

# TextDiffuser-2: Unleashing the Power of Language Models for Text Rendering

Jingye Chen<sup>\*13</sup>, Yupan Huang<sup>\*23</sup>, Tengchao Lv<sup>3</sup>, Lei Cui<sup>3</sup>, Qifeng Chen<sup>1</sup>, Furu Wei<sup>3</sup>

<sup>1</sup>HKUST <sup>2</sup>Sun Yat-sen University <sup>3</sup>Microsoft Research

qwerty.chen@connect.ust.hk, huangyp28@mail2.sysu.edu.cn, cqf@ust.hk  
{tengchaolv, lecu, fuwei}@microsoft.com



Figure 1. Text-to-image results generated by TextDiffuser-2. Alongside accurate text generation, TextDiffuser-2 offers reasonable text layouts and exhibits diversity in text style powered by the strong capability of language models.

## Abstract

The diffusion model has been proven a powerful generative model in recent years, yet remains a challenge in generating visual text. Several methods alleviated this issue by incorporating explicit text position and content as guidance on where and what text to render. However, these methods still suffer from several drawbacks, such as limited flexibility and automation, constrained capability of layout prediction, and restricted style diversity. In this paper, we present TextDiffuser-2, aiming to unleash the power of language models for text rendering. Firstly, we fine-tune a large language model for layout planning. The large language model is capable of automatically generating keywords for text rendering and also supports layout modification through chatting. Secondly, we utilize the language model within the diffusion model to encode the position and texts at the line level. Unlike previous methods that employed tight character-level guidance, this approach generates more diverse text images. We conduct extensive experiments and incorporate user studies involving human participants as well as GPT-4V, validating TextDiffuser-2’s capacity to achieve a more rational text layout and generation with enhanced diversity. The code and model will be available at <https://aka.ms/textdiffuser-2>.

## 1. Introduction

In recent years, diffusion models [16, 20, 42, 43, 47, 56, 58] have successfully revolutionized the field of image synthesis, outperforming earlier methods based on GANs [14, 38] and VAEs [23, 41] in terms of fidelity and diversity. Despite showcasing impressive performance, most existing diffusion models still fall short in rendering visual text. Specifically, existing diffusion models often generate unintended symbols or artifacts during the text rendering process [10], which significantly impairs the visual quality of the generated images. Notably, text is ubiquitous in daily life, encompassing logos, banners, book covers, newspapers, etc. In this case, how to generate images with accurate, visually appealing, and coherent visual text is a crucial problem.

Through investigation, there has been a few research works focusing on visual text rendering [1, 4, 11, 30, 34, 44, 52]. Some works [1, 11, 30, 44] validate that using powerful language models [40, 51] as text encoders benefits the text rendering process. Nevertheless, they lack controllability since users may request to place text in a specific position. In this context, several works utilize explicit text position and content guidance, such as single-line segmentation masks and glyph images for GlyphDraw [34], glyph images with multiple text lines for GlyphControl [52], as well as character-level segmentation masks for TextDiffuser [4]. Although showing impressive rendering accuracy, we have noticed several drawbacks in these meth-

\*Work done during internship at Microsoft Research.

ods: (1) **Limited flexibility and automation.** GlyphControl [52] needs users to design glyph images to provide layout guidance, while GlyphDraw [34] and TextDiffuser [4] rely on the manual specification of keywords. These requirements hinder the direct conversion of natural user prompts into corresponding images, thereby narrowing the flexibility and automation capabilities; (2) **Constrained capability of layout prediction.** GlyphDraw [34] can only render images with a single text line, constraining its applicability for scenarios involving multiple text lines. For TextDiffuser [4], the produced text layouts are not visually appealing, which is primarily attributed to the limited capability of the Layout Transformer; (3) **Restricted style diversity.** For TextDiffuser [4], the utilization of character-level segmentation masks as control signals implicitly imposes constraints on the position of each character, thereby restricting the diversity of text styles and posing challenges when rendering handwritten or artistic fonts.

Given these observations, we introduce TextDiffuser-2 in this paper, taking advantage of two language models for text rendering (samples are shown in Figure 1). Firstly, we **tame a language model into a layout planner** to transform user prompt into a layout using the caption-OCR pairs in the MARIO-10M dataset [4]. The language model demonstrates flexibility and automation by inferring keywords from user prompts or incorporating user-specified keywords to determine their positions. In addition, through chatting, users can guide the language model to alter the layout, such as regenerating, adding, or moving keywords. Secondly, we **leverage the language model in the diffusion model as the layout encoder** to represent the position and content of text at the line level. Contrary to prior methods that utilized tight character-level guidance, this approach enables diffusion models to generate text images with broader diversity. Through comprehensive experiments and user studies that engaged both human participants and GPT-4V, we validate that TextDiffuser-2 can generate reasonable and visually pleasing text layouts, and it enhances the style diversity of the generated text. We will release the code and model to promote future research.

## 2. Related work

**Visual text rendering.** Despite the significant advancements in diffusion models [16, 20, 42, 56], the generation of visual text rendering remains a persistent challenge. The advancement in visual text rendering will significantly enhance the efficiency of designers in executing text-related creative tasks, such as logo or poster design. Some work [1, 30, 44] leverage the large language models [40, 51] to enhance the spelling capabilities of generative models. The other line of works [4, 34, 52] attempts to explicitly control the position and content of the text to be rendered. For instance, GlyphDraw [34] comprises two diffusion models,

including one for single-line text position prediction and the other one for image rendering guided by glyph images. GlyphControl [52] utilizes glyph images with multiple text lines as prior to guide diffusion models render accurate and coherent text. Notably, TextDiffuser [4] uses character-level segmentation masks as more fine-grained signals for better rendering. Besides, TextDiffuser is a versatile model that can tackle text-to-image, text-to-image with template, and text inpainting tasks. TextDiffuser-2 falls in the latter line but distinguishes them by employing a language model for layout planning and using another language model to encode line-level text information to enable diverse rendering.

**Language model for layout generation.** Layout generation [17, 22, 24] has a wide range of applications, including document formatting [18, 35], screen UI design [12], and image synthesis [13, 26]. Previous methods [22, 24] usually model layout generation as a regression task, representing bounding boxes using continuous coordinates. Recent advancements, such as Pix2Seq [6, 7], have explored alternative methods by treating coordinates as discrete language tokens. Another representative work, LayoutGPT [13], carefully designs prompts to guide GPT-4 [15] generating formatted layout information to assist in image synthesis. Recently, some multimodal large language models [3, 31, 36, 53, 57, 60] have also adopted this design for grounding specific objects in images. In line with these designs, TextDiffuser-2 aims to leverage language models as layout planners for visual text rendering.

**Optical character recognition.** Images naturally serve as carriers of textual information. Optical Character Recognition (OCR) has been extensively studied in academia. Specifically, text recognition [2, 25, 46, 54, 55] and detection [19, 27, 32, 33, 61] techniques play a crucial role, aiming to locate and extract textual information and further facilitate high-level understanding tasks. Our method leverages caption-OCR pairs [4] to fine-tune a large language model for generating visual text layouts. Additionally, we employ OCR tools to conduct a comprehensive evaluation.

## 3. Methodology

The architecture of TextDiffuser-2 is depicted in Figure 2, where the language model  $M_1$  and the diffusion model are trained in two stages. We introduce the role of two language models, including a language model for layout planning and another language model for layout encoding. We focus on introducing the text-to-image process, while the functions of text-to-image with templates and text inpainting will be introduced in the experiment section.

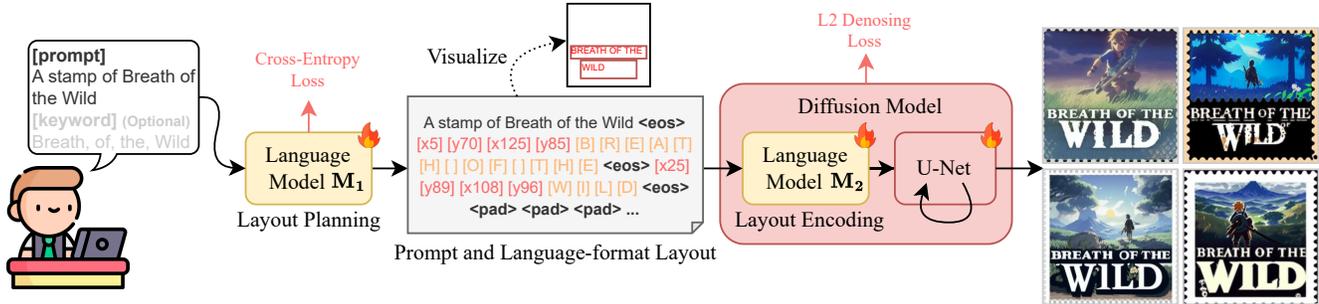


Figure 2. The architecture of TextDiffuser-2. The language model  $M_1$  and the diffusion model are trained in two stages. The language model  $M_1$  can convert the user prompt into a language-format layout and also allows users to specify keywords optionally. Further, the prompt and language-format layout is encoded with the trainable language model  $M_2$  within the diffusion model for generating images.  $M_1$  is trained via the cross-entropy loss in the first stage, while  $M_2$  and U-Net are trained using the denoising L2 loss in the second stage.

