# DuetSVG: Unified Multimodal SVG Generation with Internal Visual Guidance

Peiying Zhang[1*]   Nanxuan Zhao[2]   Matthew Fisher[2]
Yiran Xu[2]   Jing Liao[1]   Difan Liu[2]

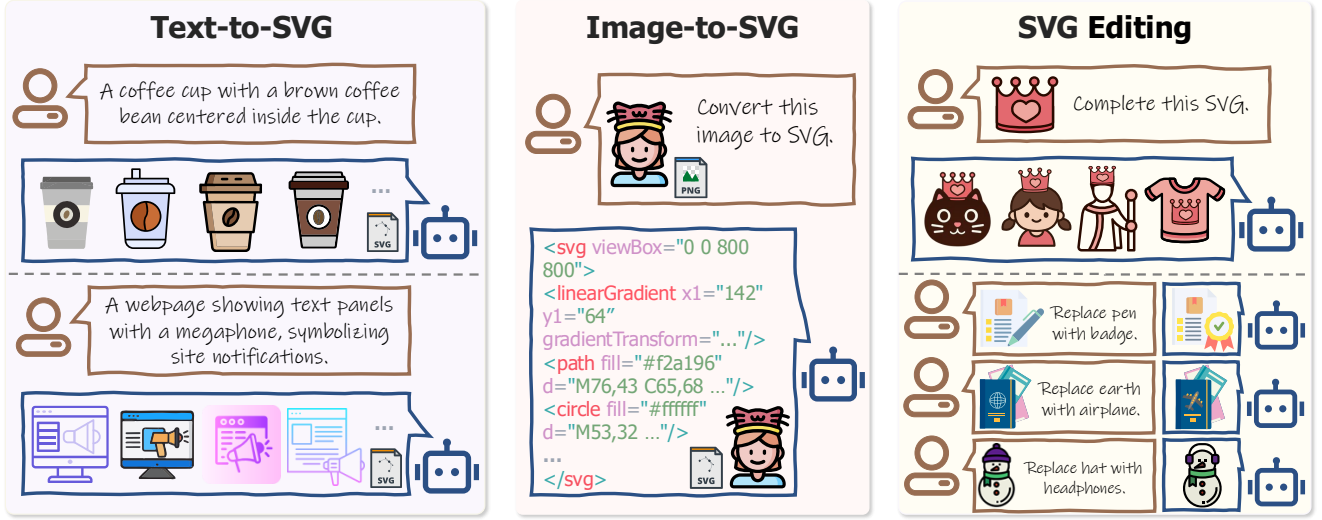[1]City University of Hong Kong   [2]Adobe Research

Figure 1. **We propose a unified multimodal model, DuetSVG, for SVG generation.** DuetSVG acts as a versatile framework across text-to-SVG, image-to-SVG, and SVG editing tasks, demonstrating strong semantic alignment and high-quality SVG generation.

## Abstract

*Recent vision-language model (VLM)-based approaches have achieved impressive results on SVG generation. However, because they generate only text and lack visual signals during decoding, they often struggle with complex semantics and fail to produce visually appealing or geometrically coherent SVGs. We introduce DuetSVG, a unified multimodal model that jointly generates image tokens and corresponding SVG tokens in an end-to-end manner. DuetSVG is trained on both image and SVG datasets. At inference, we apply a novel test-time scaling strategy that leverages the model's native visual predictions as guidance to improve SVG decoding quality. Extensive experiments show that our method outperforms existing methods, producing visually faithful, semantically aligned, and syntactically clean SVGs across a wide range of applications. The project page is* `https://intchous.github.io/DuetSVG-site`.

## 1. Introduction

Scalable Vector Graphics (SVG) are widely used in graphic design, digital art, publishing, and motion graphics. Compared to raster images, SVGs offer resolution-independent rendering, efficient storage, and intuitive editability via control point manipulation. However, creating high-quality vector graphics remains a complex and time-consuming process, even for experienced designers.

With the increasing prominence of large language models (LLMs) and vision-language models (VLMs), recent approaches [31, 44, 57, 58] have leveraged the textual nature of SVGs by formulating SVG generation as a text generation problem and finetuning large text generation models [2, 3, 64]. These methods produce impressive results on SVG generation tasks such as text-to-SVG. However, SVGs differ fundamentally from plain text, as they contain an additional dimension—the visual aspect. Formulating SVG generation purely as a text generation task introduces inherent limitations and hampers the performance of existing LLM-based methods. First, existing approaches

*Work done during internship at Adobe Research.

1

are unimodal generative models that produce only text to-kens. Minor errors, such as incorrect predictions of path coordinates, may appear negligible in text space but can lead to catastrophic failures in the rendered SVG. The absence of visual guidance during text generation represents a main limitation for these models on tasks such as text-to-SVG and SVG completion. Second, LLM-based SVG generation models exhibit poor generalization beyond their training distribution. As unimodal generative models, they are restricted to training on the relatively small amount of available SVG data and cannot leverage the abundance of high-quality raster image datasets (such as text-image pairs and image editing data), which significantly constrains their generalization capability.

To address the above-mentioned challenges, we propose a novel multimodal generative model that produces native image outputs for SVG generation and editing. Our model generates a multimodal sequence consisting of image tokens and SVG text tokens. The generated image tokens serve as internal visual guidance during SVG token generation, enabling more coherent and visually grounded SVG results. Moreover, this multimodal generation framework unlocks new capabilities beyond those of purely text-based vector models. For example, by jointly training on tasks such as text-to-image and text-to-SVG, the model can leverage large-scale text-image datasets for pretraining, which greatly improves generalization and text-SVG alignment. The multimodal generative nature of our method also simplifies verifier design for test-time scaling. We further introduce a novel and efficient scaling strategy that enhances the model's reliability and robustness during inference. On the input side, our model accepts multimodal conditions, including raster images, SVG code, and text prompts, and supports a wide range of tasks such as text-to-SVG, image-to-SVG, SVG completion, and SVG editing.

**Contributions.** We present the first unified multimodal generative model for SVG generation. Our native visual guidance enables visually grounded SVG generation and allows SVG tasks to be trained on image datasets. We further introduce a novel test-time scaling strategy that efficiently improves model reliability. We demonstrate that DuetSVG achieves much higher quality than state-of-the-art (SoTA) methods on multiple benchmarks.

## 2. Related Work

### 2.1. Optimization-based SVG Generation

Classic image vectorization methods [4, 6, 12, 13, 19, 22, 26, 37, 42] rely on fitting algorithms to reconstruct vector graphics from raster images. Although these methods can accurately reproduce the overall appearance of an image, they often generate redundant paths, imprecise control points, and struggle to handle path occlusions.

Recent approaches utilize pre-trained vision-language models (VLMs), such as CLIP [30] and diffusion models [34], to directly optimize SVG paths through differentiable rendering [24]. CLIP-based methods [14, 35, 38, 43] optimize SVG representations by maximizing image-text alignment within CLIP's latent space. To leverage the strong visual and semantic priors of text-to-image diffusion models, several methods employ score distillation sampling [28, 49] to optimize static [20, 21, 54, 55, 60] or animated [15, 51] SVGs to align with textual descriptions. However, these methods often require tens of minutes to optimize a single SVG, making them impractical for real-world applications. More importantly, as they are not trained on vector graphics data, they often produce fragmented paths, inconsistent topology, and redundant control points, which complicate subsequent editing and manipulation.

