

## 1 A MLMM Evaluation Details

### 2 A.1 Comparison of Judgement Templates.

3 In our study, we explored two distinct methodologies for evaluating the output of MLLMs. The first  
4 approach consolidates all generated images into a single composite image, facilitating an evaluative  
5 process from a human-centric perspective. This approach is delineated in the judgment template  
6 shown in Figure 1. Alternatively, the second method involves separately encoding each generated  
7 image, which is illustrated in the judgment template of Figure 2.

**### Instruction:**  
I'm converting text content (e.g., "A" to "Y") while aiming to maintain the original style of the typography consistently. I have generated four different typographies (images 1-4) based on an original typography (Ref). Please note that evaluations should ignore the background of each image and focus solely on the typography.

**### Task Description:**  
1. Evaluate the similarity of each given image pair: (image 1, Ref), (image 2, Ref), (image 3, Ref), and (image 4, Ref).  
2. Provide a score for each comparison based on a Likert scale from 1 to 5, where 1 denotes 'not similar at all' and 5 denotes 'very similar'.  
3. Your evaluation should focus solely on the following aspects of typography:  
- Color Consistency: Assess how closely the colors in each image's typography match those in Ref.  
- Texture Quality: Compare the surface quality and visual texture of the typography in each image with that of Ref.  
- Font Fidelity: Determine the extent to which the font style, thickness, and sharpness are preserved relative to Ref.  
4. Provide clear, specific justifications for each score, focusing on the degree of preservation or change in each specific aspect (color, design, texture, font) without implying degradation unless it directly affects the similarity score.

**### Format for Your Evaluation:**  
[Score of image 1]: [score], [Justification of your rating]  
[Score of image 2]: [score], [Justification of your rating]  
[Score of image 3]: [score], [Justification of your rating]  
[Score of image 4]: [score], [Justification of your rating]

Note: Ensure that the background of the images does not influence the scores; focus only on the typography itself.

**### Input Image:**






Ref	1	2	3	4
				

Figure 1: Judgment Template from a Human Perspective.

**### Instruction:**  
I'm converting text content (e.g., "A" to "Y") while aiming to maintain the original style of the typography consistently. I have generated four different typographies (images 1-4) from the original typography(given image 5). Please note that evaluations should ignore the background of each image and focus solely on the typography.

**### Task Description:**  
1. Evaluate the similarity of each given image pair: (image 1, image 5), (image 2, image 5), (image 3, image 5), and (image 4, image 5).  
2. Provide a score for each comparison based on a Likert scale from 1 to 5, where 1 denotes 'not similar at all' and 5 denotes 'very similar'.  
3. Your evaluation should focus solely on the following aspects of typography:  
- Color Consistency: Assess how closely the colors in each image's typography match those in image 5.  
- Texture Quality: Compare the surface quality and visual texture of the typography in each image with that of image 5.  
- Font Fidelity: Determine the extent to which the font style, thickness, and sharpness are preserved relative to image 5.  
4. Provide clear, specific justifications for each score, focusing on the degree of preservation or change in each specific aspect (color, design, texture, font) without implying degradation unless it directly affects the similarity score.

**### Format for Your Evaluation:**  
[Score about between image 5 and target image 1]: [score], [Justification of your rating]  
[Score about between image 5 and target image 2]: [score], [Justification of your rating]  
[Score about between image 5 and target image 3]: [score], [Justification of your rating]  
[Score about between image 5 and target image 4]: [score], [Justification of your rating]

Note: Ensure that the background of the images does not influence the scores; focus only on the typography itself.

**### Input Image 1** **### Input Image 2** **### Input Image 3** **### Input Image 4** **### Input Image 5**






				
---	---	---	--	---

Figure 2: Judgment template for encoding each image.