### 3.1. Language model for layout planning

Recent research has revealed that benefiting from the extensive training data across various domains, large language models [15, 49, 50] exhibit expertise beyond the language domain, such as layout planning [13, 28]. Inspired by this, we try to tame a large language model into a layout planner.

Specifically, we seek to fine-tune a pre-trained large language model  $M_1$ , which functions as a decoder, using caption-OCR pairs. As demonstrated in Figure 2, we consider two options: (1) If users do not explicitly provide keywords, the language model should infer the text and layout to be drawn on the image; (2) If users provide keywords (marked in gray color), the language model only needs to determine the corresponding layout for the keywords. Specifically, the input follows the format “[description] Prompt: [prompt] Keywords: [keywords]”<sup>1</sup>. For the output, we expect each line to follow the format “*textline*  $x_0, y_0, x_1, y_1$ ”, where  $(x_0, y_0)$  and  $(x_1, y_1)$  represent the coordinates of the top-left corner and bottom-right corner, respectively. We optimize the language model with cross-entropy loss, training simultaneously for scenarios with and without keywords. We use all the text detected in the OCR results as keywords to formulate the input. Moreover, we expect the fine-tuned language model can be guided to alter the generated layout through chatting. We will delve into this aspect in the discussion section.

### 3.2. Language model for layout encoding

Based on the layouts generated by  $M_1$ , we leverage the latent diffusion models [42] for image generation. Different from TextDiffuser [4] which incorporates text information

using segmentation masks and GlyphControl [52] which duplicates backbone parameters to accommodate the glyph image conditions, we introduce a simple and parameter-free strategy by combining the prompt and the layout for the language model  $M_2$ , *i.e.*, the text encoder within the latent diffusion model. In contrast to character-level segmentation masks that regulate the position of individual characters, the line-level bounding box offers greater flexibility during generation and does not constrain the diversity of styles.

Previous work [30] demonstrates that fine-grained tokenization can enhance the spelling capability of diffusion models. Inspired by this, we design a hybrid-granularity tokenization method that not only improves the spelling capability of the model but also keeps the sequence from getting too long. Specifically, on the one hand, we maintain the original BPE tokenization method [45] for the prompt. On the other hand, we introduce new character tokens and decompose each keyword into the character-level representation. For example, the word “WILD” is decomposed into tokens “[W]”, “[I]”, “[L]”, “[D]”. Additionally, we introduce new coordinate tokens to encode the position. For instance, the tokens “[x5]” and “[y70]” correspond to an x-coordinate of 5 and a y-coordinate of 70, respectively. Each keyword information is separated by the end-of-sentence token “<eos>”, and any remaining space to the maximum length  $L$  will be filled with padding tokens “<pad>”. We train the whole diffusion model, including the language model  $M_2$  and U-Net, using the L2 denoising loss.

## 4. Experiments

**Implementation details.** For *layout planning*, we fine-tune the vicuna-7b-v1.5 [8] model based on the FastChat framework [59]. The caption-OCR pairs are derived from the MARIO-10M dataset [4], and we use 5k samples for fine-tuning. We normalize the positions to the range of 0~128 to increase the compactness of the coordinate feature space. The learning rate is set to  $2e-5$ , and we conduct

<sup>1</sup>Task description: Given a prompt that will be used to generate an image, plan the layout of visual text for the image. The size of the image is 128x128. Therefore, none of the properties of the positions should exceed 128, including the coordinates of the top, left, right, and bottom. You don’t need to specify the details of font styles. At each line, the format should be textline left, top, right, and bottom. So let us begin.

#Data	Acc $\uparrow$	Pre $\uparrow$	Rec $\uparrow$	F $\uparrow$	IOU $\downarrow$
0k-2shot	49.65	84.18	69.69	76.25	19.69
2.5k	61.10	82.20	85.18	83.67	<b>3.21</b>
5k	<b>64.85</b>	84.98	<b>86.38</b>	<b>85.67</b>	3.25
10k	<b>64.85</b>	84.38	86.23	85.29	4.27
50k	63.72	<b>85.32</b>	85.78	85.55	3.68
100k	62.87	85.26	85.98	85.62	4.31

Table 1. Ablation studies on the amount of fine-tuning data. The “0k-2shot” setting denotes the use of two examples for few-shot learning, without any additional fine-tuning. When using 5k data, the language model  $M_1$  performs better. The percentage sign is omitted, as is consistent with the following tables. ‘Pre’, ‘Rec’, and ‘F’ denote precision, recall, and f-measure, same as follows.

Representation	Acc $\uparrow$	Pre $\uparrow$	Rec $\uparrow$	F $\uparrow$
Center (Char)	35.19	61.75	62.71	62.23
LT (Char)	28.32	54.94	55.64	55.29
LT+RB (Subword)	15.48	41.74	42.53	42.13
LT+RB (Char)	<b>57.58</b>	<b>74.02</b>	<b>76.14</b>	<b>75.06</b>

Table 2. Ablation studies on the representation of coordinates and the tokenization level. ‘L’, ‘T’, ‘R’, and ‘B’ denote left, top, right, and bottom. “Char” refers to tokenizing keywords into individual characters, whereas “Subword” refers to the use of BPE for tokenizing into subwords. Using the top-left and bottom-right corners and character-level tokenization achieves better performance.

a total of 6 epochs of fine-tuning with a batch size of 256. It takes one day to train with 8 A100 GPU cards. During the inference stage, when using a single A100 GPU card, the average time to generate a layout for each prompt is 1.1 seconds. For *layout encoding*, we utilize SD 1.5 [42] and use the built-in CLIP text encoder with base size [39]. The whole model consists of 922M parameters. We incorporate special tokens, including 256 coordinate tokens and 95 character tokens. The alphabet contains 26 uppercase and 26 lowercase letters, 10 numbers, 32 punctuation marks, and a space. The size of the input image is  $512 \times 512$ . The model is trained for 6 epochs on the MARIO-10M dataset [4] with a learning rate of  $1e-4$  and a batch size of 576. The maximum length  $L$  is set to 128. More details about the choice of  $L$  are shown in Appendix A. It takes one week to train the whole diffusion model with 8 A100 GPU cards. When sampling with 50 steps, the generation for a single image costs 6 seconds.

#### 4.1. Ablation studies

**How much data is needed for fine-tuning  $M_1$ ?** As illustrated in Table 1, we conduct experiments with different data amount, including 0k, 2.5k, 5k, 10k, 50k, and 100k. Particularly, in the 0k setting, we provide two examples

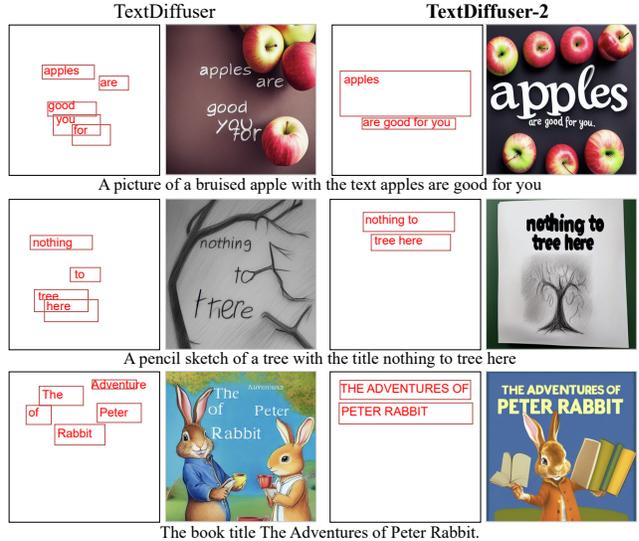


Figure 3. Visualizations of layouts. TextDiffuser-2 generates more visually pleasing and rational layouts compared with TextDiffuser.

of few-shot learning. In the absence of examples, the result often fails to conform to the appropriate format. We evaluate our approach using the MARIO-Eval benchmark [4], which consists of prompt and keyword pairs. Besides, the quotation marks in the prompt are removed for evaluation. Since the LAIONEval subset contains some noise in its keyword ground truth, it is unreliable for accurate assessments. So, we decided not to use it in this keyword extraction experiment. For evaluation, we use accuracy, precision, recall, and F-measure to assess the model’s ability to extract keywords. Additionally, we introduce an IoU metric to measure the maximum IoU value between the generated boxes for each sample (only those samples with more than one predicted box are calculated). The experimental results showcase that the model achieved optimal performance in the majority of metrics when fine-tuned with 5k data, and we visualize some samples in Figure 3. We notice that the language model exhibits flexibility in generating keywords, such as determining the case of the keyword, or introducing appropriate words beyond the provided prompt. More samples are in Appendix B. In subsequent experiments, we employ the model fine-tuned on 5k data for layout planning.