### 2.2. Learning-based SVG Generation

Early approaches to SVG generative modeling formulated SVGs as sequences of geometric primitives, using VAEs [5, 46, 47], or diffusion models [10, 41] as the foundational generative models. VecFusion [41] first employs a raster diffusion model to generate an image, followed by a vector diffusion model conditioned on the raster output to produce vector graphics. However, the lack of end-to-end training limits generalization between the raster and vector models, often leading to inaccurate control points and suboptimal geometry.

Recent advances in large language models have inspired SVG generation approaches [31, 33, 40, 44, 50, 52, 56–58] to represent SVG scripts as discrete text tokens through specialized tokenization schemes, enabling the autoregressive generation of SVG command sequences. These methods finetune LLMs or VLMs on SVG datasets for tasks such as Text-to-SVG, Image-to-SVG and SVG editing, achieving impressive results. However, these models exhibit limited generalization because they are trained on relatively small SVG datasets. Since SVG generation is formulated as a text generation task, they also lack visual guidance during inference, which further constrains their output quality. Concurrent work RoboSVG [45] relies on external VLMs to generate additional multimodal conditions as input, which may introduce inconsistencies between models. In contrast, our unified multimodal generative model co-generates image and SVG tokens within an end-to-end architecture, enabling the use of large-scale text-image data and improving visual grounding during SVG decoding.

### 2.3. Unified Multimodal Generation

Recent unified multimodal models have made substantial progress, showing that a single architecture can both understand and generate multiple modalities, including

fully autoregressive [7, 9] and AR-diffusion fused [11, 53] paradigms. We refer readers to [62] for a more comprehensive review. However, as SVG scripts differ significantly from natural language in structure, SoTA VLMs still fail to produce high-quality vector graphics without SVG-specialized training.

## 3. Method

### 3.1. Task Definition

An SVG file contains a sequence of paths with drawing commands (e.g., `M`, `C`, `Q`, `Rect`), numeric coordinates, and style attributes (e.g., `fill`, `stroke`). Previous work [31, 58] formulate SVG generation as a code generation task and fine-tunes language models to output SVG as text tokens. However, these methods are limited to SVG-based training data and generalize poorly to complex or out-of-distribution inputs. The absence of visual guidance during SVG generation further constrains their capabilities. In contrast, we train a unified multimodal generative model that produces both image and SVG tokens. The image modality captures appearance, and the SVG modality learns shape geometry and layer structure. Formally, given a multimodal conditioning input $\mathbf{x}$, which may include text prompts, images, and SVG code, our model generates a mixed-modality target sequence $\mathbf{z}$ consisting of image tokens $z^{\text{img}}$ and SVG tokens $z^{\text{svg}}$, where $\mathbf{z} = [\langle\text{IMG}\rangle, z^{\text{img}}_{1:I}, \langle/\text{IMG}\rangle, \langle\text{SVG}\rangle, z^{\text{svg}}_{1:S}, \langle/\text{SVG}\rangle]$. Training maximizes the unified next-token objective:

$$P_\theta(\mathbf{z} \mid \mathbf{x}) = \prod_{t=1}^{T} p_\theta\left(z_t \mid z_{<t}, \mathbf{x}\right) \quad (1)$$

where $\theta$ denotes the model parameters.

### 3.2. Unified Multimodal SVG Generation Model

**Model Architecture.** Building on recent unified autoregressive models, our architecture follows Janus-Pro [7], which supports multimodal generation of both image and text tokens, as illustrated in Figure 2. Our model accepts multimodal inputs, including text prompts, SVG code, and images. We allow different input configurations for different tasks, such as text prompts for text-to-SVG generation and all three modalities for SVG completion. Text and SVG inputs are embedded using Janus-Pro's text tokenizer. For image inputs, we employ SigLIP [59] as the understanding encoder to extract semantic features, and a VQ tokenizer [39] as the generation encoder to convert images into compact discrete embeddings. Two separate MLP-based aligners, one for understanding and one for generation, map image embeddings into the LLM feature space. The concatenated multimodal sequence is input to a unified autoregressive transformer with causal attention, enabling the model
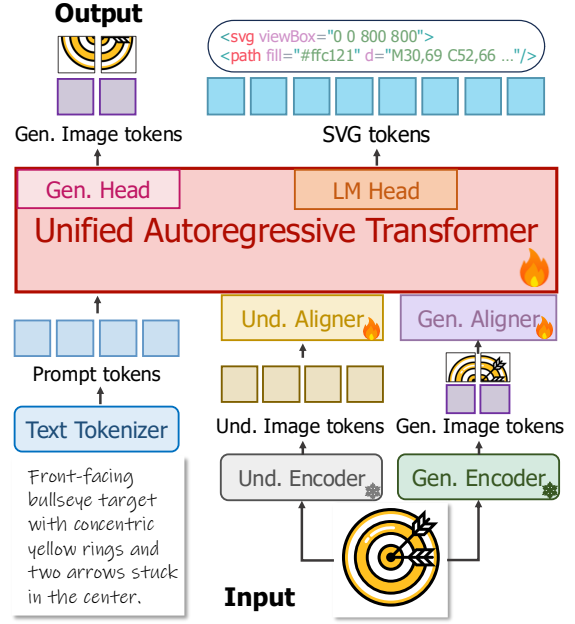


Figure 2. **Model architecture of DuetSVG.** As a unified model, DuetSVG accepts multimodal inputs, including text prompts, SVG code and raster images. We use Janus-Pro text tokenizer [7] for text prompts. For images, an Understanding (Und.) Encoder extracts semantic features, while a Generation (Gen.) Encoder converts images into discrete visual tokens. Two MLP aligners project the encoder outputs into the same feature space as the text embeddings. A Generation (Gen.) head predicts image tokens, and a language modeling (LM) head predicts SVG tokens.

to learn cross-modal alignment and next-token prediction over visual and SVG modalities. We use a generation head to predict image tokens from a visual codebook and an LM head to predict SVG tokens from a text vocabulary. During training, the parameters of the understanding and generation image encoders are frozen, while all other parameters are trainable.

**Training Stages.** Leveraging its multimodal generative design, DuetSVG can be trained on both image and SVG datasets. Training proceeds in two stages: text-to-image pretraining and multi-task supervised fine-tuning (SFT).

Because the Janus-Pro base model has limited ability to generate SVG-style images, we begin with large-scale text-to-image (T2I) pretraining in **Stage 1**. The objective of this stage is to strengthen the model's capacity to produce visually appealing and clean images characterized by clear geometric primitives and flat colors. We train on a hybrid corpus of real and synthetic T2I data, including: (i) rendered SVG images paired with captions from curated SVG datasets, and (ii) synthetic images generated by FLUX.1 [23], which takes a text prompt and an SVG reference to produce images that match the reference's style. This pre-

"A lightbulb contains a green plant, symbolizing eco-friendly or sustainable energy ideas"

Input text

CFG=1  CFG=1.5  CFG=2  CFG=2.5

(a) Visual candidate selection

Guidance

Iterative SVG token sampling
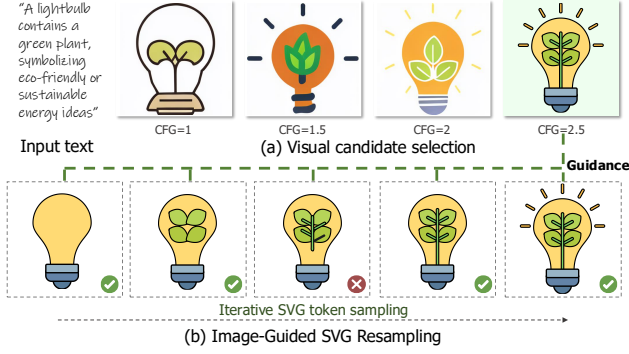
(b) Image-Guided SVG Resampling

Figure 3. **Test-time scaling with image-guided SVG resampling.** (a) We first generate $N$ raster image candidates with CFG. Since image-token sequences are much shorter than SVG-token sequences, this step is efficient. A CLIP-based verifier [30] selects the best candidate $I^*$. (b) Using the selected image tokens as internal guidance, we iteratively generate SVG tokens. At each iteration, we render a raster image $R_t$ from the current SVG codes and accept the update only if the LPIPS distance $d(R_t, I^*)$ does not increase; otherwise, we reject and resample.

training stage enables the model to learn strong semantic and visual priors from large-scale T2I data, providing a robust initialization for subsequent SVG generation tasks. As shown in our experiments, T2I pretraining helps the model generalize better to complex text prompts and out-of-distribution inputs when generating SVGs.

