black-Forest-Labs 2024b] and FLUX [Black-Forest-Labs 2024a] as the base model for customization and generation. The visual encoder is based on siglipatch14 [Zhai et al. 2023] and Qformer [Liu et al. 2023]. In the training phase, we freeze the FLUX model parameters to maintain its powerful generation ability. The localized style injection module is trained for 100,000 steps using 8 Tesla A800 GPUs, taking approximately 10 days. The AdamW optimizer is employed with a learning rate of 2×10^{-5} and a batch size of 32. In the inference phase, we employ the flow-matching Euler scheduler [Esser et al. 2024] with sampling steps of 50 and a guidance scale of 30.0.

4.2 Settings and Applications

Self-reference text image customization. One of the applications of our method is to modify the text content in the input image following the original text style. As shown in Fig. 5(a), our approach allows for the editing of text content while preserving the original text style, achieved by simply modifying the relevant descriptions in the input text prompt. For example, given an input image, our approach manages to inpaint the “Eugenia” to “Infatuate”, “Willie” to “Garden” (Fig. 5(a)). The background of the input and output images remains consistent. Considering previous works are only enable to perform this mentioned task of self-reference customization (inpainting), we conduct quantitative and qualitative comparisons under this setting in the latter sections. We also demonstrate the further unique capabilities of our model beyond this setting.

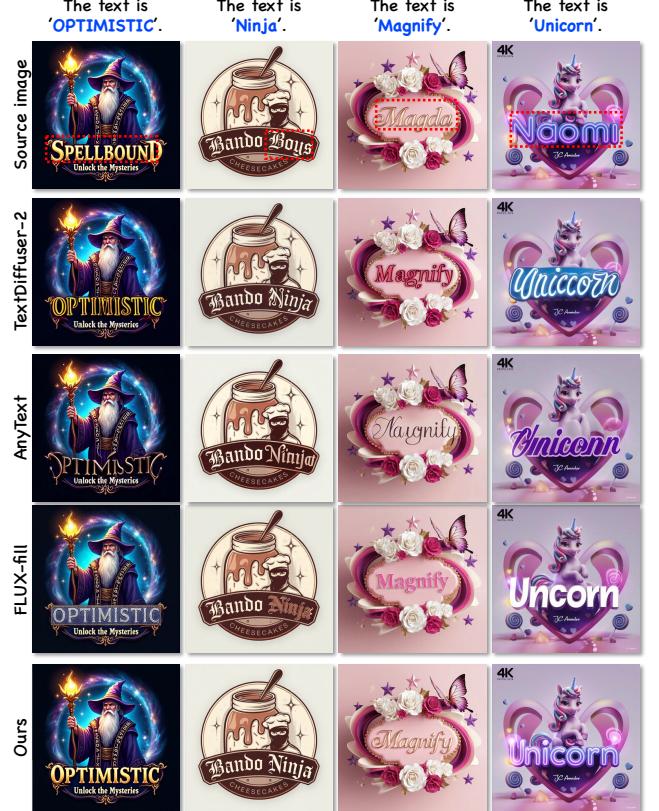


Fig. 6. **Qualitative comparisons on self-reference customization.** Calligrapher achieves better performance in terms of style sync and quality.

Cross-reference text image customization. Cross-reference text customization aims to edit the text content using the reference with different style, which has never been demonstrated in previous methods [Chen et al. 2023a; Tuo et al. 2023]. In Fig. 5(b), we present various customization results given different styles of reference text images. Our approach is capable of generating style-aligned images while ensuring the controllability of the text. On the other hand, we empirically find that, the text customization model also works well when non-text images serve as the reference, such as images of fire, rainbows, and lightning. As in Fig. 5(c), our approach generates text that well aligns with these styles. The generated image also maintains a high level of background consistency and achieves impressive aesthetic quality.



Fig. 7. **Ablation studies** on the self-distillation (left) and in-context generation (right) to validate their effectiveness.

Reference-based text image generation. Furthermore, we make efforts to achieve additional global controllable tasks of reference-based text image generation, where the input \mathbf{x} in Eq. (2) only includes the noise latent. We find the style encoder trained based on the original main branch (FLUX-fill [Black-Forest-Labs 2024b]) works with the new main branch of FLUX [Black-Forest-Labs 2024a], which could enable reference-based text image generation without further training and suggests the generalization of the learned model as in Fig. 5(d). This may be attributed to the parameter similarity between these two base models.