**How to represent the position of text lines?** Apart from utilizing the top-left and bottom-right corners to represent a text line, we also investigate alternative single-point representations, such as employing the top-left point or the center point. Intuitively, using a single point to represent a text line provides more flexibility, enabling the generated text to exhibit greater diversity in angles and sizes. In Appendix C, visualizations are shown to validate the diversity of the generation using single-point conditions. However, as shown

Metrics	SD-XL [37]	PixArt- $\alpha$ [5]	GlyphControl [52]	TextDiffuser [4]	TextDiffuser-2
<b>Quantitative Results</b>					
FID $\downarrow$	62.54	87.09	50.82	<u>38.76</u>	<b>33.66</b>
CLIPScore $\uparrow$	31.31	27.88	<b>34.56</b>	34.36	<u>34.50</u>
OCR (Accuracy) $\uparrow$	0.31	0.02	32.56	<u>56.09</u>	<b>57.58</b>
OCR (F-measure) $\uparrow$	3.66	0.03	64.07	<b>78.24</b>	<u>75.06</u>
<b>User Studies by Humans / GPT-4V</b>					
Layout Aesthetics $\uparrow$	-	-	-	<u>28.43 / 0.00</u>	<b>71.57 / 100.00</b>
Style Diversity $\uparrow$	-	-	<u>31.37 / 33.33</u>	27.45 / <b>33.33</b>	<b>41.18 / 33.33</b>
Text Quality $\uparrow$	14.58 / 7.69	3.65 / 0.00	21.35 / 15.38	<u>23.44 / 30.77</u>	<b>36.98 / 46.15</b>
Text-Image Matching $\uparrow$	7.14 / 0.00	3.30 / 0.00	<u>29.67 / 18.18</u>	19.23 / <u>36.36</u>	<b>40.66 / 45.45</b>
Inpainting Ability $\uparrow$	-	-	-	<u>25.49 / 33.33</u>	<b>74.51 / 66.67</b>

Table 3. Demonstration of the quantitative results and user studies. We also incorporate GPT-4V [15] into the user studies. The best and second-best results are indicated in bold and underlined formats. TextDiffuser-2 achieves the best results under the majority of metrics.

in Table 2, we notice that there is a considerable decline in the OCR accuracy of the single-point representation on the MARIO-Eval benchmark [4]. For example, compared with the LT-RB setting, the accuracy of the center and LR settings declined by 22.39% and 29.26%. Hence, we leverage the top-left and bottom-right corners to represent the box in the following experiments. We also explore the inclusion of angle information in Appendix D.

#### Should text be tokenized at the character or subword level?

We also explore Byte Pair Encoding (BPE) to tokenize keywords into the subword level. As shown in Table 2, we observe that using subword-level tokenization significantly underperforms character-level representation, *i.e.*, it is lower by 42.1% on the accuracy metric. When using subword-level tokenization, the model becomes insensitive to the spelling of each token, which poses significant challenges to the text rendering process.

## 4.2. Experimental results

**Quantitative results.** As shown in Table 3, we conduct quantitative experiments on the MARIO-Eval benchmark [4]. For comparisons, we leverage two state-of-the-art text-to-image models including SD-XL [37] and PixArt- $\alpha$  [5], and two models incorporating specific guidance for generating text images including TextDiffuser [4] and GlyphControl [52]. Details of these compared methods are shown in Appendix E. For all methods, we employ 50 sampling steps and set the classifier-free guidance to 7.5. The experimental results demonstrate that TextDiffuser-2 outperforms other methods in terms of the FID evaluation metric. Having not been specifically trained on text images, SD-XL and PixArt- $\alpha$  exhibit a larger divergence against the ground truth, resulting in higher FID and lower CLIP scores. For the OCR metrics, it is observed that only models incorporating guid-

ance can effectively render text. Furthermore, TextDiffuser-2 outperforms GlyphControl and has an OCR performance comparable to TextDiffuser. It is noteworthy that the TextDiffuser renders text in a standardized font (see Figure 5), thereby reducing the complexity of the rendering process. This strategy sacrifices font style diversity to enhance the accuracy of text rendering. By contrast, while maintaining the ability to generate accurate text, TextDiffuser-2 can generate text with a greater diversity.

**Qualitative results.** The visualizations are demonstrated in Figure 4. We compare our method with SD-XL [37], PixArt- $\alpha$  [5], Ideogram [21], DALLE-3 [9], GlyphControl [52] and TextDiffuser [4]. Although some of the latest text-to-image models (*e.g.*, DALLE-3 and PixArt- $\alpha$ ) showcase superior image quality, they do not perform as well in rendering text compared with models that incorporate explicit guidance. Compared to TextDiffuser, our method generates more aesthetically pleasing layouts, avoiding misalignment or discordant font sizes. Furthermore, since TextDiffuser-2 utilizes the more flexible line-level guidance, it offers better control over font style, such as when rendering the handwritten style “deep learning”. In contrast, TextDiffuser, which uses character-level guidance, can mainly render rigid font styles. We offer some visualizations in Figure 5. For instance, when rendering “Winter”, our method demonstrates greater diversity in terms of perspective angle and font style compared to other methods. In addition, we adopt the same layout to validate the performance of GlyphControl, which also uses line-level guidance. We observe that TextDiffuser-2 achieves a higher accuracy than GlyphControl. We also conduct comparisons with some methods that are neither open-source nor offer APIs, such as GlyphDraw [34] and Character-Aware Model [30] using the samples shown in their corresponding papers in Appendix F.