In **Stage 2**, we perform SFT across multiple tasks including T2I, T2SVG, and I2SVG under a unified next-token prediction objective with cross-entropy loss over interleaved multimodal outputs. For SVG-generation tasks, we arrange the target sequence as image tokens followed by SVG tokens so that, during autoregressive decoding, the image tokens can guide the SVG tokens via causal attention. To strengthen robustness and improve the model's understanding of structural relationships in the I2SVG task, we apply SVG-specific data augmentations. In particular, we randomly modify the rotation, translation, scaling, and color attributes of the target SVG, optionally remove a subset of its paths, and render the modified SVG into an image. We further apply random dropout to the text and image inputs with probability $10\%$, enabling classifier-free guidance [18] during inference. Multi-task SFT allows the model to share knowledge across modalities and tasks. For example, T2I and I2SVG can enhance T2SVG in different ways, leading to better model quality and stronger generalization. In an **optional Stage 3**, DuetSVG can be further finetuned for downstream applications such as SVG completion, as detailed in Section 5.

## 3.3. Test-Time Scaling with SVG Resampling

For complex SVGs containing thousands of tokens, autoregressive decoding can accumulate sampling errors—such as spurious loops or weakened grounding—often producing suboptimal geometry or even invalid SVGs. In pure text-based generation model, a common test-time scaling approach is best-of-$N$ sampling: the model produces $N$ complete SVG rollouts, renders each, and a verifier (*e.g.*, CLIP) selects the best result. This strategy is computationally expensive and only reranks after full outputs are generated, providing no visual guidance during decoding. In contrast, our multimodal model jointly generates image and SVG tokens, enabling a more efficient test-time scaling method with image-guided resampling during inference (see Figure 3). Our procedure has two stages: (1) selecting the best visual candidate at the image level, and (2) performing image-guided resampling as the SVG generation continues from the stage-1 visual.

**Visual Candidate Selection.** We first generate $N$ visual candidates using classifier-free guidance (CFG):

$$z_t^{\text{CFG}} = z_t^{\text{uncond}} + \gamma\big(z_t^{\text{cond}} - z_t^{\text{uncond}}\big) \tag{2}$$

where $z_t^{\text{cond}}$ and $z_t^{\text{uncond}}$ are the model predictions for image tokens with and without conditioning, respectively, and $\gamma$ is the guidance scale. Because image-token sequences are much shorter than SVG-token sequences, sampling $N$ candidate images is relatively efficient. We then score each candidate image using CLIP as the verifier and keep the best one, denoted by $I^*$ with corresponding image tokens $z_{\text{img}}^*$, which can serve as visual guidance during SVG decoding.

**Image-Guided SVG Resampling.** We continue the SVG token generation from the best image tokens $z_{\text{img}}^*$. We generate SVG tokens in small chunks with image-guided resampling. More specifically, at each iteration, we generate $K$ SVG tokens, append them to the current SVG script, and render a provisional raster $R_t$. We then compute its perceptual distance to the best visual candidate $d(R_t, I^*)$ using LPIPS [61]. If $d(R_t, I^*)$ is less than or equal to $d(R_{t-1}, I^*)$, we accept the newly generated tokens; otherwise, we reject them and resample, allowing up to $M$ rejections per SVG. This image-guided resampling encourages the decoded SVG to remain consistent with the selected visual candidate while avoiding the high cost of best-of-$N$ sampling over long SVG-token sequences.

By coupling image-level search with chunked, image-guided SVG resampling, our test-time scaling strategy improves semantic alignment and SVG validity at much lower computational cost than naive best-of-$N$ sampling.

## 3.4. SVG-Hub Dataset

Existing SVG datasets and benchmarks suffer from limited quality and diversity. On one hand, many are constructed by

vectorizing raster images (e.g., MMSVG [58], InternSVG [44]), which introduces irregular paths and visual artifacts that compromise SVG structure and regularity. On the other hand, the accompanying text descriptions are often short and generic, lacking the fine-grained semantics required for high-quality, complex T2SVG generation.

To address these issues, we introduce SVG-Hub-1M, curated from diverse public SVG sources (MMSVG[58], SVGX [57], and Iconfont[1]) with data cleaning and standardization. We remove duplicates, auto-vectorized and blank-rendering SVGs. The SVG-Hub-1M dataset will be released to support future research. We also conducted experiments on an internal large-scale dataset SVG-Hub-5M. Both datasets consist of high-quality SVGs that are not vectorized results of raster images.

**SVG Captioning.** Previous SVG datasets typically provide only simple text descriptions, which are insufficient for training models to understand complex semantics and generate semantically aligned SVGs. To support T2SVG from semantically rich prompts, we rasterize each SVG and use open-source VLMs (InternVL3 [64] and Qwen2.5-VL [3]) to produce captions at three levels of detail: (1) short prompts that capture core semantics, (2) medium descriptions that enumerate semantic elements along with their layout and style, and (3) detailed annotations covering fine-grained shapes, strokes, and colors. Our dataset pairs each SVG with a comprehensive set of captions. Additional details of the captioning pipeline are provided in Section D.

**SVG Tokenization.** An SVG file contains geometric primitives and paths with drawing commands and attributes. We construct a compact and regular representation by (i) removing redundant or invisible elements, (ii) normalizing the canvas to an $800\times800$ `viewBox`, and (iii) restricting the command vocabulary to $\{$`M, L, C, Q, A, Z, Ellipse, Circle, Polygon, Rect`$\}$. We then quantize all coordinates $(x, y)$ and serialize the normalized SVG into a sequence of discrete tokens that includes command, attribute, and quantized coordinate tokens. Gradient definitions (`<defs>`) and group-level transformations (`<g>`) are retained to preserve expressiveness. These steps standardize the SVG script structure and reduce file size while remaining lossless with respect to rendering.

## 4. Experiments

### 4.1. Implementation Details

We initialize DuetSVG from Janus-Pro-7B [7]. We resize each image to $384 \times 384$, and then use the generation encoder to encode it into visual tokens with a sequence length of 576, with a codebook of size 16,384. Each SVG is tokenized and truncated to a maximum of 12,000 text tokens. During the T2I pre-training stage, we use a mixture of real

and synthetic T2I data, and train for 80K steps with a batch size of 512. In the multi-task SFT stage, we jointly train on T2I, T2SVG, and I2SVG data, sampling them with a ratio of 1:5:4, respectively. We train this stage for 300K steps with a batch size of 128. The full training process takes about 14 days on 64 NVIDIA-A100 GPUs. We use AdamW [25] optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.95$ with a learning rate of $1 \times 10^{-5}$ in all stages. For the test-time scaling strategy, we set $N = 3$ and $M = 3$ in our experiments.