8 The results of these evaluations are presented in Figure 3 and 4, respectively. When employing  
 9 the first methodology, wherein multiple images are assembled within a single frame, the visual  
 10 comparisons and subsequent justifications of ratings demonstrate a high degree of consistency.  
 11 Conversely, the approach of encoding and evaluating each image individually reveals variations in the  
 12 visual assessments and the rationales provided for the ratings, indicating a divergence in consistency  
 compared to the composite image evaluation method.

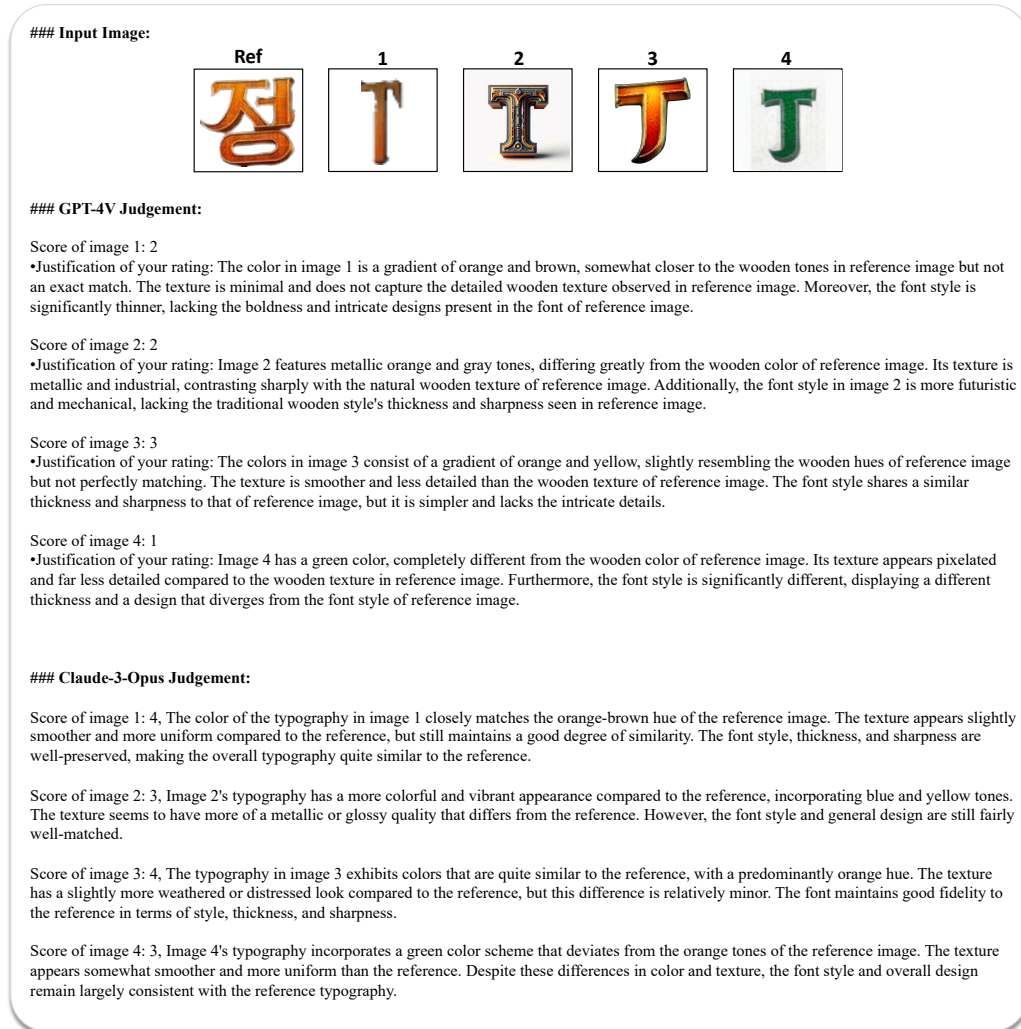


Figure 3: Judgment about Human Perspective template from GPT-4V and Claude-3-Opus.

