

# HDLayout: Hierarchical and Directional Layout Planning for Arbitrary Shaped Visual Text Generation

Tonghui Feng<sup>1</sup> Chunsheng Yan<sup>2</sup> Qianru Wang<sup>1</sup> Jiangtao Cui<sup>1</sup> Xiaotian Qiao<sup>1, 2 \*</sup>

<sup>1</sup>School of Computer Science and Technology, Xidian University, China <sup>2</sup>Guangzhou Institute of Technology, Xidian University, China



AAAI-25 / IAAI-25 / EAAI-25  
FEBRUARY 25 – MARCH 4, 2025 | PHILADELPHIA, USA

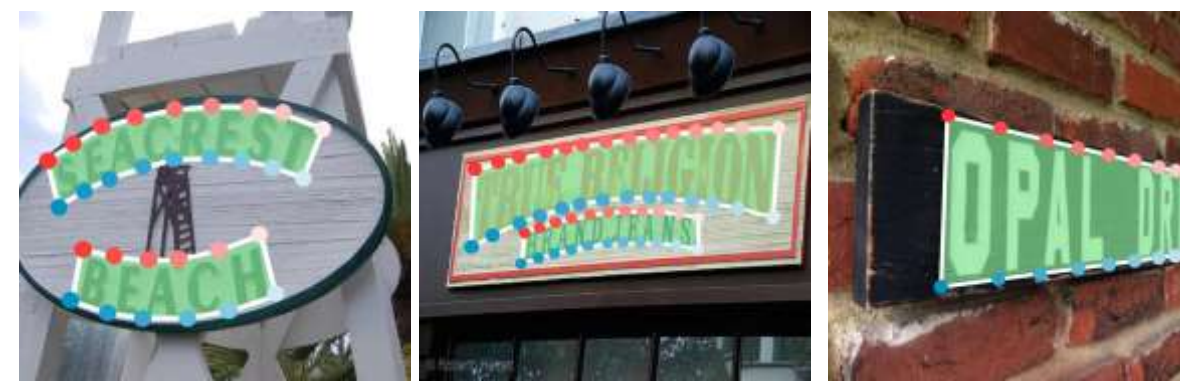
## Introduction

### Problem

- Arbitrary Shaped Visual Text Generation
  - Real-world visual text often appears in various shapes and layouts (e.g., curved, multi-oriented).



Prompt



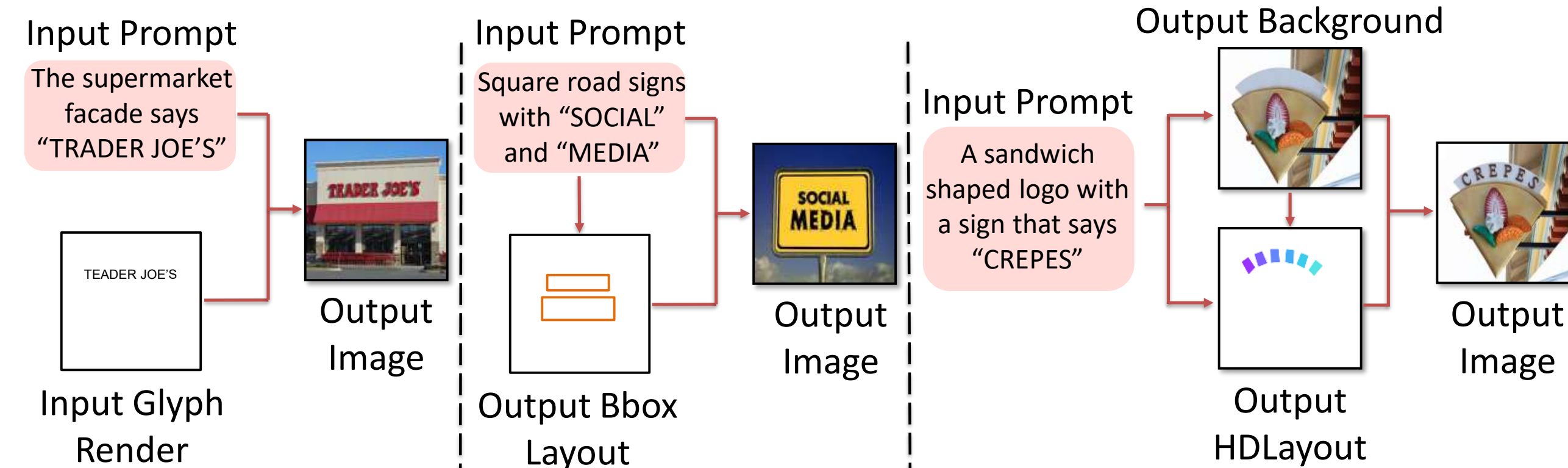
Visual Text Image

### Motivation

- Multi-Granularity Perception
  - Text in scenes inherently exhibits correlations across different granularities (character, word, line).
- Flexibility and User Accessibility
  - Design a framework that allows users to generate diverse and coherent visual text directly from text prompts.