4.3 Comparison with baselines

Quantitative results. For quantitative evaluation, we compare our method with state-of-the-art methods on the test set of our typography benchmark, which includes 100 text images with masks, prompts, and corresponding references. FID [Heusel et al. 2017] is adopted to evaluate the general quality and similarity of the whole images following prior arts. We also compute the style similarity of the text images within masked regions respectively, with the CLIP ViT-base [Radford et al. 2021] and DINO-v2 [Oquab et al. 2023] models. For the OCR metrics, we utilize the Google Cloud text detection API [Google-Cloud-API 2025] to recognize the content and calculate the accuracy of generated text. Results shown in Table 1 demonstrate that the proposed method achieves the best in terms of all metrics. The user study, conducted with 30 participants yielding over 1000 votes, provides results including three sub-domain scores (on a scale of 1-4) and an overall preference percentage, which further demonstrate that our approach achieves the best performance.

Qualitative results. Qualitative comparisons with TextDiffuser-2 [Chen et al. 2024], AnyText [Tuo et al. 2023], and FLUX-fill [Black-Forest-Labs 2024b] are shown in Fig. 6, TextDiffuser-2 struggles in synthesizing the correct characters and styles. AnyText also generates text images in an undesirable style and low visual quality. It occasionally generates incorrect characters such as “Ninja” and “Magnify”. FLUX-fill [Black-Forest-Labs 2024b] demonstrates competent lexical accuracy but suffers from stylistic inconsistency, whereas the proposed method achieves substantial superiority in both dimensions. Compared to existing methods, Calligrapher demonstrates

significant advantages in terms of textual correctness and style consistency. A notable example is the distinctive pattern of the “D” letters in the reference word “SPELLBOUND” where our method maintains superior glyph integrity and stylistic coherence during generation.

4.4 Ablation studies

Effectiveness of self-distillation. We evaluate the impact of self-distillation on style similarity in text image customization. For comparison, we show generated results from the model with and without the self-distillation training method. As shown in Fig. 7 (left), the model with self-distillation achieves significantly higher style consistency between generated images. This demonstrates that self-distillation leverages the generative model’s internal knowledge to create stylistically coherent training pairs, circumventing the scarcity of manually curated paired data and enabling the model to robustly learn and transfer nuanced style characteristics.

Effectiveness of in-context generation. We also evaluate the effectiveness of the in-context strategy during the inference stage. As shown in the right subfigure in Fig. 7, it is clear to observe that the generated results achieves better style consistency the in-context strategy. We analyze that this is because the DiT structure incorporates self-attention, which is calculated on all tokens. The in-context strategy helps enhance the interaction of attention between reference text images and generated results.

5 Conclusion

Automating typography customization is critical for advertising. This work addresses labor-intensive manual font tuning by proposing a diffusion-based framework for automated typography customization with style consistency. Our key contributions include a self-distillation dataset construction pipeline, local style injection via trainable encoders, and in-context generation integrating references. A style-centric benchmark is also constructed to facilitate text customization. Experiments show our model enables accurate style replication for arbitrary text or non-text inputs of diverse styles. This advances efficient, artistic typography design, reducing manual effort and enhancing workflow consistency in creative industries.