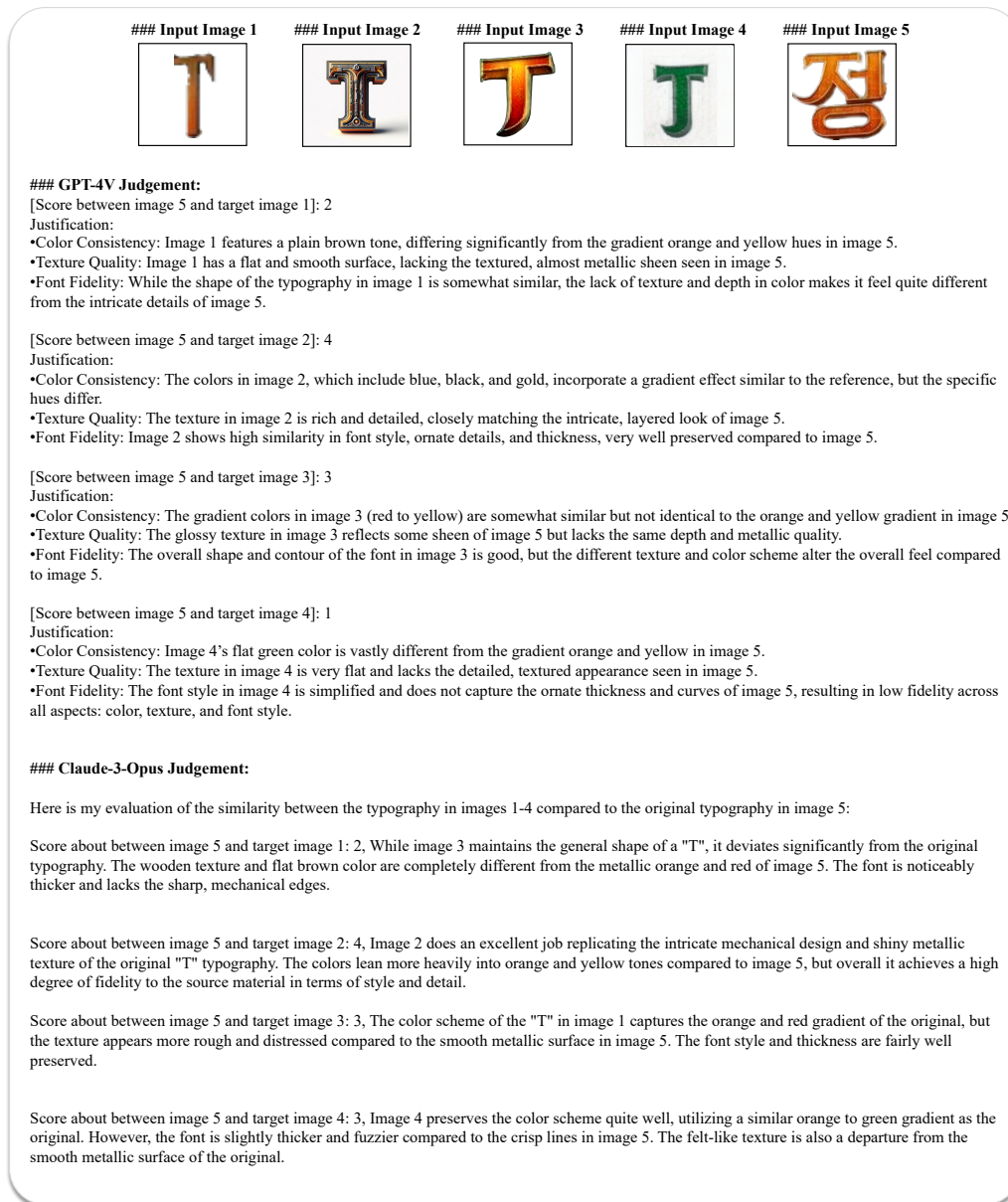


Figure 4: Judgment about encoding each image template from GPT-4V and Claude-3-Opus.