### 4.2. Experiment Setup

**Dataset and Benchmark.** We evaluate our method on two benchmarks. One is the test split from SVG-Hub-5M which has 9,000 samples. The other one is the SArena-Icon Benchmark [44] which contains 6,000 samples.

**Evaluation Metrics.** We evaluate the quality of our results from both vector-level and image-level perspectives. For **vector-level** evaluation, we measure *Path Semantics* [60] by randomly removing 30% of the SVG paths and computing the drop in CLIP score [30] between the original and modified renderings. A smaller drop indicates that the generated paths are redundant or carry limited semantic meaning. In I2SVG, we measure *SVG code similarity* by encoding the generated and ground-truth SVG code with Qwen3-Embedding-8B [63] and computing the cosine similarity between their embeddings, which reflects the syntactic quality of the generated SVGs. For **image-level** evaluation, we assess the visual fidelity and quality of T2SVG using *FID* [17], *FID-CLIP* [50], *CLIP* score [30] and the *Aesthetic* score [36]. For I2SVG, we measure the visual similarity between the rendered SVGs and the input images using *DINO* [27], *SSIM* [48], *PSNR*, and *LPIPS* [61].

**Baselines.** We compare our DuetSVG against a diverse set of baselines, including both optimization-based and learning-based SVG generation methods. For optimization-based methods, in the T2SVG task, one baseline uses a raster-then-vectorize pipeline: images are generated by FLUX.1-dev [23] and then vectorized by VTracer [29, 37]. We also evaluate three text-guided SVG optimization methods, VectorFusion [21], SVGDreamer [55], and T2I-NPR [60], each optimized with 64 paths. For the I2SVG task, we use VTracer as a baseline. Learning-based methods include SoTA proprietary and open-source VLMs. For proprietary VLMs, we use GPT-5-Thinking [1], Gemini-3-Pro [16] and Gemini-2.5-Pro [8]. For open-source SVG-specific VLMs, we compare with StarVector [32], LLM4SVG [57], and OmniSVG [58], and we also include the recently released VLM Qwen3-VL-8B [2]. Since these models adopt different backbones and training data, we further **fine-tune** all open-source VLM baselines on our SVG-Hub-5M dataset to ensure a fair comparison.

Table 1. **Quantitative comparison with existing methods on the SVG-Hub-5M test set. Bold** scores indicate the best results among all methods, and underlined scores indicate the best results among VLM-based methods without test-time scaling (w/o TTS).

| | Method | Text-to-SVG | | | | Image-to-SVG | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | FID ↓ | CLIP ↑ | Aesthetic ↑ | Path Semantics ↑ | DINO ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SVG Code Similarity ↑ | Path Semantics ↑ |
| Optimization | VTracer | - | - | - | - | **0.968** | **0.936** | **0.071** | **23.623** | 0.788 | 0.982 |
| | FLUX.1-dev + VTracer | 46.990 | 25.326 | 5.52 | 1.216 | - | - | - | - | - | - |
| | Vectorfusion | 59.354 | 24.519 | 5.05 | 1.405 | - | - | - | - | - | - |
| | SVGDreamer | 56.743 | 24.871 | 5.34 | 1.468 | - | - | - | - | - | - |
| | T2V-NPR | 54.420 | 25.024 | 5.38 | 1.924 | - | - | - | - | - | - |
| VLM | GPT-5-Thinking | 50.122 | 24.950 | 5.39 | 2.295 | 0.904 | 0.806 | 0.212 | 11.485 | 0.916 | 2.496 |
| | Gemini-3-Pro | 48.765 | 25.146 | 5.46 | 2.412 | 0.921 | 0.880 | 0.116 | 13.858 | 0.908 | 2.516 |
| | Gemini-2.5-Pro | 57.597 | 24.572 | 5.20 | 2.154 | 0.885 | 0.697 | 0.275 | 10.237 | 0.887 | 2.028 |
| | StarVector-8B (w/o FT) | - | - | - | - | 0.691 | 0.489 | 0.388 | 8.920 | 0.862 | 1.246 |
| | StarVector-8B (FT) | - | - | - | - | 0.896 | 0.813 | 0.186 | 18.646 | 0.932 | 1.830 |
| | LLM4SVG-7B (FT) | 49.318 | 23.304 | 5.24 | 2.315 | 0.938 | 0.882 | 0.099 | 19.843 | 0.944 | 2.245 |
| | OmniSVG-3B (w/o FT) | 93.172 | 21.225 | 4.32 | 1.780 | 0.812 | 0.701 | 0.294 | 11.577 | 0.893 | 1.505 |
| | OmniSVG-3B (FT) | 58.237 | 22.726 | 5.17 | 2.286 | 0.933 | 0.854 | 0.105 | 19.256 | 0.932 | 2.164 |
| | Qwen3-VL-8B (FT) | 43.720 | 23.935 | 5.35 | 2.525 | 0.947 | 0.910 | 0.090 | 20.915 | 0.950 | 2.492 |
| | Ours-7B (w/o TTS) | 35.066 | 25.584 | 5.52 | 2.772 | 0.955 | 0.920 | 0.082 | 22.024 | 0.959 | 2.558 |
| | Ours-7B (TTS) | **33.574** | **26.106** | **5.56** | **2.910** | 0.962 | 0.928 | 0.075 | 23.590 | **0.964** | **2.724** |

Table 2. **Quantitative comparison with VLM-based baselines on the SArena-Icon benchmark [44].**

| Method | Text-to-SVG | | | Image-to-SVG | | | |
|---|---|---|---|---|---|---|---|
| | FID ↓ | FID-C ↓ | CLIP ↑ | DINO ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ |
| GPT-5-Thinking | 14.892 | 4.993 | 25.125 | 0.916 | 0.815 | 0.134 | 11.512 |
| Gemini-3-Pro | 14.324 | 4.847 | 25.382 | 0.938 | 0.874 | 0.125 | 14.366 |
| Gemini-2.5-Pro | 15.484 | 5.016 | 24.910 | 0.895 | 0.690 | 0.261 | 10.878 |
| StarVector-8B (w/o FT) | - | - | - | 0.871 | 0.623 | 0.206 | 13.595 |
| StarVector-8B (FT) | - | - | - | 0.920 | 0.844 | 0.118 | 19.076 |
| LLM4SVG-7B (FT) | 18.815 | 6.784 | 23.138 | 0.946 | 0.910 | 0.092 | 22.160 |
| OmniSVG-3B (w/o FT) | 28.292 | 11.318 | 21.679 | 0.894 | 0.756 | 0.186 | 12.669 |
| OmniSVG-3B (FT) | 19.941 | 7.477 | 22.450 | 0.937 | 0.908 | 0.099 | 21.984 |
| Qwen3-VL-8B (FT) | 16.780 | 5.192 | 23.704 | 0.955 | 0.921 | 0.081 | 22.750 |
| Ours-7B (w/o TTS) | 12.105 | 4.205 | 25.416 | 0.969 | 0.929 | 0.069 | 23.482 |
| Ours-7B (TTS) | **11.712** | **3.928** | **25.734** | **0.972** | **0.938** | **0.060** | **24.016** |

## 4.3. Comparisons

We evaluate the performance of DuetSVG against all baselines qualitatively and quantitatively on both T2SVG and I2SVG tasks. We report quantitative results on the SVG-Hub-5M test set in Table 1, and on the SArena-Icon benchmark in Table 2. We also show qualitative comparisons in Figure 4 for T2SVG and Figure 5 for I2SVG. Overall, DuetSVG achieves consistently stronger performance than existing methods. We refer readers to Section B.2 for additional comparisons and results.