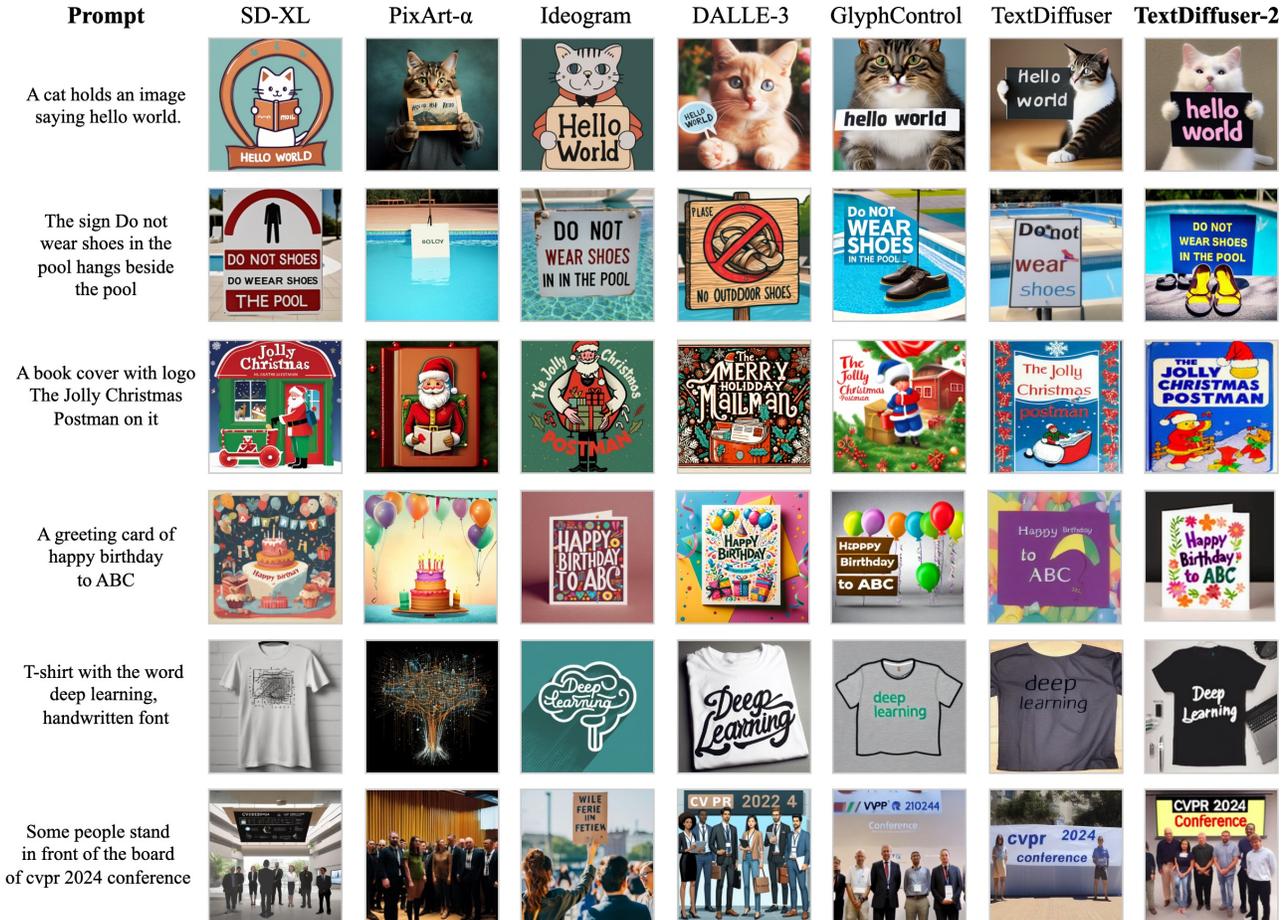


Figure 4. Visualizations of text-to-image results compared with existing methods. TextDiffuser-2 can automatically extract keywords from prompts for accurate rendering. Additionally, the fonts generated by TextDiffuser-2 exhibit a wide range of diversity.

**User studies.** As shown in Table 3, we design questions covering five aspects: layout aesthetics, style diversity, text quality, text-image matching, and inpainting ability, each of which contains 6, 3, 6, 6, 6 questions. We involved a total of 17 human participants in our study. Additionally, we employ GPT-4V [15] to carry out the user study. We devised instruction guidance for each task, prompting GPT-4V to think step by step to arrive at the final answer. Ultimately, each method’s score is calculated as the number of votes it receives divided by the total number of votes. Based on the results, TextDiffuser-2 has achieved optimal performance in four out of five metrics in studies involving human participation and GPT-4V. More details are in Appendix G.

**More applications.** In addition to text-to-image, we also explore other applications of TextDiffuser-2. (1) *Text-to-image generation with template.* When a template image (e.g., printed, handwritten, or scene text image) is provided, TextDiffuser-2 can use existing OCR tools to extract text information and directly feed it into the diffusion model as the

condition, eliminating the need for layout prediction from the language model  $M_1$ . We showcase some samples in Figure 8. (2) *Text inpainting.* Similar to TextDiffuser, the architecture of TextDiffuser-2 adapts well for training on text inpainting tasks. We only need to modify the channel of the input convolution kernel in the U-Net. Specifically, we augment the original 4-dimensional latent feature with 5 additional dimensions, including 4 dimensions of non-inpaint area features and 1 dimension for the mask. Moreover, only the text position and content from the inpaint area are required as conditions for the diffusion model. More details are shown in Appendix H. We compare the performance with TextDiffuser, and the experimental results are shown in Figure 6. It is worth noting that TextDiffuser requires a text mask as a condition to specify the position of each character, which can be cumbersome in practical applications. Additionally, the text mask may limit the style of the generated results. For example, when rendering the word “Curve”, the generated result cannot produce a visually curved effect due to the constraints of the character-



Figure 5. Visualization of diversity in generating multiple images under the same prompt. TextDiffuser-2 is capable of generating more artistic fonts, with increased diversity in the positioning of characters and the inclination angle of text lines.

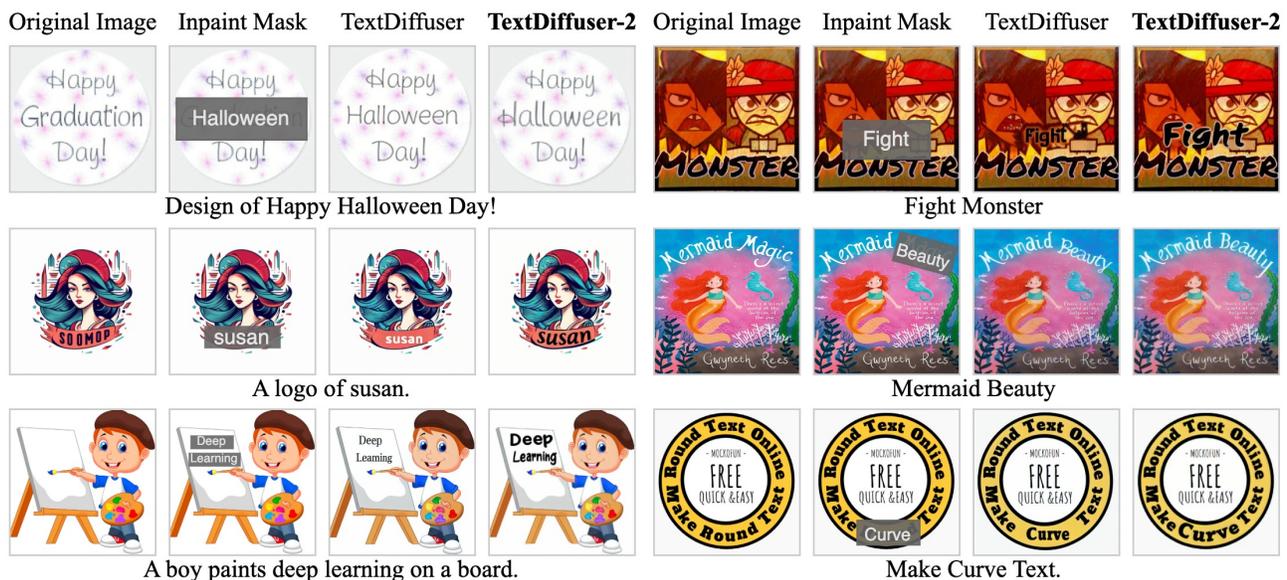


Figure 6. Visualizations of the text inpainting task compared with TextDiffuser. TextDiffuser-2 can generate more coherent text.

level segmentation mask. In contrast, the inpainting process of TextDiffuser-2 is more flexible, thus resulting in a better user experience. (3) **Natural image generation without text.** Given that we train TextDiffuser-2 on images with text, we are curious about its capability to generate images without text. Specifically, by omitting the text position and content guidance, TextDiffuser-2 can generate images without text. We randomly select 10,000 prompts from the Microsoft COCO dataset [29] for generation and compare the results with those generated by SD 1.5 [42]. The visualization results are shown in Figure 9. Although TextDiffuser-2 is fine-tuned on domain-specific data, it still maintains its generative capabilities in the original domain. We calculate

the FID score of the generated results. When the sampling steps are set to 50, and the classifier-free guidance is set to 7.5, TextDiffuser-2’s FID score is 24.06, versus 23.03 for SD 1.5. Although the FID score slightly increases, the overall difference is not significant according to the visualizations.

### 4.3. Discussions

**Operating layout through multi-round chat.** Since the language model used to generate the layout is fine-tuned based on a chat model, we are curious whether we can manipulate the layout through multi-round chat. We demonstrate the results in Figure 7. Experimental results showcase

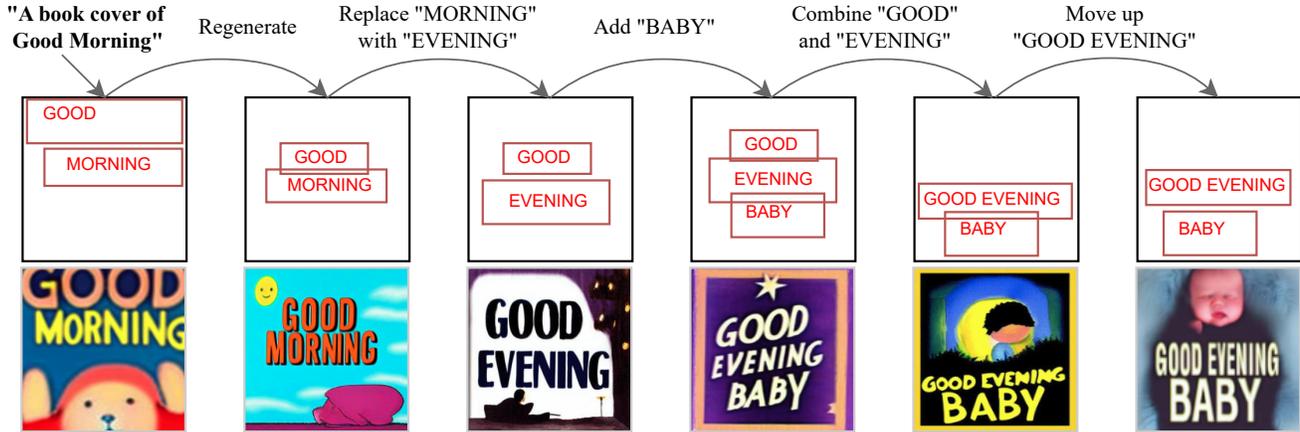


Figure 7. Visualizations of operating layout through multi-round chat and the corresponding images.