## 14 A.2 Safety Misclassification of GPT-4V

15 In Figure 5, we present examples that, despite containing content deemed safe, elicited a 'bad request' error from GPT-4V.

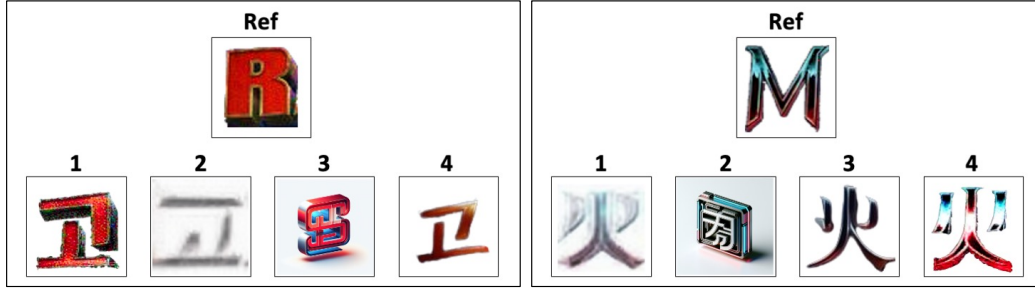


Figure 5: Samples of GPT-4V safety misclassification.

16

## 17 B More Ablation Studies

18 **Visual Text Generation Prompts.** The prompt for visual text generation can be expressed in  
 19 various ways, as shown in Figure 6. In this ablation study, the seed was fixed at 100, and the prompt  
 20 was changed for experimentation. While the initial images generated for each prompt vary greatly,  
 21 the final generated images at the end of training maintain a variety of style fidelity without losing  
 legibility.

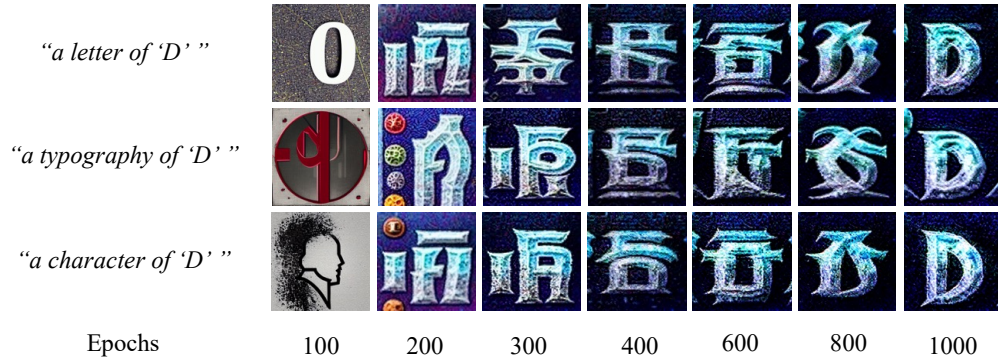


Figure 6: Various prompts for generating visual text 'D' and the corresponding generation results.

22

23 **Effectiveness of Latent Space Distances.** In the generator component of SIGIL, if glyph images  
 24 are used directly without glyph latent guidance, the generator training results are as shown in Figure  
 25 7 (a). In comparison, by using the glyph latent guidance we propose, the generation results can be  
 26 obtained as shown in Figure 7 (b).

## 27 C Dataset Curation Details

28 **Collecting Multilingual Film Poster Images.** We collected multilingual posters for movie titles  
 29 from a movie database website. Then, we kept only those titles that had pairs in different languages,  
 30 creating pairs of multilingual posters.