**Text-to-SVG Task.** As shown in Table 1 and Table 2, DuetSVG consistently outperforms all T2SVG baselines across all metrics.

Optimization-based methods are time-consuming and often yield SVGs with redundant paths, artifacts, and disorganized layers. Regarding learning-based methods, SoTA proprietary VLMs exhibit strong semantic understanding, but they mainly produce combinations of oversimplified primitives (*e.g.*, circles and rectangles) and struggle with spatial layout and fine-grained geometric details (see Figure 4). It indicates that precise SVG generation remains challenging for general VLMs. For open-source VLMs,

evaluating their public checkpoints reveals severe overfitting: they often fail to produce valid SVGs for inputs outside their training distribution. To enable a fair comparison, we fine-tune all VLMs on our SVG-Hub-5M training set. As shown by the "w/o FT" and "FT" in Table 1 and Table 2. However, these fine-tuned VLMs still struggle to produce semantically accurate SVGs with the complex geometric detail, as shown in Figure 4. A major reason is their text-centric design: SVGs are treated purely as text sequences, so training mainly enforces syntactic correctness of the SVG code while providing no supervision on the visual appearance.

In contrast, DuetSVG produces diverse, visually appealing SVGs with clear structure and semantically rich content. We further analyze the model's generalization in Section B.1.

**Image-to-SVG Task.** Classical image vectorization techniques achieve strong reconstruction scores by densely fitting paths to the image, but they generate verbose and inefficient code (reflected by low *SVG Code Similarity* in Table 1), lack layer organization (low *Path Semantics* in Table 1), and struggle with overly complex vector elements (shown in the green boxes in Figure 5). VLM-based baselines struggle to accurately reconstruct complex shapes and fine details, often resulting in inaccurate spatial layouts and simplified fine-grained geometry. In contrast, DuetSVG produces visually faithful and syntactically clean SVGs, as shown in Figure 5.

## 4.4. Ablation Study

**Benefits of Internal Visual Guidance.** To assess the effect of multimodal SVG generation, we train a variant that shares the same training configuration as DuetSVG but only decodes SVG tokens, without generating image tokens. As shown in Figure 6, this SVG-only baseline exhibits failure

Figure 4. **Qualitative comparison on text-to-SVG generation task.** Our DuetSVG aligns better with the text prompts and generates high-quality visual outputs with detailed structures.
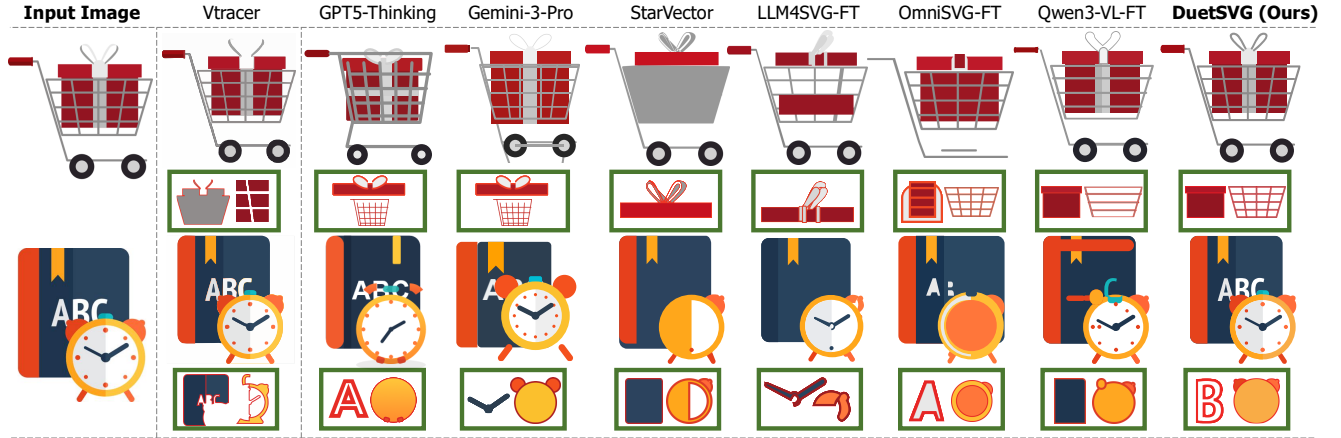


Figure 5. **Qualitative comparison on image-to-SVG conversion task.** Ours preserves more visual details than other previous baselines.

modes similar to prior text-centric VLMs, struggling to produce semantically accurate and structurally coherent SVGs. We further observe that it performs worse than the fine-tuned Qwen3-VL-8B in Table 3, indicating that our backbone is weaker as a pure language model than Qwen3-VL-8B. Yet DuetSVG still outperforms fine-tuned Qwen3-VL-8B when equipped with image modality, showing that vi-

sual guidance is crucial for high-quality SVG generation.

**Ablation of T2I Pretraining.** To validate the effectiveness of the T2I pre-training stage, we conduct an ablation in which this stage is removed. In Table 3, we observe that T2I pre-training helps the model acquire stronger visual and semantic priors, enabling it to generate more visually appealing SVG-style images and more complex SVGs.

7

Table 3. **Ablation study.** (a) Removing visual output brings an notable degradation for all metrics, indicating the effectiveness of our unified multimodal generation design. (b) Without T2I pretraining, we observe a clear drop in T2SVG quality, showing that T2I alignment learned during pretraining is crucial. (c) Test-time scaling further boosts the quality in both T2SVG and I2SVG tasks.

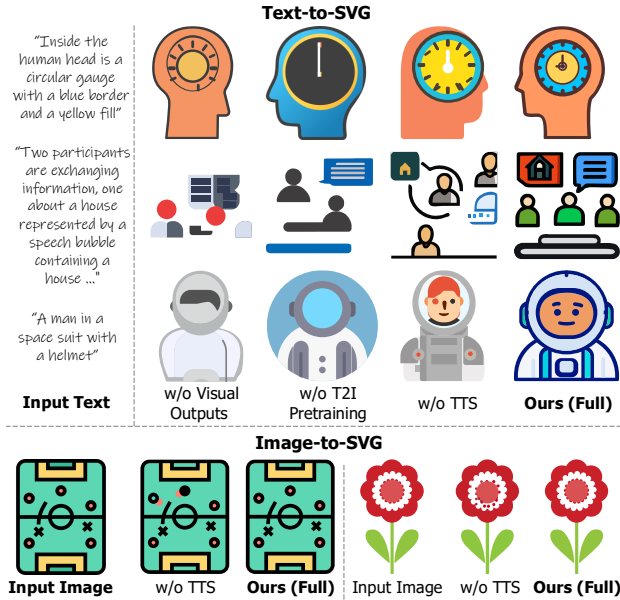| Method | Text-to-SVG | | Image-to-SVG | | |
|---|---|---|---|---|---|
| | FID ↓ | CLIP ↑ | DINO ↑ | SSIM ↑ | LPIPS ↓ |
| w/o Internal Visual Guidance | 51.482 | 23.260 | 0.939 | 0.878 | 0.096 |
| w/o T2I Pretraining | 36.953 | 25.122 | - | - | - |
| w/o Test-time Scaling | 35.066 | 25.584 | 0.955 | 0.920 | 0.082 |
| Ours (Full) | **33.574** | **26.106** | **0.962** | **0.928** | **0.075** |



Figure 6. **Ablation results.** (a) Without **visual output**, SVG-only generation struggles to produce semantically accurate and structurally coherent SVGs. (b) Without **T2I pretraining**, the model lacks strong visual and semantic priors, resulting in simpler and less visually appealing SVGs. (c) Without **test-time scaling (TTS)**, autoregressive decoding may introduce local geometric distortions when handling complex inputs.