### Related Work

- GlyphControl<sup>[3]</sup> (left), TextDiffuser<sup>[4]</sup> (middle), Ours (right).

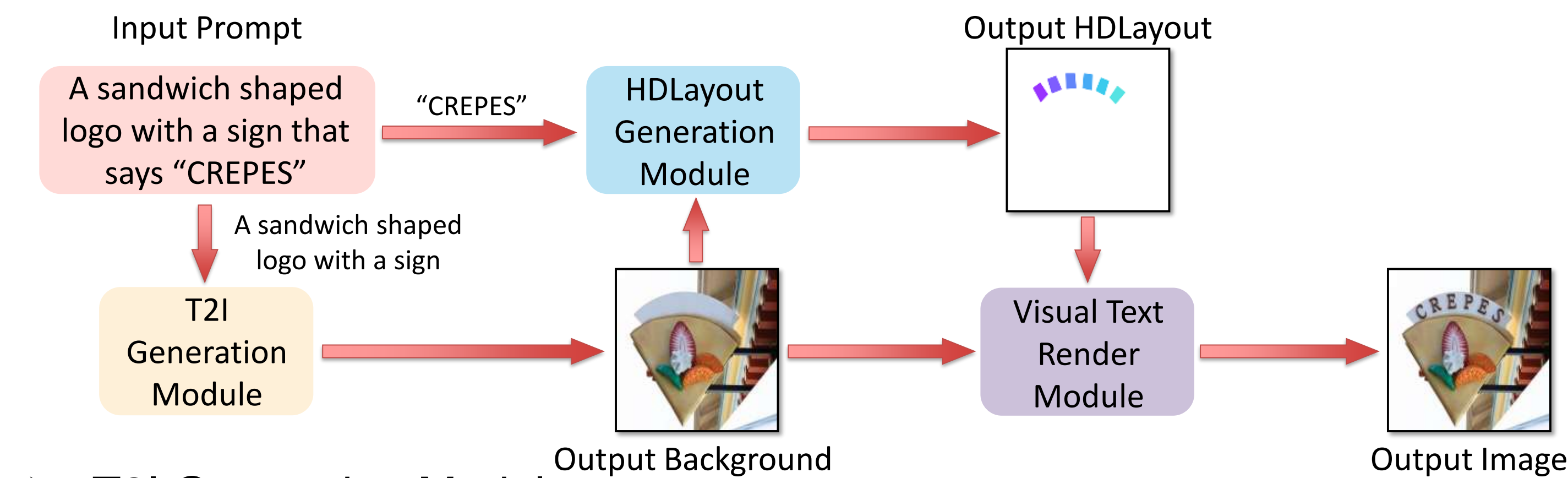


### Contribution

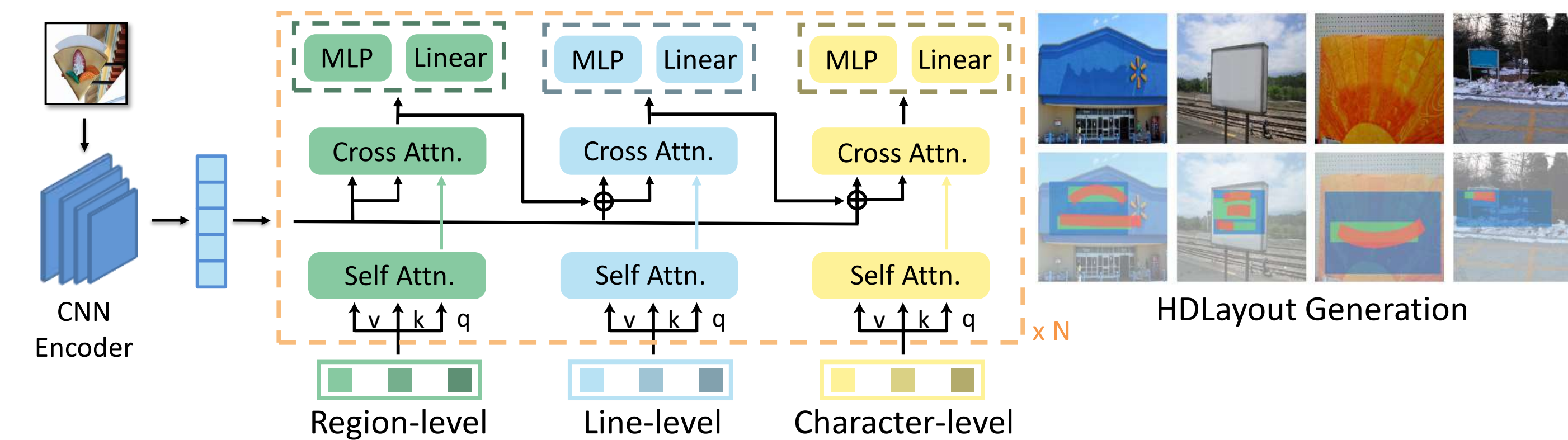
- A hierarchical and directional layout representation to model the unique characteristics of visual text.
- A new separation and composition framework to generate both textual and visual information.
- A new HDLayout3k dataset with diverse and arbitrarily shaped text layout.