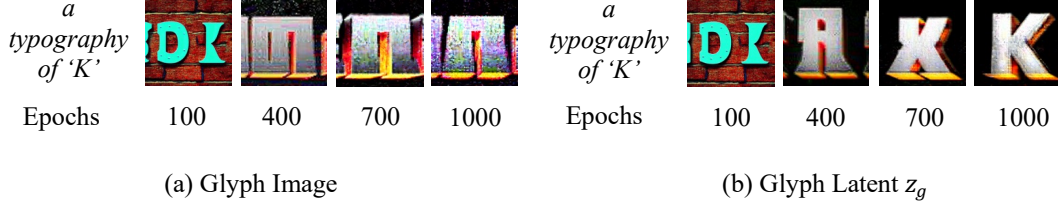


Figure 7: Generation results for the prompt "a typography of 'K'".

31 **Filtering by Style Similarity.** We hired three AI researchers to perform filtering tasks by determin-  
 32 ing the style similarity of the text in two poster images displayed on the screen. A screenshot of the  
 33 full text of instructions for this task can be found in Figure 8. Figure 8 also shows the annotation tool  
 34 we developed specifically for this task.

35 **Character-level Bounding Box Annotation.** The three human annotators who participated in the  
 36 prior filtering task volunteered for this task as well. They also carried out bounding box labeling on  
 37 the filtered images. The instructions for this task are introduced in Figure 8. Each annotator receives  
 38 a wage ranging from 12-15 dollars per hour, the total amount spent on participant compensation is  
 39 648 dollars.

#### 40 C.1 Ethical Considerations

41 In this paper, we present the MuST-Bench, which incorporates copyrighted film posters designed by  
 42 human experts. To address the issue of copyright, instead of distributing raw data, we provide down-  
 43 load links for each poster image along with bounding box annotations. Furthermore, accompanying  
 44 code is made available, enabling easy transformation of the posters into the format proposed in the  
 45 paper. This approach ensures compliance with copyright laws while maintaining the utility of the  
 46 dataset for research purposes.

## 47 D Implementation and Training Details

48 SIGIL comprises two main components, the generator and the corrector. The implementation and  
 49 training details for each are as follows:

50 **Generator.** For the pre-training configuration, we employed the runwayml/stable-diffusion-v1-5  
 51 model publicly available weights from the Huggingface Hub (<https://huggingface.co/models>). The  
 52 training dataset utilized was derived from the MuST-Bench style subject, where each style subject  
 53 allowed for fine-tuning on two glyph combinations. Training was conducted for approximately 1,000  
 54 epochs. The total epochs were adjusted based on the dataset volume; for instance, a dataset containing  
 55 15 training samples warranted an extension to 1,005 epochs to ensure thorough model training. The  
 56 learning rate was maintained at  $1e-4$  throughout the training process. Regarding the framework  
 57 and computational resources, we extended the LoRA-based implementation of Dreambooth to  
 58 process multi-subject inputs, allowing for concurrent fine-tuning on two glyph subjects. Training was  
 59 performed on a single NVIDIA A100 40GB GPU, with each session completing in about 20 minutes.

60 **Corrector.** The corrector was initiated from the last checkpoint obtained after the preliminary fine-  
 61 tuning phase (generator’s last checkpoint). The training dataset comprised images sampled during  
 62 the generator’s operation, serving as the primary data for further training. Training was conducted  
 63 with an emphasis on efficiency, incorporating an early stopping mechanism to curtail the training  
 64 as soon as the model reached a satisfactory level of performance. On average, the training duration

**Instruction:**

- **Selection Task:** Select a similar multilingual visual text image pair
- **Bounding Box Annotation Task:** Apply bounding box labeling to each individual character and the entire text for the selected similar image pair

**Process:**- **Selection Task**

1. View the two visual text images displayed on the screen, and if the style of the visual text matches, press 'K' (Keep)
2. In all other cases, press 'D' (Discard)

**Note:** Please choose based on the similarity of the visual text, not the overall similarity of the images presented.

- **Bounding Box Annotation Task**

1. Mark the area of the visual text's each individual character by placing a dot at the top left and bottom right.
2. Once the areas for each individual character are marked, finally indicate the area for the entire word.