References

- Jaided AI. 2023. EasyOCR. <https://github.com/JaidedAI/EasyOCR>.
- Qingyan Bai, Hao Ouyang, Yinghao Xu, Qiuyu Wang, Ceyuan Yang, Ka Leong Cheng, Yujun Shen, and Qifeng Chen. 2024. Edicho: Consistent Image Editing in the Wild. *arXiv preprint arXiv:2412.21079* (2024).
- Black-Forest-Labs. 2024a. FLUX. <https://github.com/black-forest-labs/flux>.
- Black-Forest-Labs. 2024b. FLUX.1 Tools. <https://blackforestlabs.ai/flux-1-tools/>.
- Tim Brooks, Aleksander Holynski, and Alexei A Efros. 2023. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Dongdong Chen, Lu Yuan, Jing Liao, Nenghai Yu, and Gang Hua. 2017. Stylebank: An explicit representation for neural image style transfer. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*.
- Haibo Chen, Zhizhong Wang, Huiming Zhang, Zhiwen Zuo, Ailin Li, Wei Xing, Dongming Lu, et al. 2021. Artistic style transfer with internal-external learning and contrastive learning. *Advances in Neural Information Processing Systems* (2021).
- Haoxing Chen, Zhuoer Xu, Zhangxuan Gu, Yaohui Li, Changhua Meng, Huijia Zhu, Weiqiang Wang, et al. 2023b. Diffute: Universal text editing diffusion model. *Advances in Neural Information Processing Systems* (2023).
- Jingye Chen, Yupan Huang, Tengchao Lv, Lei Cui, Qifeng Chen, and Furu Wei. 2023a. Textdiffuser: Diffusion models as text painters. *Advances in Neural Information Processing Systems* (2023).
- Jingye Chen, Yupan Huang, Tengchao Lv, Lei Cui, Qifeng Chen, and Furu Wei. 2024. Textdiffuser-2: Unleashing the power of language models for text rendering. In *European Conference on Computer Vision*.
- SiXiang Chen, Jianyu Lai, Jialin Gao, Tian Ye, Haoyu Chen, Hengyu Shi, Shitong Shao, Yunlong Lin, Song Fei, Zhaohua Xing, Yeying Jin, Junfeng Luo, Xiaoming Wei, and Lei Zhu. 2025. PosterCraft: Rethinking High-Quality Aesthetic Poster Generation in a Unified Framework. *arXiv preprint arXiv:2506.10741* (2025).
- Jiwoo Chung, Sangeek Hyun, and Jae-Pil Heo. 2024. Style injection in diffusion: A training-free approach for adapting large-scale diffusion models for style transfer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. 2024. Scaling rectify flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*.
- Yarden Frenkel, Yael Vinker, Ariel Shamir, and Daniel Cohen-Or. 2024. Implicit style-content separation using B-LoRA. In *European Conference on Computer Vision*.
- Rinon Gal, Or Patashnik, Haggai Maron, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. 2022. Stylegan-nada: Clip-guided domain adaptation of image generators. *ACM Transactions on Graphics (TOG)* (2022).
- Google-Cloud-API. 2025. Google Cloud text detection API. <https://cloud.google.com/vision/docs/ocr>.
- Hideaki Hayashi, Kohtarao Abe, and Seiichi Uchida. 2019. GlyphGAN: Style-consistent font generation based on generative adversarial networks. *Knowledge-Based Systems* (2019).
- Haibin He, Xinyuan Chen, Chaoyue Wang, Juhua Liu, Bo Du, Dacheng Tao, and Qiao Yu. 2024a. Diff-font: Diffusion model for robust one-shot font generation. *International Journal of Computer Vision* (2024).
- Jun-Yan He, Zhi-Qi Cheng, Chenyang Li, Jingdong Sun, Qi He, Wangmeng Xiang, Hanqian Chen, Jin-Peng Lan, Xianhui Lin, Kang Zhu, et al. 2024b. MetaDesigner: Advancing artistic typography through AI-Driven, user-Centric, and multilingual wordArt synthesis. *arXiv preprint arXiv:2406.19859* (2024).
- Jun-Yan He, Zhi-Qi Cheng, Chenyang Li, Jingdong Sun, Wangmeng Xiang, Yusen Hu, Xianhui Lin, Xiaoyang Kang, Zengke Jin, Bin Luo, et al. 2024c. WordArt Designer API: User-Driven artistic typography synthesis with large language models on modelScope. *arXiv preprint arXiv:2401.01699* (2024).
- Jun-Yan He, Zhi-Qi Cheng, Chenyang Li, Jingdong Sun, Wangmeng Xiang, Xianhui Lin, Xiaoyang Kang, Zengke Jin, Yusen Hu, Bin Luo, et al. 2023. WordArt designer: user-driven artistic typography synthesis using large language models. *arXiv preprint arXiv:2310.18332* (2023).
- Amir Hertz, Andrey Voynov, Shlomi Fruchter, and Daniel Cohen-Or. 2024. Style aligned image generation via shared attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. GANs trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*.
- Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2022. Large language models can self-improve. *arXiv preprint arXiv:2210.11610* (2022).