**Ablation of TTS with SVG Resampling.** In this ablation study, we evaluate DuetSVG with our test-time scaling (TTS) strategy, DuetSVG with best-of-$N$ TTS, and LLM4SVG with best-of-$N$ TTS. The *CLIP vs Compute* graph is shown in Figure 7. Compared to best-of-$N$, our novel two-stage TTS with image-guided SVG resampling is much more efficient. Text-centric VLMs still perform much worse with TTS enabled.

## 5. Applications

To support downstream applications, DuetSVG can be further fine-tuned for SVG completion and semantic SVG editing. For **SVG completion**, we randomly mask a subset of
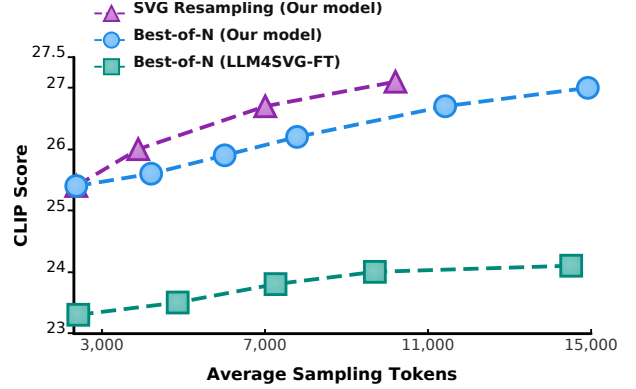


Figure 7. **Sampling efficiency comparison of test-time scaling strategies.** Our strategy leverages image-level candidate selection and SVG resampling to achieve efficient test-time scaling, with substantially lower sampling cost than best-of-$N$.
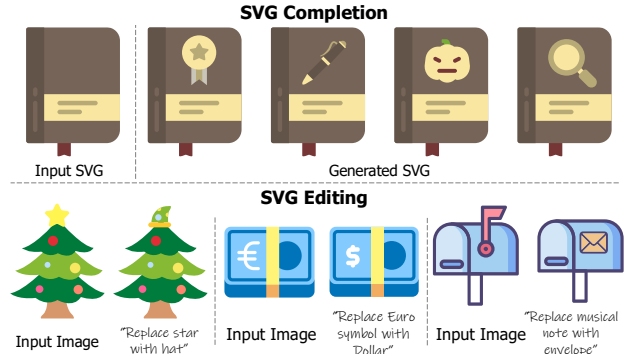


Figure 8. **Downstream applications.** After fine-tuning, DuetSVG supports: (a) *SVG completion*: given the partial SVG and its rendered image, DuetSVG completes coherent and visually appealing SVGs. (b) *SVG editing*: DuetSVG enables instruction-based editing for SVGs.

paths in an SVG and provide the model with both the partial SVG and its rendered image as inputs. The model is trained to infer the complete image and the full SVG script. As shown in Figure 8, DuetSVG can effectively complete coherent and visually appealing SVGs. For semantic **SVG editing**, since no large-scale paired SVG editing dataset is available, we construct a synthetic dataset using a powerful image editing model. Concretely, for each SVG rendering $I_a$ from our dataset and a text instruction, we use Gemini-2.5-Flash [8] to produce an edited image $I_e$ that reflects the instruction. During training, we take $I_e$ and an inverse editing instruction as inputs, and train DuetSVG to reconstruct the original image $I_a$ and its corresponding SVG script. As shown in Figure 8, DuetSVG can perform high-level semantic SVG editing, modifying SVG content according to the text instruction.

## 6. Conclusion

In this paper, we presented DuetSVG, a unified multimodal model that generates both image tokens and SVG tokens rather than treating SVGs as pure text. By combining large-scale T2I pre-training with multi-task SFT on T2I, T2SVG, and I2SVG, DuetSVG leverages rich visual priors to achieve stronger text-SVG alignment and to produce high-quality SVGs. We further introduced a vision-aware test-time scaling strategy that uses the internal visual predictions to guide SVG decoding, improving robustness and reliability. DuetSVG supports a wide range of SVG generation and editing tasks. With advanced VLM techniques, our framework could further accelerate generation, and the proposed training and inference strategies can be used to train larger-scale unified multimodal models (*e.g.*, Emu3.5 [9]), which we leave for future work.

## References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 5

[2] Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, , et al. Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*, 2025. 1, 5

[3] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 1, 5, 2

[4] Mikhail Bessmeltsev and Justin Solomon. Vectorization of line drawings via polyvector fields. *ACM Transactions on Graphics (TOG)*, 38(1):1–12, 2019. 2

[5] Alexandre Carlier, Martin Danelljan, Alexandre Alahi, and Radu Timofte. Deepsvg: A hierarchical generative network for vector graphics animation. *Advances in Neural Information Processing Systems*, 33:16351–16361, 2020. 2

[6] Souymodip Chakraborty, Vineet Batra, Ankit Phogat, Vishwas Jain, Jaswant Singh Ranawat, Sumit Dhingra, Kevin Wampler, Matthew Fisher, and Michal Lukáč. Image vectorization via gradient reconstruction. In *Computer Graphics Forum*, page e70055. Wiley Online Library, 2025. 2

[7] Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. Januspro: Unified multimodal understanding and generation with data and model scaling. *arXiv preprint arXiv:2501.17811*, 2025. 3, 5

[8] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. 5, 8

[9] Yufeng Cui, Honghao Chen, Haoge Deng, Xu Huang, Xinghang Li, Jirong Liu, Yang Liu, Zhuoyan Luo, Jinsheng Wang, Wenxuan Wang, et al. Emu3. 5: Native multimodal models are world learners. *arXiv preprint arXiv:2510.26583*, 2025. 3, 9

[10] Ayan Das, Yongxin Yang, Timothy Hospedales, Tao Xiang, and Yi-Zhe Song. Chirodiff: Modelling chirographic data with diffusion models. *arXiv preprint arXiv:2304.03785*, 2023. 2

[11] Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao Yu, Xiaonan Nie, Ziang Song, et al. Emerging properties in unified multimodal pretraining. *arXiv preprint arXiv:2505.14683*, 2025. 3

[12] Edoardo Alberto Dominici, Nico Schertler, Jonathan Griffin, Shayan Hoshyari, Leonid Sigal, and Alla Sheffer. Polyfit: Perception-aligned vectorization of raster clip-art via intermediate polygonal fitting. *ACM Transactions on Graphics (TOG)*, 39(4):77–1, 2020. 2

[13] Jean-Dominique Favreau, Florent Lafarge, and Adrien Bousseau. Photo2clipart: Image abstraction and vectorization using layered linear gradients. *ACM Transactions on Graphics (TOG)*, 36(6):1–11, 2017. 2

[14] Kevin Frans, Lisa Soros, and Olaf Witkowski. Clipdraw: Exploring text-to-drawing synthesis through language-image encoders. *Advances in Neural Information Processing Systems*, 35:5207–5218, 2022. 2

[15] Rinon Gal, Yael Vinker, Yuval Alaluf, Amit H Bermano, Daniel Cohen-Or, Ariel Shamir, and Gal Chechik. Breathing life into sketches using text-to-video priors. *arXiv preprint arXiv:2311.13608*, 2023. 2

[16] Google DeepMind. Gemini 3 Pro. https://deepmind.google/models/gemini/pro, 2025. 5