**Example for Similar Image Pair:**

# Example 1



# Example 2

**Example for Bounding Box Annotation:**

# Top left



# Bottom right



# Annotation Result

**Individual character****Entire word**

Figure 8: A screenshot of the human annotation interface for MuST-Bench dataset curation.

was approximately 30 minutes using the same GPU as the generator. The learning rate was set to  $3e-4$ , which was determined to be optimal for achieving convergence while maintaining training stability. Additionally, the EasyOCR tool was employed to facilitate multilanguage text recognition. While utilizing this off-the-shelf OCR model, we adapted its output mechanism to provide confidence scores for characters presented in the prompt instead of predicted characters.

## E Image Input to DALL-E3

SIGIL can accept style images as input. However, DALL-E 3 only accepts text inputs. Therefore, to ensure fairness in evaluation, GPT-4V is used as a bridge to enable DALL-E 3 to "see" images. As shown in Figure 9, GPT-4V views the style image and creates a description of that style. This



74 description is then combined with the textual prompt and sent to DALL-E 3. The importance of using  
 75 GPT-4V for style description can be seen by comparing Figure 10 parts (a) and (b).

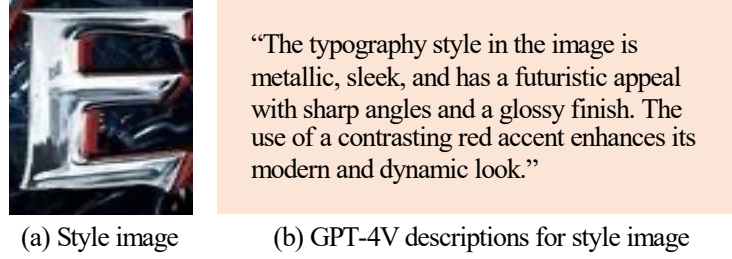


Figure 9: GPT-4V description about style image.

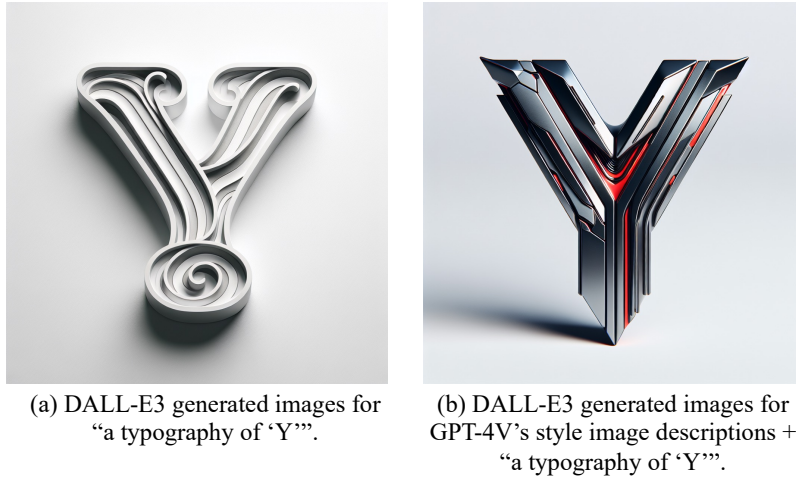


Figure 10: Comparison of DALL-E 3 generated outputs based on the presence of GPT-4V descriptions.

## 76 F User Studies

77 In the user study, assessments were conducted on two distinct parameters, style fidelity and legibility.  
 78 60 participants were provided with detailed guidelines for each instruction, as outlined in Figure 11,  
 79 before initiating the study. Figure 12 is a sample of a user study. Figure 13 exemplifies a scoring sheet  
 80 used for ranking the outcomes of comparison methods based on style fidelity, and illustrates a scoring  
 81 sheet designed for the evaluation of legibility.

## 82 G Discussion and Limitations

83 Our method can optimize generated images to match the input style image even with a small amount  
 84 of data, without relying on extensive image-caption datasets. In this study, we have successfully  
 85 combined two specific glyphs per style during training. In future work, we plan to explore methods  
 86 that enable the combination of a greater number of glyphs.