- Lianghua Huang, Di Chen, Yu Liu, Yujun Shen, Deli Zhao, and Jingren Zhou. 2023. Composer: Creative and controllable image synthesis with composable conditions. *arXiv preprint arXiv:2302.09778* (2023).
- Lianghua Huang, Wei Wang, Zhi-Fan Wu, Yupeng Shi, Huanzhang Dou, Chen Liang, Yutong Feng, Yu Liu, and Jingren Zhou. 2024. In-Context LoRA for Diffusion Transformers. *arXiv preprint arxiv:2410.23775* (2024).
- Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. 2017. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Jiabao Ji, Guanhua Zhang, Zhaowen Wang, Bairu Hou, Zhifei Zhang, Brian Price, and Shiyu Chang. 2023. Improving diffusion models for scene text editing with dual encoders. *arXiv preprint arXiv:2304.05568* (2023).
- Bowen Jiang, Yuan Yuan, Xinyi Bai, Zhuoqun Hao, Alyson Yin, Yaojie Hu, Wenyu Liao, Lyle Ungar, and Camillo J Taylor. 2025. ControlText: Unlocking controllable fonts in multilingual text rendering without font annotations. *arXiv preprint arXiv:2502.10999* (2025).
- Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2020. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Tetta Kondo, Shumpei Takezaki, Daichi Haraguchi, and Seiichi Uchida. 2024. Font style interpolation with diffusion models. In *International Conference on Document Analysis and Recognition*.
- Myungkyu Koo, Subin Kim, Sangkyung Kwak, Jaehyun Nam, Seojin Kim, and Jinwoo Shin. 2025. FontAdapter: Instant Font Adaptation in Visual Text Generation. *arXiv preprint arXiv:2506.05843* (2025).
- Dmytro Kotovenko, Artsiom Sanakoyeu, Pingchuan Ma, Sabine Lang, and Bjorn Ommer. 2019. A content transformation block for image style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Gihyun Kwon and Jong Chul Ye. 2022. Clipstyler: Image style transfer with a single text condition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International Conference on Machine Learning*.
- Gongyu Liu, Menghan Xia, Yong Zhang, Haoxin Chen, Jinbo Xing, Xintao Wang, Yujiu Yang, and Ying Shan. 2023. StyleCrafter: Enhancing Stylized Text-to-Video Generation with Style Adapter. *arXiv preprint arXiv:2312.00330* (2023).
- Rosanne Liu, Dan Garrette, Chitwan Saharia, William Chan, Adam Roberts, Sharad Narang, Irina Blok, RJ Mical, Mohammad Norouzi, and Noah Constant. 2022. Character-aware models improve visual text rendering. *arXiv preprint arXiv:2212.10562* (2022).
- Zeyu Liu, Weicong Liang, Zhanhao Liang, Chong Luo, Ji Li, Gao Huang, and Yuhui Yuan. 2024a. Glyph-byt5: A customized text encoder for accurate visual text rendering. In *European Conference on Computer Vision*.
- Zeyu Liu, Weicong Liang, Yiming Zhao, Bohan Chen, Lin Liang, Lijuan Wang, Ji Li, and Yuhui Yuan. 2024b. Glyph-byt5-v2: A strong aesthetic baseline for accurate multilingual visual text rendering. *arXiv preprint arXiv:2406.10208* (2024).
- Jian Ma, Yonglin Deng, Chen Chen, Nanyang Du, Haonan Lu, and Zhenyu Yang. 2024. Glyphdraw2: Automatic generation of complex glyph posters with diffusion models and large language models. *arXiv preprint arXiv:2407.02252* (2024).
- Jian Ma, Mingjun Zhao, Chen Chen, Ruichen Wang, Di Niu, Haonan Lu, and Xiaodong Lin. 2023. Glyphdraw: Seamlessly rendering text with intricate spatial structures in text-to-image generation. *arXiv preprint arXiv:2303.17870* (2023).
- mdn web docs. 1996. <https://developer.mozilla.org/en-US/docs/Web/CSS/font-family>
- Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, and Ying Shan. 2024. Ti2-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 38. 4296–4304.
- Maxime Oquab, Timothée Darcret, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. 2023. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193* (2023).
- Ziheng Ouyang, Zhen Li, and Qibin Hou. 2025. K-LoRA: Unlocking training-free fusion of any subject and style loras. *arXiv preprint arXiv:2502.18461* (2025).
- Shubham Paliwal, Arushi Jain, Monika Sharma, Vikram Jamwal, and Lovekesh Vig. 2024. CustomText: Customized textual image generation using diffusion models. *arXiv preprint arXiv:2405.12531* (2024).
- Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. 2021. Styleclip: Text-driven manipulation of stylegan imagery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation. In *International Conference on Machine Learning*.

- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022a. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022b. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lali Li, Jay Whang, Emily L Denton, Kamiar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems* (2022).
- Viraj Shah, Nataniel Ruiz, Forrester Cole, Erika Lu, Svetlana Lazebnik, Yuanzhen Li, and Varun Jampani. 2024. Ziplora: Any subject in any style by effectively merging loras. In *European Conference on Computer Vision*.
- Kihyuk Sohn, Nataniel Ruiz, Kimin Lee, Daniel Castro Chin, Irina Blok, Huiwen Chang, Jarred Barber, Lu Jiang, Glenn Entis, Yuanzhen Li, et al. 2023. Styledrop: Text-to-image generation in any style. *arXiv preprint arXiv:2306.00983* (2023).
- Tonghua Su, Fuxiang Yang, Xiang Zhou, Donglin Di, Zhongjie Wang, and Songze Li. 2023. Scene style text editing. *arXiv preprint arXiv:2304.10097* (2023).
- Yuxiang Tuo, Yifeng Geng, and Liefeng Bo. 2024. AnyText2: Visual text generation and editing with customizable attributes. *arXiv preprint arXiv:2411.15245* (2024).
- Yuxiang Tuo, Wangmeng Xiang, Jun-Yan He, Yifeng Geng, and Xuansong Xie. 2023. Anytext: Multilingual visual text generation and editing. *arXiv preprint arXiv:2311.03054* (2023).
- Alex Jinpeng Wang, Dongxing Mao, Jiawei Zhang, Weiming Han, Zhubai Dong, Linjie Li, Yiqi Lin, Zhengyuan Yang, Libo Qin, Fuwei Zhang, et al. 2025. TextAtlas5M: A large-scale dataset for dense text image generation. *arXiv preprint arXiv:2502.07870* (2025).
- Pei Wang, Yijun Li, and Nuno Vasconcelos. 2021. Rethinking and improving the robustness of image style transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Yizhi Wang, Yue Gao, and Zhouhui Lian. 2020. Attribute2font: Creating fonts you want from attributes. *ACM Transactions on Graphics* (2020).
- Zhizhong Wang, Lei Zhao, and Wei Xing. 2023. Stylediffusion: Controllable disentangled style transfer via diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*.
- Zhenhua Yang, Dezheng Peng, Yuxin Kong, Yuyi Zhang, Cong Yao, and Lianwen Jin. 2024. Fontdiffuser: One-shot font generation via denoising diffusion with multi-scale content aggregation and style contrastive learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. 2023. IP-Adapter: Text compatible image prompt adapter for text-to-image diffusion models. *arXiv preprint arxiv:2308.06721* (2023).
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid Loss for Language Image Pre-Training. *arXiv preprint arXiv:2303.15343* (2023).
- Lingjun Zhang, Xinyuan Chen, Yaohui Wang, Yue Lu, and Yu Qiao. 2024. Brush your text: Synthesize any scene text on images via diffusion model. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. 2023b. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*.
- Yuxin Zhang, Nisha Huang, Fan Tang, Haibin Huang, Chongyang Ma, Weiming Dong, and Changsheng Xu. 2023a. Inversion-based style transfer with diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*.
- Yuxin Zhang, Fan Tang, Weiming Dong, Haibin Huang, Chongyang Ma, Tong-Yee Lee, and Changsheng Xu. 2022. Domain enhanced arbitrary image style transfer via contrastive learning. In *ACM SIGGRAPH 2022 Conference Proceedings*.
- Zechuan Zhang, Ji Xie, Yu Lu, Zongxin Yang, and Yi Yang. 2025. In-Context Edit: Enabling Instructional Image Editing with In-Context Generation in Large Scale Diffusion Transformer. *arXiv preprint arxiv:2410.23775* (2025).
- Yiming Zhao and Zhouhui Lian. 2023. Uddifftext: A unified framework for high-quality text synthesis in arbitrary images via character-aware diffusion models. *arXiv preprint arXiv:2312.04884* (2023).
- Zhen Zhao, Jingqun Tang, Binghong Wu, Chunhui Lin, Shu Wei, Hao Liu, Xin Tan, Zhizhong Zhang, Can Huang, and Yuan Xie. 2024. Harmonizing visual text comprehension and generation. *arXiv preprint arXiv:2407.16364* (2024).
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision*.