Figure 8. Visualizations of the text-to-image with template task.



Figure 9. Visualizations of generating images without text.

that through interactive conversation,  $M_1$  can not merely regenerate the layout, but also possess the ability to add or modify keywords, as well as manipulate the location of the box. This further enhances the flexibility and versatility of the proposed TextDiffuser-2.

**Generation based on overlapping layouts.** Occasionally, we notice that there exist overlapping boxes during the layout prediction stage. We present TextDiffuser-2, as well as the results generated by GlyphControl and TextDiffuser using overlapping layouts in Appendix I. Experimental results indicate that TextDiffuser-2 demonstrates greater robustness towards overlapping boxes. Conversely, the results generated by the other two methods will produce scrambled text, thereby impacting the overall quality of the image.

## 5. Conclusion

In this paper, we introduce TextDiffuser-2, aiming to unleash the power of language models for the text rendering

task. Specifically, we attempt to tame two language models, one for layout planning and the other for layout encoding. Experimental results validate that TextDiffuser-2 is capable of generating more diverse images while maintaining the accuracy of the generated text. For the *limitation*, TextDiffuser-2 faces challenges when rendering complex languages, as it expands the renderable character table by adding new tokens. For instance, when rendering Chinese text, TextDiffuser-2 faces difficulties due to the extensive character set, potentially leading to few-shot or even zero-shot scenarios. For the *broader impact*, TextDiffuser-2 can be used to enhance creativity in fields like graphic design, advertising, and art. Besides, it can be used to generate informative images for teaching and learning. For example, it could create diagrams with explanatory text for educational materials. However, if used maliciously, TextDiffuser-2 could be employed to create images containing false text information. For *future work*, we seek to explore the rendering of characters in multiple languages and enhance the resolution of generated text images.

# Appendix

## A. Choice of the maximum length $L$

During the training process, the composed sequence (*i.e.*, the prompt combined with text content and position) has a maximum length limit. As shown in Figure 10, by analyzing the MARIO-10M dataset [4], we notice that the composed sequence for 94.0% of the samples is less than 128 in length, and all samples are below this threshold during the evaluation. Obviously, we can increase the maximum length to a larger value, such as 256, accommodating 99.2% of training samples. However, enlarging the length limit would also result in increased computational costs, such as raising the single sample inference time from 6 to 7 seconds. Therefore, the choice of length limit should be made based on practical considerations, balancing between the model’s capability and efficiency.

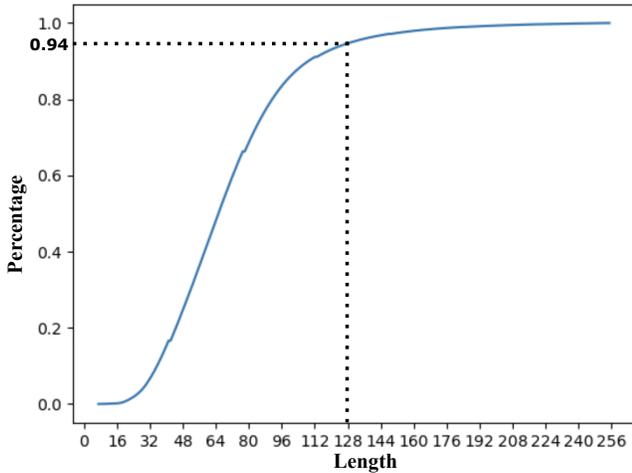


Figure 10. The cumulative distribution function to analyze the length of the composed sequences. When setting the maximum length  $L$  at 128, the vast majority of samples (94%) are covered.

## B. More visualizations of layout prediction

As depicted in Figure 11, we showcase more layout prediction results. We specify the keywords to be rendered in the first two rows. The language model has the capacity to organize the specified keywords, placing related keywords in the same line and generating aesthetically pleasing layouts. Notably, the final row of predictions includes words not present within the prompt. For instance, the model substitutes “200g” for “200gram”. It is a logical substitution given that both terms convey the same meaning. Additionally, the model replaces the misspelled term “RRAINBOW” in the prompt with the correct term “RAINBOW”. This further showcases the flexibility of the layout planner  $M_1$ .



Figure 11. More visualizations of the layout predictions. The specified keywords are marked in blue color in the first two rows.

## C. Generation guided by single-point condition

We retrain TextDiffuser-2 and implement a single-point supervision strategy during the training process, such as using the center and top-left points. As illustrated in Figure 12, despite the diversity in text size and angle generated by the single-point conditions, we observe a significant portion of the text to be inaccurate. Given the observation of a decline in accuracy over 20% (as shown in Table 2 in the main paper), we ultimately employ the top-left and bottom-right points as the condition.

## D. Generation with additional angle condition

As shown in Figure 13, we demonstrate samples generated with different angle conditions. Specifically, we retrain TextDiffuser-2 and add 181 angle tokens, ranging from  $-90^\circ$  to  $90^\circ$ . When constructing the language-format layout, the angle token is placed after the four coordinate tokens. The results show that the generated results align well with the angle instructions.

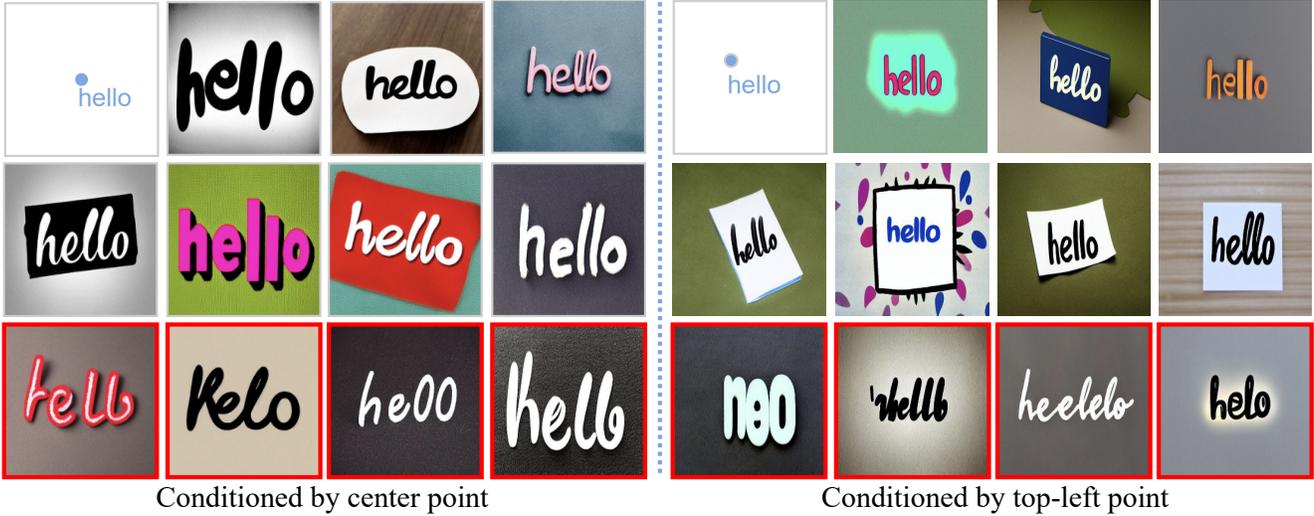


Figure 12. Visualizations of generation guided by single-point conditions, including the center point and the top-left point. The prompt is “A text image of hello”. The samples highlighted by red boxes in the last row denote the rendered text is incorrect.