[17] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*, 2017. 5

[18] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 4

[19] Shayan Hoshyari, Edoardo Alberto Dominici, Alla Sheffer, Nathan Carr, Zhaowen Wang, Duygu Ceylan, and I-Chao Shen. Perception-driven semi-structured boundary vectorization. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018. 2

[20] Shir Iluz, Yael Vinker, Amir Hertz, Daniel Berio, Daniel Cohen-Or, and Ariel Shamir. Word-as-image for semantic typography. *arXiv preprint arXiv:2303.01818*, 2023. 2

[21] Ajay Jain, Amber Xie, and Pieter Abbeel. Vectorfusion: Text-to-svg by abstracting pixel-based diffusion models. *arXiv preprint arXiv:2211.11319*, 2022. 2, 5

[22] Johannes Kopf and Dani Lischinski. Depixelizing pixel art. In *ACM SIGGRAPH 2011 papers*, pages 1–8. 2011. 2

[23] Black Forest Labs, Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril Diagne, Tim Dockhorn, Jack English, Zion English, Patrick Esser, et al. Flux. 1 kontext: Flow matching for in-context image generation and editing in latent space. *arXiv preprint arXiv:2506.15742*, 2025. 3, 5, 1

[24] Tzu-Mao Li, Michal Lukáč, Michaël Gharbi, and Jonathan Ragan-Kelley. Differentiable vector graphics rasterization

for editing and learning. *ACM Transactions on Graphics (TOG)*, 39(6):1–15, 2020. 2

[25] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019. 5

[26] Xu Ma, Yuqian Zhou, Xingqian Xu, Bin Sun, Valerii Filev, Nikita Orlov, Yun Fu, and Humphrey Shi. Towards layer-wise image vectorization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16314–16323, 2022. 2

[27] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 5

[28] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Milden-hall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 2

[29] Sanford Pun and Chris Tsang. VTracer. https://github.com/visioncortex/vtracer, 2025. 5

[30] 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 *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2, 4, 5

[31] Juan A Rodriguez, Shubham Agarwal, Issam H Laradji, Pau Rodriguez, David Vazquez, Christopher Pal, and Marco Pedersoli. Starvector: Generating scalable vector graphics code from images. *arXiv preprint arXiv:2312.11556*, 2023. 1, 2, 3

[32] Juan A Rodriguez, Abhay Puri, Shubham Agarwal, Issam H Laradji, Pau Rodriguez, Sai Rajeswar, David Vazquez, Christopher Pal, and Marco Pedersoli. Starvector: Generating scalable vector graphics code from images and text. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 16175–16186, 2025. 5

[33] Juan A Rodriguez, Haotian Zhang, Abhay Puri, Aarash Feizi, Rishav Pramanik, Pascal Wichmann, Arnab Mondal, Mohammad Reza Samsami, Rabiul Awal, Perouz Taslakian, et al. Rendering-aware reinforcement learning for vector graphics generation. *arXiv preprint arXiv:2505.20793*, 2025. 2

[34] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 2

[35] Peter Schaldenbrand, Zhixuan Liu, and Jean Oh. Styleclip-draw: Coupling content and style in text-to-drawing translation. *arXiv preprint arXiv:2202.12362*, 2022. 2

[36] Christoph Schuhmann. Improved aesthetic predictor. https://github.com/christophschuhmann/improved-aesthetic-predictor, 2022. 5

[37] Peter Selinger. Potrace: a polygon-based tracing algorithm, 2003. 2, 5

[38] Yiren Song, Xning Shao, Kang Chen, Weidong Zhang, Minzhe Li, and Zhongliang Jing. Clipvg: Text-guided im-age manipulation using differentiable vector graphics. *arXiv preprint arXiv:2212.02122*, 2022. 2

[39] Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan. Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint arXiv:2406.06525*, 2024. 3

[40] Zecheng Tang, Chenfei Wu, Zekai Zhang, Mingheng Ni, Shengming Yin, Yu Liu, Zhengyuan Yang, Lijuan Wang, Zicheng Liu, Juntao Li, et al. Strokenuwa: Tokeniz-ing strokes for vector graphic synthesis. *arXiv preprint arXiv:2401.17093*, 2024. 2

[41] Vikas Thamizharasan, Difan Liu, Shantanu Agarwal, Matthew Fisher, Michaël Gharbi, Oliver Wang, Alec Jacob-son, and Evangelos Kalogerakis. Vecfusion: Vector font gen-eration with diffusion. In *Proceedings of the IEEE/CVF Con-ference on Computer Vision and Pattern Recognition*, pages 7943–7952, 2024. 2

[42] Xingze Tian and Tobias Günther. A survey of smooth vec-tor graphics: Recent advances in repr esentation, creation, rasterization, and image vectorization. *IEEE Transactions on Visualization and Computer Graphics*, 30(3):1652–1671, 2022. 2

[43] Yael Vinker, Ehsan Pajouheshgar, Jessica Y Bo, Ro-man Christian Bachmann, Amit Haim Bermano, Daniel Cohen-Or, Amir Zamir, and Ariel Shamir. Clipasso: Semantically-aware object sketching. *ACM Transactions on Graphics (TOG)*, 41(4):1–11, 2022. 2

[44] Haomin Wang, Jinhui Yin, Qi Wei, Wenguang Zeng, Lixin Gu, Shenglong Ye, Zhangwei Gao, Yaohui Wang, Yanting Zhang, Yuanqi Li, et al. Internsvg: Towards unified svg tasks with multimodal large language models. *arXiv preprint arXiv:2510.11341*, 2025. 1, 2, 5, 6

[45] Jiuniu Wang, Gongjie Zhang, Quanhao Qian, Junlong Gao, Deli Zhao, and Ran Xu. Robosvg: A unified framework for interactive svg generation with multi-modal guidance. *arXiv preprint arXiv:2510.22684*, 2025. 2

[46] Yizhi Wang and Zhouhui Lian. Deepvecfont: Synthesizing high-quality vector fonts via dual-modality learning. *ACM Transactions on Graphics (TOG)*, 40(6):1–15, 2021. 2

[47] Yuqing Wang, Yizhi Wang, Longhui Yu, Yuesheng Zhu, and Zhouhui Lian. Deepvecfont-v2: Exploiting transformers to synthesize vector fonts with higher quality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18320–18328, 2023. 2

[48] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Si-moncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 5

[49] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distilla-tion. *arXiv preprint arXiv:2305.16213*, 2023. 2

[50] Ronghuan Wu, Wanchao Su, Kede Ma, and Jing Liao. Icon-shop: Text-guided vector icon synthesis with autoregressive transformers. *ACM Transactions on Graphics (TOG)*, 42(6): 1–14, 2023. 2, 5, 1

[51] Ronghuan Wu, Wanchao Su, Kede Ma, and Jing Liao. Aniclipart: Clipart animation with text-to-video priors. *International Journal of Computer Vision*, pages 1–17, 2024. 2

[52] Ronghuan Wu, Wanchao Su, and Jing Liao. Chat2svg: Vector graphics generation with large language models and image diffusion models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 23690–23700, 2025. 2

[53] Jinheng Xie, Zhenheng Yang, and Mike Zheng Shou. Showo2: Improved native unified multimodal models. *arXiv preprint arXiv:2506.15564*, 2025. 3

[54] Ximing Xing, Chuang Wang, Haitao Zhou, Jing Zhang, Qian Yu, and Dong Xu. Diffsketcher: Text guided vector sketch synthesis through latent diffusion models. *arXiv preprint arXiv:2306.14685*, 2023. 2