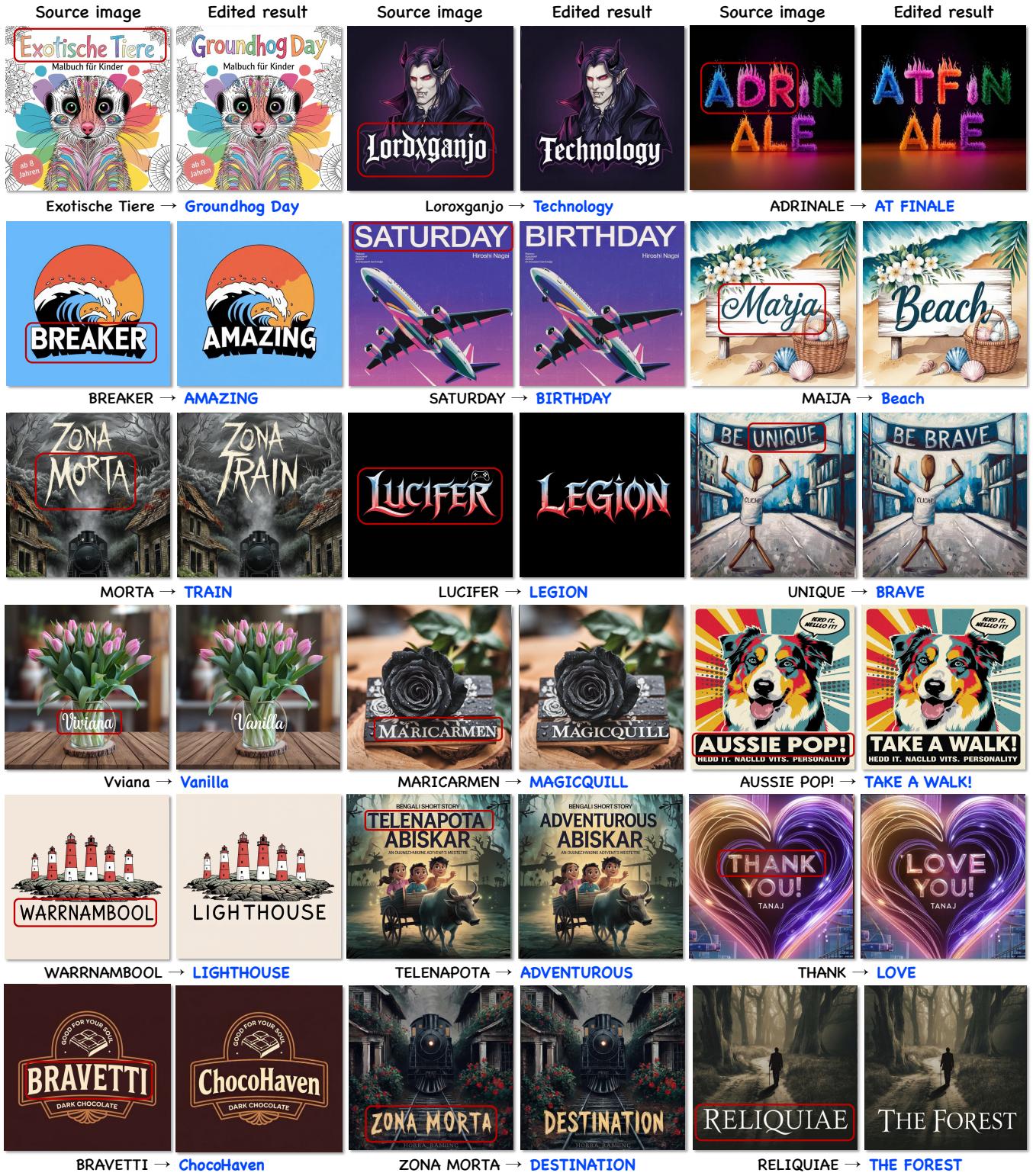


Fig. 8. Self-reference text image customization results.