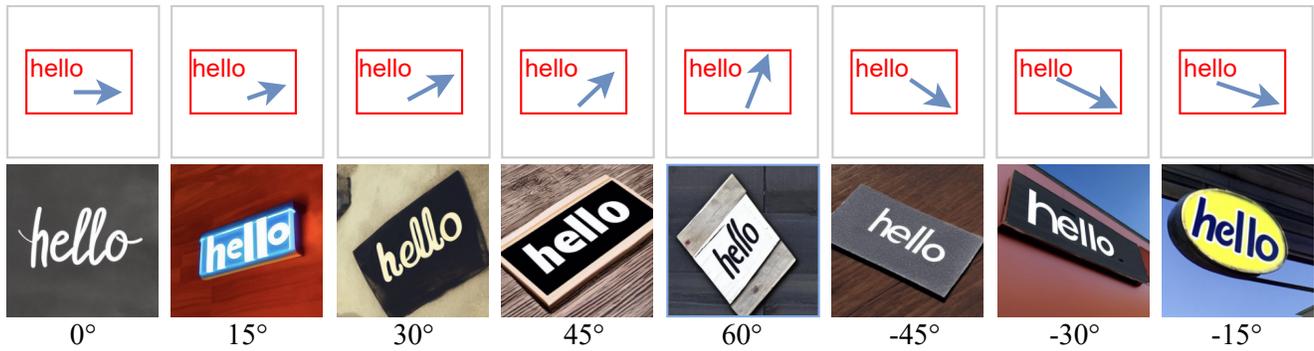


Figure 13. Visualizations of generation with different angle guidance. The prompt is “A text image of hello”.

## E. Details of compared methods and evaluation

We introduce all the baselines and their experimental settings as follows.

**SD-XL** [37] is an improved version of the latent diffusion model [42] with stronger backbone and powerful text embedding. SD-XL comprises 5.8B parameters and the resolution of the output images is  $1024 \times 1024$ .

**PixArt- $\alpha$**  [5] is a powerful Transformer-based text-to-image diffusion model and is training-efficient. It consists of 0.6B parameters. The output resolution is of size  $1024 \times 1024$ .

**Ideogram** [21] is an online website that can produce attractive logos, posters, and other natural images based on prompts. We use the typography mode and manually quote keywords to be rendered. The resolution is  $1024 \times 1024$ .

**DALLE-3** [9] exhibits robust text-to-image capabilities, producing images that precisely conform to the given prompt. It generates high-resolution outputs with

$1024 \times 1024$  resolution. We leverage the official API for the generation process.

**GlyphControl** [52] utilizes the framework of ControlNet [56] and the pre-trained model of SD 2.1 [42], producing the output image of size  $768 \times 768$ . It takes glyph images with multiple text lines as the condition. GlyphControl has 1.3B parameters. Specifically, since it can not generate images from prompts, we use the layouts produced by TextDiffuser-2 to make the glyph image.

**TextDiffuser** [4] is a two-stage framework that can convert user prompts into images. It relies on users to specify keywords for rendering. TextDiffuser is pre-trained based on SD 1.5 [42], and the resolution of the generated images is  $512 \times 512$ . It consists of 884M parameters in total.

For evaluation, we utilize the metrics employed in TextDiffuser [4] and also use Microsoft Read API to evaluate the OCR performance.



Figure 14. Comparisons with Character-Aware Model (CA Model) [30] and GlyphDraw [34] using samples in their papers.

## F. Comparisons with samples in other papers

Since the source code, pre-trained weight, or demo is not available for Character-Aware Model [30] and GlyphDraw [34], we conduct comparisons with samples in their corresponding papers. As demonstrated in Figure 14, we visualize four samples for each compared method. Notably, TextDiffuser-2 shows better rendering accuracy compared with the Character-Aware Model, which contains several typos, including the missing “r” in “from” and the incorrect spelling of “Chimpanzees”. Besides, the Character-Aware Model enhances visual text rendering by utilizing language models with a larger parameter size (*e.g.*, T5-XXL [40] with 11B parameters). We have demonstrated that even with a smaller-scaled CLIP text encoder with 63M parameters, superior text rendering performance can be achieved by virtue of explicit positional and content supervision. Besides, TextDiffuser-2 outperforms GlyphDraw as TextDiffuser-2 can render images with multiple text lines.

## G. More details about user studies

We conduct comprehensive user studies on five aspects, including layout aesthetics, style diversity, text quality, text-image matching, and inpainting ability. The details of the questions are displayed in Figure 18. In addition to human involvement, we incorporate GPT-4V [15] in our user studies. Specifically, we design prompts to encourage GPT-4V to proceed step-by-step, deriving the final answer through logical analysis. The dialogue record is shown in Figure 19 and Figure 20. It suggests that GPT-4V exhibits impressive literacy skills, and its logical chain is reasonable.

## H. More details about text inpainting

Similar to TextDiffuser [4], by appending another five-dimension feature, including the one-dimension mask and

four-dimension non-inpainted area features, to the input of U-Net, TextDiffuser-2 can be trained for the text inpainting task. Specifically, 14,400 parameters will be added, which accounts for a small proportion of the whole architecture containing 922M parameters. We set the classifier-free guidance to 7.5 and used 50 sampling steps, which cost 6 seconds for generation using one A100 GPU card.

## I. Generation based on overlapping layouts

The results of the layout predictor will inevitably contain overlapping boxes. In this section, we analyze the robustness of three methods, including GlyphControl [52], TextDiffuser [4], and the proposed TextDiffuser-2. In terms of explicit guidance, GlyphControl employs glyph images, TextDiffuser uses character-level segmentation masks, and TextDiffuser-2 harnesses bounding boxes with corresponding text. The visualization results are demonstrated in Figure 15. The results reveal that the proposed TextDiffuser-2 is more robust when using overlapping layouts for a generation. By contrast, GlyphControl and TextDiffuser will generate incorrect text, resulting from the occluded glyph images and segmentation masks.

## J. Conformity of the positional guidance

To verify whether the text generated in the images adheres to positional guidance, we employ a widely-used text detection tool [48] for assessment. The precision, recall, and F-measure are 0.9524, 0.9635, and 0.9579, respectively. This indicates that in the vast majority of cases, the generated text is able to comply with the provided positional constraints.



Figure 15. Comparative visualizations of generation results using overlapping layouts. TextDiffuser-2 demonstrates enhanced robustness compared with other methods.



Figure 16. Comparative visualizations of TextDiffuser-2 with different versions of the Stable Diffusion model. Utilization of SD 2.1 exhibits improved detail rendering and a more accurate depiction of small-scale characters compared to SD 1.5.

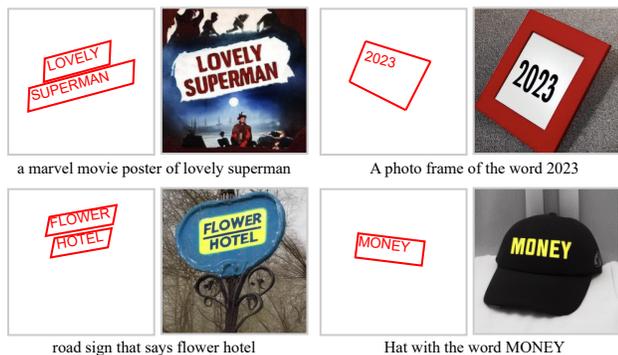


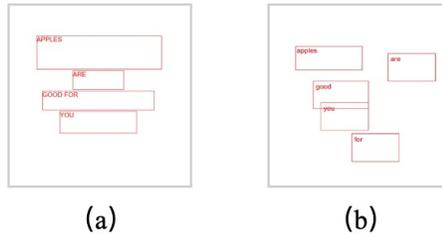
Figure 17. Demonstration of TextDiffuser-2’s generation guided by quadrilateral bounding boxes, showcasing the model’s ability to align text accurately within the specified geometrical constraints.

## K. TextDiffuser-2 based on SD 2.1

TextDiffuser-2 can be trained based on different pre-trained models. In Figure 16, we visualize the comparisons of results based on SD 1.5 and SD 2.1. We notice that the model based on SD 2.1 provides more details thanks to the stronger power of the pre-trained model. Meanwhile, it can correctly render characters with small sizes because the resolution of the latent space is higher.