87 In the real world, there is abundant high-resolution typography data. However, for collecting multi-  
 88 language pairs with the same style, movie poster data has proven to be the most effective. When  
 89 extracting typography from movie posters, the resulting size is relatively small. If the input style

# Multi-language Visual Text Generation User Study

[English]

In this task, two types of questions will alternate.

**# Question 1** (Style Similarity): Look at the reference image and rank images A, B, C, D in order of similarity to the reference image's style.

- Note: Rank them based on style similarity, not personal preference.

- Style includes color, texture, font, etc.

**# Question 2** (Legibility of Text): Based on the given prompt, view the text in images A, B, C, D and rank them in order of accuracy and legibility.

Figure 11: User study instruction and user interface.

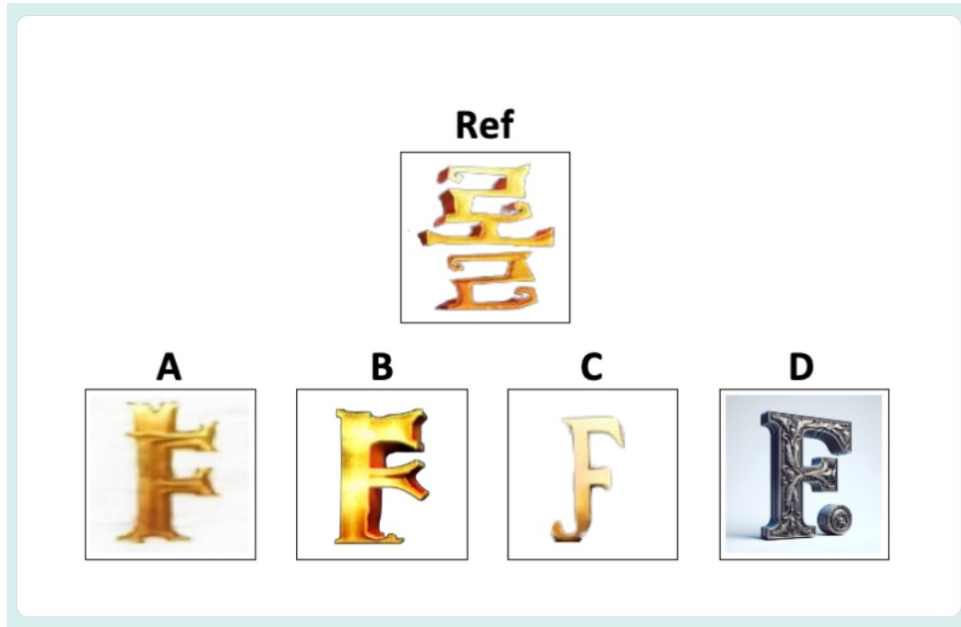


Figure 12: Evaluation sample of the user study.

90 image is low resolution, it can complicate VAE-based pixel-level encoding, potentially affecting the  
91 generation results. In future research, we aim to explore methods to generate typography in different  
92 languages from single-language style images, thereby utilizing more real-world data.

93 Unlike English, some glyphs with complex strokes in Chinese and Korean present challenges in  
94 generation. Applying fine-grained image generation methods could enable the accurate creation of all  
95 glyphs.

96 Lastly, we observed that the EasyOCR model used as a reward model in the Reinforcement Learning  
97 process sometimes exhibits False Negative issues with generated images. Additionally, if the model



Q 22-1: Style Similarity \*

	A	B	C	D
1	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
3	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
4	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

Q 22-2: Legibility (Prompt: "F") \*

	A	B	C	D
1	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
2	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
3	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
4	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

Figure 13: Questionnaire examples of style fidelity and legibility.

98 generates words not in EasyOCR's vocabulary, it can complicate the process. Future work will attempt  
 99 to use more robust OCR models with larger vocabularies and higher prediction accuracy.