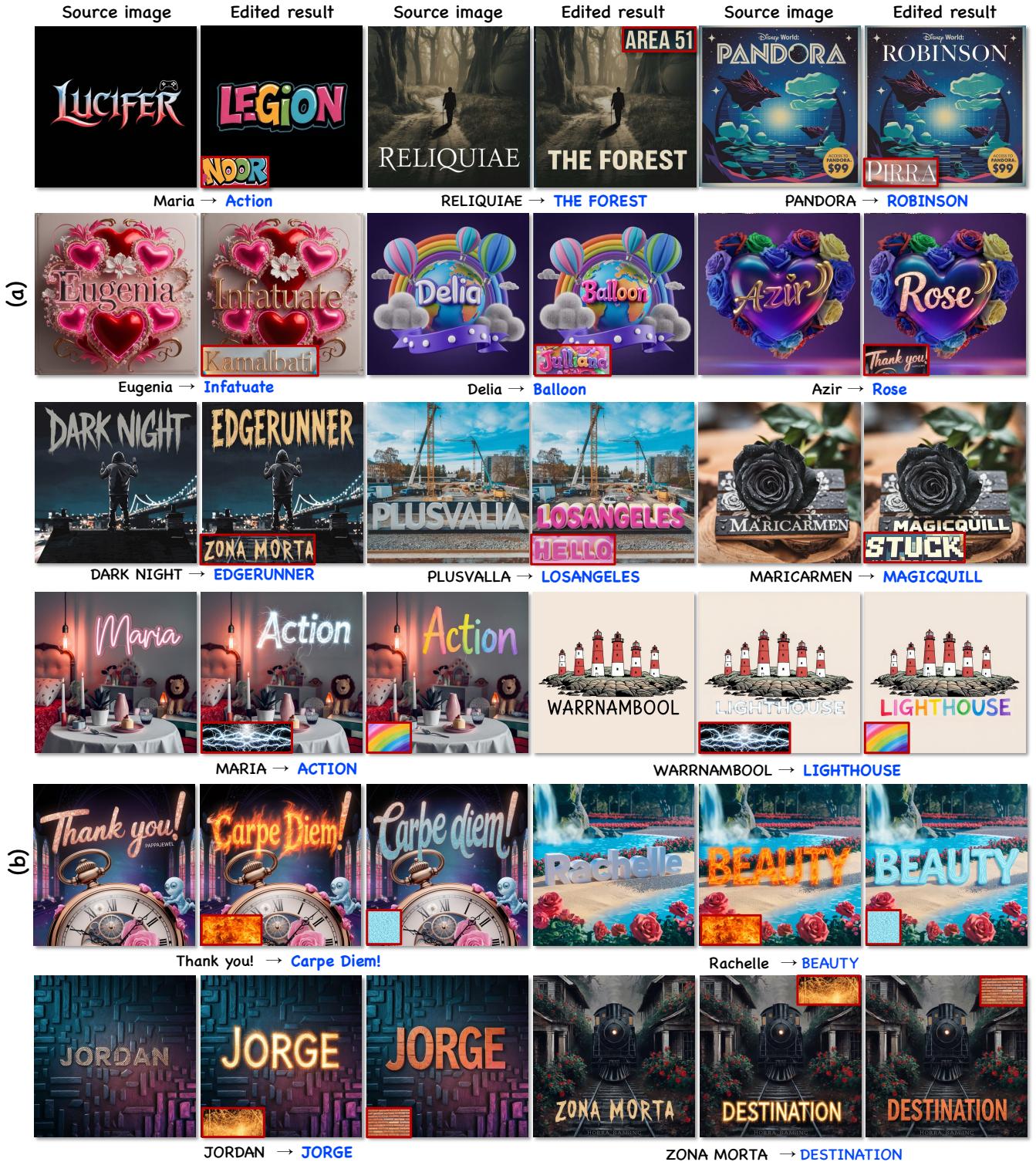


Fig. 9. Cross-style customization based on (a) text reference and (b) non-text reference images.