## L. Generation guided by quadrilateral boxes

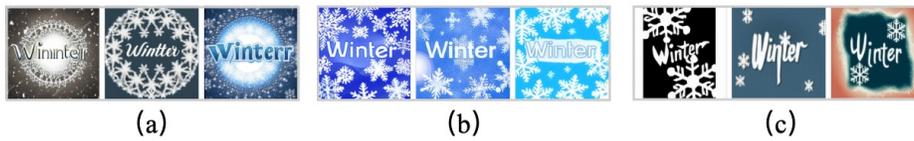
In addition to using horizontal boxes to provide positional information, we remain curious whether TextDiffuser-2 can be guided by quadrilateral boxes, which could more accurately describe slanted text. To investigate this, we make two modifications. First, we train a layout planner M1 to output each line in the format of “*textline*  $x_0, y_0, x_1, y_1, x_2, y_2, x_3, y_3$ ”. Secondly, we adapt the layout encoder M2 to encode this sequence. We set the maximum length limit  $L$  to 256 to accommodate longer input sequences. Visualizations are shown in Figure 17. We notice that the generated results align well with the guidance of quadrilateral boxes. For future work, we plan to use more control points to represent the boxes, allowing for rendering more artistic text.



Prompt: a picture of a bruised apple with the text apples are good for you in a fancy font

**Question:** The above two figures represent layout planning based on prompt. Which figure best accomplishes this task?

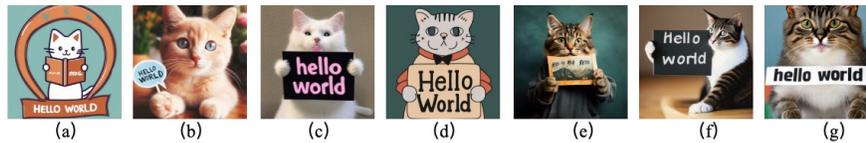
### 1. Layout Aesthetics



Prompt: a logo of Winter in artistic font, made by snowflake

**Question:** Which set of images exhibits a richer style and more diverse text?

### 2. Style Diversity



Prompt: a cat holds an image saying hello world

**Question:** Which of the following images has the best text rendering quality? (Multiple Choice)

**Question:** Which of the following images best matches the text description? (Multiple Choice)

### 3. Text Quality

### 4. Text-Image Matching



Prompt: Design of Happy Halloween Day!

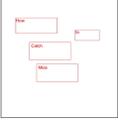
**Question:** Which image exhibits better result of text inpainting?

### 5. Inpainting Ability

Figure 18. The user studies comprise five questions. Specifically, users should make a single selection for questions pertaining to layout aesthetics, style diversity, and inpainting ability. For questions regarding text quality and text-image matching, users have the freedom to select multiple options.

Let's conduct a user study. I will show you an image and a question each time. Please start your response with 'Let's think step by step.' and end with 'Therefore, the answer is'. Please treat each case individually and do not be influenced by previous cases. If you have understood the requirements. Please answer "Yes" and we will start.

 Yes

(a) (b)

Prompt: A cat is reading a book with the title 'How to catch mice'

The above two figures represent layout planning based on prompt. Which figure best accomplishes this task? (Single Choice)

 Let's think step by step. The prompt given is "A cat is reading a book with the title 'How to catch mice'." We have two figures, (a) and (b), that are meant to represent layout planning based on this prompt. Figure (a) shows a more consolidated layout where the title is in a single block, which could represent the cover of a book. Figure (b), however, breaks the title into separate blocks which seems more disjointed and less representative of a book title. For a book cover layout, which is what the prompt suggests, having a cohesive and clear title is important. Therefore, the answer is figure (a) best accomplishes this task.










(a) (b) (c) (d) (e) (f) (g)

Prompt: A cat holds an image saying hello world.

Which of the following images has the best text rendering quality? (Multiple Choice) Tips: Focus more on the accuracy of the text, and whether the text is coherent with backgrounds

 Let's think step by step. We are looking for the image(s) with the best text rendering quality, focusing on the accuracy of the text and its coherence with the background. · Image (a) has clear and coherent text, but it's part of an illustration rather than a real-world image. · Image (b) shows text on a physical object the cat is holding. The text is clear and seems coherent with the lighting and perspective. · Image (c) also shows clear text on a sign the cat is holding. The text is coherent with the sign's perspective and the cat's pose. · Image (d) is similar to (a), being part of an illustration, with clear text. · Image (e) shows a cat with a sign, but the text is obscured and distorted, making it less accurate than others. · Image (f) has a clear text on a sign held by a cat, with accurate perspective and lighting. · Image (g) has clear text, but it appears to be superimposed on the image without correct perspective. Considering the accuracy and coherence of the text with the background, the best text rendering quality appears to be in images (b), (c), and (f). Therefore, the answers (B/C/F) are.










(a) (b) (c) (d) (e) (f) (g)

Prompt: A book cover with logo The Jolly Christmas Postman on it

Which of the following images best matches the text description? (Multiple Choice)

Figure 19. The illustration of using GPT-4V [15] to conduct user studies (Pages: 1/2).



Let's think step by step. We need to determine which image(s) best match the text description: "A book cover with logo The Jolly Christmas Postman on it."  Image (a) matches the description as it shows a book cover with the logo "The Jolly Christmas Postman."  Image (b) also matches the description with a book cover that includes the logo.  Image (c) is consistent with the description, featuring a book cover with the logo.  Image (d) does not have the logo "The Jolly Christmas Postman," hence it does not match the description.  Image (e) depicts Santa Claus and does not include the logo as described.  Image (f) has "The Merry Holiday Mailman," which is not the logo specified in the description.  Image (g) matches the description with a book cover that has the logo "The Jolly Christmas Postman." The images that best match the description, featuring a book cover with the logo "The Jolly Christmas Postman," are (a), (b), (c), and (g). Therefore, the answers (A/B/C/G) are.



Which set of images exhibits a richer style and more diverse text? (Single Choice)



Let's think step by step. The term "richer style" could refer to the complexity and variety in the design, including factors such as color variety, texturing, and overall visual impact. More "diverse text" would refer to the variety in the typography, including font style, size, and perhaps the effects applied to the text. Set (a) displays neon signs with different background textures and colors but the font style remains the same across all three. Set (b) shows different styles of neon signs with varying background textures, lighting effects, and the addition of elements like an outline and a hanging sign, indicating a richer style. The text also appears to have slight variations in style. Set (c) has consistent font styles but exhibits a range of colors and a glowing effect that varies in intensity and spread. Considering the elements of design and typography, set (b) has more variations in the way the sign is presented and slight variations in the text style itself. Therefore, the answer (B) is.



Which image exhibits better result of text inpainting? (Single Choice)



Let's think step by step. We need to assess which image, (a) or (b), shows a better result of text inpainting. The original image has "MONSTER" written in stylized text, and the inpaint mask indicates that the word "Fight" has been added above it. In image (a), the word "Fight" is inpainted with a font style and size that seems to match the original text in terms of its boldness and impact. The text aligns well with the pre-existing design and maintains the graphic consistency. In image (b), the word "fight" is in a smaller, less stylized font that does not seem to match the original design's bold and impactful style. It appears less integrated with the overall design and is less prominent. Considering the consistency in style and the impact of the text with the original design, the answer (A/B) is: A.

Figure 20. The illustration of using GPT-4V [15] to conduct user studies (Pages: 2/2).