## Approach

### Visual Text Generation Architecture

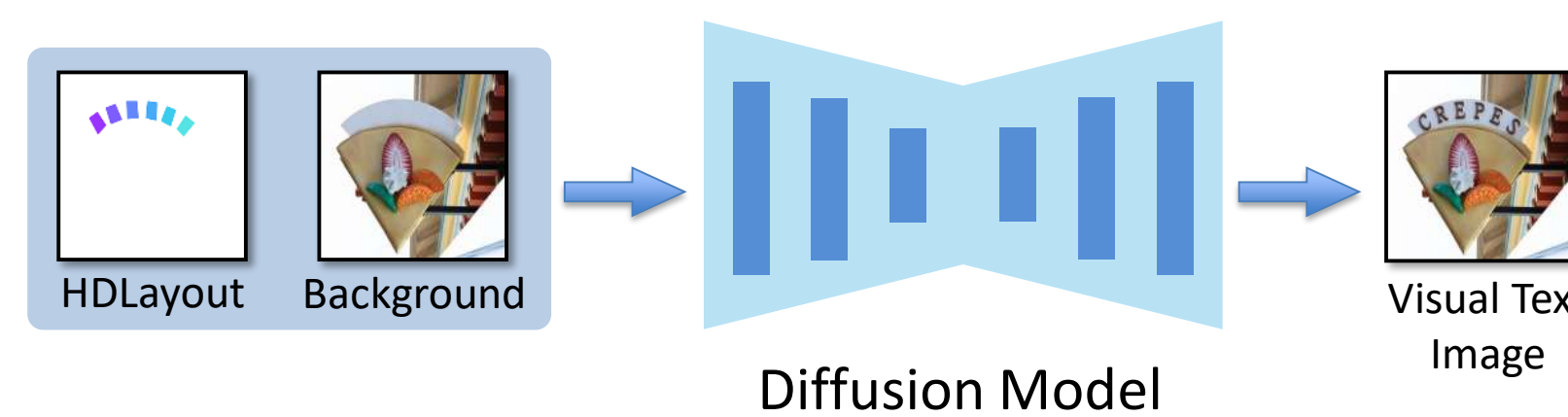


- T2I Generation Module
  - Separating description and keywords, and generating initial visual content.
- HDLayout Generation Module
  - Predicting fine-grained visual text structures.



### Text Rendering Module

- Visual text image generation.



### Training

- Bezier Point Loss  $L_{L1} = \frac{1}{N} \sum_{i=1}^N |\hat{B} - B|$
  - Bounding Box Loss  $L_{bbox} = c_1 L_{L1} + c_2 L_{GIoU} + c_3 L_{ol}$
  - Confidence Loss  $L_{conf} = -p \cdot \log(q)$
- $$L_{ol} = \frac{2}{M(M-1)} \sum_{i=1}^M \sum_{j=i+1}^M \text{IoU}(B_i, B_j) \quad L_{GIoU} = \frac{1}{N} \sum_{i=1}^N (1 - \text{GIoU}(\hat{B}, B))$$

## Dataset

### HDLayout3k

- We construct the HDLayout3k dataset, which consists of 2,749 training samples and 813 test samples.



## Evaluation

### Qualitative Results

	GlyphControl <sup>[3]</sup>		TextDiffuser <sup>[4]</sup>		Ours	
Input Prompt	Input Glyph	Output Image	Output Layout	Output Image	Output Layout	Output Image
A circle logo of a mermaid, with the words "GOOD" and "MORNING"	GOOD MORNING					
The curved arch bridge with words "Rainbow"	Rainbow					
A nice drawing of meadow, houses and sun made by a child with crayons with words "Beautiful Village"	Beautiful Village					

### Quantitative Results

Metric	SD-XL <sup>[1]</sup>	ControlNet <sup>[2]</sup>	GlyphControl <sup>[3]</sup>	TextDiffuser <sup>[4]</sup>	Ours
Image FID ↓	102.128	89.802	83.402	82.068	<b>78.027</b>
Layout FID ↓	-	-	-	27.135	<b>12.343</b>

### References

- [1] Podell, D , et al. Sdxl: Improving latent diffusion models for high-resolution image synthesis. In ICLR 2024.  
 [2] Zhang, L , et al. Adding conditional control to text-to-image diffusion models. In ICCV 2023.  
 [3] Yang, Y , et al. GlyphControl: Glyph Conditional Control for Visual Text Generation. In NeurIPS 2024.  
 [4] Chen, J , et al. Textdiffuser: Diffusion models as text painters. In NeurIPS 2024.