## References

- [1] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. [1](#), [2](#)
- [2] Darwin Bautista and Rowel Atienza. Scene text recognition with permuted autoregressive sequence models. In *ECCV*, 2022. [2](#)
- [3] Chi Chen, Ruoyu Qin, Fuwen Luo, Xiaoyue Mi, Peng Li, Maosong Sun, and Yang Liu. Position-enhanced visual instruction tuning for multimodal large language models. *arXiv preprint arXiv:2308.13437*, 2023. [2](#)
- [4] Jingye Chen, Yupan Huang, Tengchao Lv, Lei Cui, Qifeng Chen, and Furu Wei. Textdiffuser: Diffusion models as text painters. In *NeurIPS*, 2023. [1](#), [2](#), [3](#), [4](#), [5](#), [9](#), [10](#), [11](#)
- [5] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart- $\alpha$ : Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023. [5](#), [10](#)
- [6] Ting Chen, Saurabh Saxena, Lala Li, David J Fleet, and Geoffrey Hinton. Pix2seq: A language modeling framework for object detection. In *ICLR*, 2021. [2](#)
- [7] Ting Chen, Saurabh Saxena, Lala Li, Tsung-Yi Lin, David J Fleet, and Geoffrey E Hinton. A unified sequence interface for vision tasks. In *NeurIPS*, 2022. [2](#)
- [8] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, 2023. [3](#)
- [9] DALLE-3. Link: <https://openai.com/dall-e-3>, 2023. [5](#), [10](#)
- [10] Giannis Daras and Alexandros G Dimakis. Discovering the hidden vocabulary of dalle-2. *arXiv preprint arXiv:2206.00169*, 2022. [1](#)
- [11] DeepFloyd. Github link: <https://github.com/deep-floyd/df>, 2023. [1](#)
- [12] Biplab Deka, Zifeng Huang, Chad Franzen, Joshua Hirschman, Daniel Afegan, Yang Li, Jeffrey Nichols, and Ranjitha Kumar. Rico: A mobile app dataset for building data-driven design applications. In *UIST*, 2017. [2](#)
- [13] Weixi Feng, Wanrong Zhu, Tsu-jui Fu, Varun Jampani, Arjun Akula, Xuehai He, Sugato Basu, Xin Eric Wang, and William Yang Wang. Layoutgpt: Compositional visual planning and generation with large language models. In *NeurIPS*, 2023. [2](#), [3](#)
- [14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014. [1](#)
- [15] GPT-4. Link: <https://openai.com/gpt-4>, 2023. [2](#), [3](#), [5](#), [6](#), [11](#), [14](#), [15](#)
- [16] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *CVPR*, 2022. [1](#), [2](#)
- [17] Kamal Gupta, Justin Lazarow, Alessandro Achille, Larry S Davis, Vijay Mahadevan, and Abhinav Shrivastava. Layout-transformer: Layout generation and completion with self-attention. In *ICCV*, 2021. [2](#)
- [18] Liu He, Yijuan Lu, John Corring, Dinei Florencio, and Cha Zhang. Diffusion-based document layout generation. In *IC-DAR*, 2023. [2](#)
- [19] Minghang He, Minghui Liao, Zhibo Yang, Humen Zhong, Jun Tang, Wenqing Cheng, Cong Yao, Yongpan Wang, and Xiang Bai. Most: A multi-oriented scene text detector with localization refinement. In *CVPR*, 2021. [2](#)
- [20] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. [1](#), [2](#)
- [21] ideogram. Link: <https://ideogram.ai/>, 2023. [5](#), [10](#)
- [22] Akash Abdu Jyothi, Thibaut Durand, Jiawei He, Leonid Sigal, and Greg Mori. Layoutvae: Stochastic scene layout generation from a label set. In *ICCV*, 2019. [2](#)
- [23] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. [1](#)
- [24] Jianan Li, Jimei Yang, Aaron Hertzmann, Jianming Zhang, and Tingfa Xu. Layoutgan: Generating graphic layouts with wireframe discriminators. In *ICLR*, 2019. [2](#)
- [25] Minghao Li, Tengchao Lv, Jingye Chen, Lei Cui, Yijuan Lu, Dinei Florencio, Cha Zhang, Zhoujun Li, and Furu Wei. Trocr: Transformer-based optical character recognition with pre-trained models. In *AAAI*, 2023. [2](#)
- [26] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. In *CVPR*, 2023. [2](#)
- [27] Minghui Liao, Zhaoyi Wan, Cong Yao, Kai Chen, and Xiang Bai. Real-time scene text detection with differentiable binarization. In *AAAI*, 2020. [2](#)
- [28] Jiawei Lin, Jiaqi Guo, Shizhao Sun, Zijiang James Yang, Jian-Guang Lou, and Dongmei Zhang. Layoutprompter: Awaken the design ability of large language models. In *NeurIPS*, 2023. [3](#)
- [29] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. [7](#)
- [30] Rosanne Liu, Dan Garrette, Chitwan Saharia, William Chan, Adam Roberts, Sharan Narang, Irina Blok, RJ Mical, Mohammad Norouzi, and Noah Constant. Character-aware models improve visual text rendering. In *ACL*, 2023. [1](#), [2](#), [3](#), [5](#), [11](#)
- [31] Tengchao Lv, Yupan Huang, Jingye Chen, Lei Cui, Shuming Ma, Yaoyao Chang, Shaohan Huang, Wenhui Wang, Li Dong, Weiyao Luo, et al. Kosmos-2.5: A multimodal literate model. *arXiv preprint arXiv:2309.11419*, 2023. [2](#)
- [32] Pengyuan Lyu, Cong Yao, Wenhao Wu, Shuicheng Yan, and Xiang Bai. Multi-oriented scene text detection via corner localization and region segmentation. In *CVPR*, 2018. [2](#)
- [33] Jianqi Ma, Weiyuan Shao, Hao Ye, Li Wang, Hong Wang, Yingbin Zheng, and Xiangyang Xue. Arbitrary-oriented scene text detection via rotation proposals. *IEEE transactions on multimedia*, 2018. [2](#)

- [34] Jian Ma, Mingjun Zhao, Chen Chen, Ruichen Wang, Di Niu, Haonan Lu, and Xiaodong Lin. Glyphdraw: Learning to draw chinese characters in image synthesis models coherently. *arXiv preprint arXiv:2303.17870*, 2023. **1, 2, 5, 11**
- [35] Akshay Gadi Patil, Omri Ben-Eliezer, Or Perel, and Hadar Averbuch-Elor. Read: Recursive autoencoders for document layout generation. In *CVPRW*, 2020. **2**
- [36] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023. **2**
- [37] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. **5, 10**
- [38] Alec Radford, Luke Metz, and Soumith Chintala. Un-supervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. **1**
- [39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. **4**
- [40] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 2020. **1, 2, 11**
- [41] Jason Tyler Rolfe. Discrete variational autoencoders. *arXiv preprint arXiv:1609.02200*, 2016. **1**
- [42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. **1, 2, 3, 4, 7, 10**
- [43] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *SIGGRAPH*, 2022. **1**
- [44] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *NeurIPS*, 2022. **1, 2**
- [45] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *ACL*, 2016. **3**
- [46] 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. *IEEE transactions on pattern analysis and machine intelligence*, 2016. **2**
- [47] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021. **1**
- [48] Yipeng Sun, Zihan Ni, Chee-Kheng Chng, Yuliang Liu, Canjie Luo, Chun Chet Ng, Junyu Han, Errui Ding, Jingtuo Liu, Dimosthenis Karatzas, et al. Icdar 2019 competition on large-scale street view text with partial labeling-rrc-lsvt. In *ICDAR*, 2019. **11**
- [49] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. **3**
- [50] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. **3**
- [51] Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. Byt5: Towards a token-free future with pre-trained byte-to-byte models. *Transactions of the Association for Computational Linguistics*, 2022. **1, 2**
- [52] Yukang Yang, Dongnan Gui, Yuhui Yuan, Haisong Ding, Han Hu, and Kai Chen. Glyphcontrol: Glyph conditional control for visual text generation. In *NeurIPS*, 2023. **1, 2, 3, 5, 10, 11**
- [53] Haoxuan You, Haotian Zhang, Zhe Gan, Xianzhi Du, Bowen Zhang, Zirui Wang, Liangliang Cao, Shih-Fu Chang, and Yinfei Yang. Ferret: Refer and ground anything anywhere at any granularity. *arXiv preprint arXiv:2310.07704*, 2023. **2**
- [54] Haiyang Yu, Jingye Chen, Bin Li, Jianqi Ma, Mengnan Guan, Xixi Xu, Xiacong Wang, Shaobo Qu, and Xiangyang Xue. Benchmarking chinese text recognition: Datasets, baselines, and an empirical study. *arXiv preprint arXiv:2112.15093*, 2021. **2**
- [55] Haiyang Yu, Xiacong Wang, Bin Li, and Xiangyang Xue. Chinese text recognition with a pre-trained clip-like model through image-ids aligning. In *ICCV*, 2023. **2**
- [56] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, 2023. **1, 2, 10**
- [57] Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Kai Chen, and Ping Luo. Gpt4roi: Instruction tuning large language model on region-of-interest. *arXiv preprint arXiv:2307.03601*, 2023. **2**
- [58] Shihao Zhao, Dongdong Chen, Yen-Chun Chen, Jianmin Bao, Shaozhe Hao, Lu Yuan, and Kwan-Yee K Wong. Uni-controlnet: All-in-one control to text-to-image diffusion models. In *NeurIPS*, 2023. **1**
- [59] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. **3**
- [60] Qiang Zhou, Chaohui Yu, Shaofeng Zhang, Sitong Wu, Zhibing Wang, and Fan Wang. Regionclip: A unified multimodal pre-training framework for holistic and regional comprehension. *arXiv preprint arXiv:2308.02299*, 2023. **2**
- [61] Xinyu Zhou, Cong Yao, He Wen, Yuzhi Wang, Shuchang Zhou, Weiran He, and Jiajun Liang. East: an efficient and accurate scene text detector. In *CVPR*, 2017. **2**