[55] Ximing Xing, Haitao Zhou, Chuang Wang, Jing Zhang, Dong Xu, and Qian Yu. Svgdreamer: Text guided svg generation with diffusion model. *arXiv preprint arXiv:2312.16476*, 2023. 2, 5

[56] Ximing Xing, Yandong Guan, Jing Zhang, Dong Xu, and Qian Yu. Reason-svg: Hybrid reward rl for ahamoments in vector graphics generation. *arXiv preprint arXiv:2505.24499*, 2025. 2

[57] Ximing Xing, Juncheng Hu, Guotao Liang, Jing Zhang, Dong Xu, and Qian Yu. Empowering llms to understand and generate complex vector graphics. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 19487–19497, 2025. 1, 5

[58] Yiying Yang, Wei Cheng, Sijin Chen, Xianfang Zeng, Fukun Yin, Jiaxu Zhang, Liao Wang, Gang Yu, Xingjun Ma, and Yu-Gang Jiang. Omnisvg: A unified scalable vector graphics generation model. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. 1, 2, 3, 5

[59] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11975–11986, 2023. 3

[60] Peiying Zhang, Nanxuan Zhao, and Jing Liao. Text-to-vector generation with neural path representation. *ACM Transactions on Graphics (TOG)*, 43(4):1–13, 2024. 2, 5

[61] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 4, 5

[62] Xinjie Zhang, Jintao Guo, Shanshan Zhao, Minghao Fu, Lunhao Duan, Jiakui Hu, Yong Xien Chng, Guo-Hua Wang, Qing-Guo Chen, Zhao Xu, et al. Unified multimodal understanding and generation models: Advances, challenges, and opportunities. *arXiv preprint arXiv:2505.02567*, 2025. 3

[63] Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren Zhou. Qwen3 embedding: Advancing text embedding and reranking through foundation models. *arXiv preprint arXiv:2506.05176*, 2025. 5

[64] Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025. 1, 5, 2

# DuetSVG: Unified Multimodal SVG Generation with Internal Visual Guidance

## Supplementary Material

## A. Overview

In this supplementary material, we provide additional details and evaluations, including:
- Generalization evaluation (Section B.1).
- Additional comparisons and results (Section B.2).
- Limitation (Section C).
- Details of the captioning pipeline (Section D).

## B. Additional Evaluation

### B.1. Generalization Evaluation

We evaluate our model's generalization capability using *Novelty* and *Uniqueness* scores following IconShop [50]. Two SVGs are considered identical if the cosine similarity between their CLIP image embeddings is at least 0.98. Under this criterion, *Uniqueness* is the fraction of generated samples that appear exactly once within the generated set. *Novelty* is the fraction of generated samples that have no match in the *SVG-Hub-5M* training corpus: for each generated SVG, we render it to an image, retrieve its nearest neighbor in *SVG-Hub-5M* using CLIP cosine similarity, and label the sample as novel if the maximum similarity is below 0.98.

We randomly select 500 text prompts and generate three SVGs per prompt with different seeds, yielding 1,500 samples. On this set, our model attains 99.5% *Novelty* and 99.8% *Uniqueness*. Figure 9 shows generated SVGs alongside their nearest neighbors from *SVG-Hub-5M*, illustrating our model's ability to produce novel and diverse results.

### B.2. Additional Comparisons and Results

**Additional Results on SVG-Hub-1M.** We provide additional quantitative results on the *SVG-Hub-1M* dataset in Table 4. As expected, the smaller training corpus leads to an overall degradation in performance across metrics, yet our method consistently outperforms prior approaches on all metrics. As the dataset size increases, our model improves substantially, whereas other methods show only limited improvement on the text-to-SVG task. The results further highlight the effectiveness of our unified multimodal SVG generation model.

**Additional Comparisons with Flux + I2SVG.** For the T2SVG task, we compare against a two-stage baseline that first synthesizes images with FLUX.1-dev [23] and then performs image-to-SVG using a VLM-based model, Qwen3-VL-8B [2], fine-tuned on *SVG-Hub-5M* dataset. We



"A directional signpost with two yellow arrows on a purple pole, featuring a green spherical top and a blue base with a green bush beside it"

"A stylized illustration of a bearded king wearing a crown"

"A flat, two-dimensional illustration of a cash register with a display screen above and a keypad below"
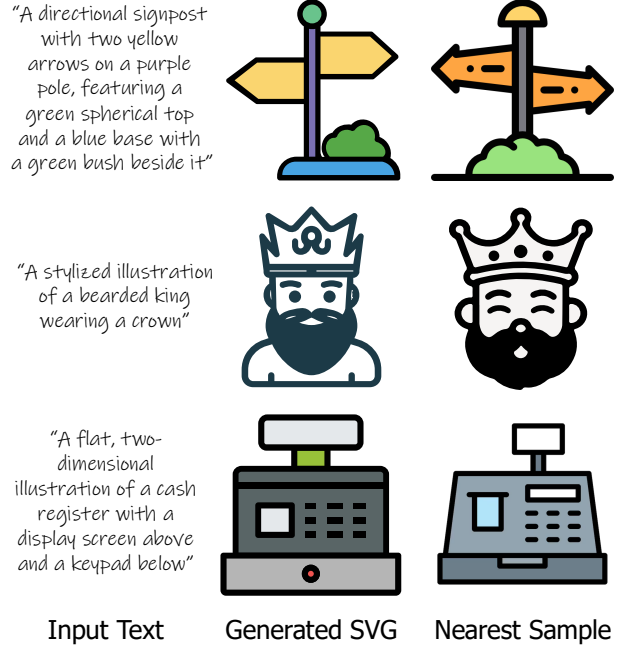
Input Text    Generated SVG    Nearest Sample

Figure 9. Generated SVGs and their nearest neighbor samples.

Table 4. **Quantitative comparison with VLM-based baselines on the SArena-Icon benchmark [44].** In this setting, "Ours-7B" and all baselines marked "FT" are fine-tuned on *SVG-Hub-1M* dataset.

| Method | Text-to-SVG | | | Image-to-SVG | | |
|---|---|---|---|---|---|---|
| | FID ↓ | FID-C ↓ | CLIP ↑ | DINO ↑ | SSIM ↑ | LPIPS ↓ |
| StarVector-8B (w/o FT) | - | - | - | 0.871 | 0.623 | 0.206 |
| StarVector-8B (FT) | - | - | - | 0.895 | 0.742 | 0.182 |
| LLM4SVG-7B (FT) | 20.896 | 7.980 | 22.108 | 0.898 | 0.760 | 0.164 |
| OmniSVG-3B (w/o FT) | 28.292 | 11.318 | 21.679 | 0.894 | 0.756 | 0.186 |
| OmniSVG-3B (FT) | 24.977 | 9.659 | 21.825 | 0.902 | 0.761 | 0.178 |
| Qwen3-VL-8B (FT) | 19.760 | 7.572 | 22.126 | 0.908 | 0.774 | 0.162 |
| Ours-7B (w/o TTS) | 17.265 | 6.388 | 22.586 | 0.917 | 0.808 | 0.151 |
| Ours-7B (TTS) | **16.952** | **6.242** | **22.645** | **0.924** | **0.814** | **0.145** |

observe cross-model inconsistencies: the synthesized intermediate images differ in style and fine detail from the real SVG images used during I2SVG training, introducing a train-test mismatch for the VLM-based I2SVG model. As a result, the I2SVG model does not generalize well to these intermediates, and the resulting SVGs often show reduced alignment with the intermediate images (see Figure 10). In contrast, our unified multimodal generative model co-generates image and SVG tokens within a single end-to-end architecture, enabling the use of large-scale text-image data and providing visual grounding during SVG decoding.
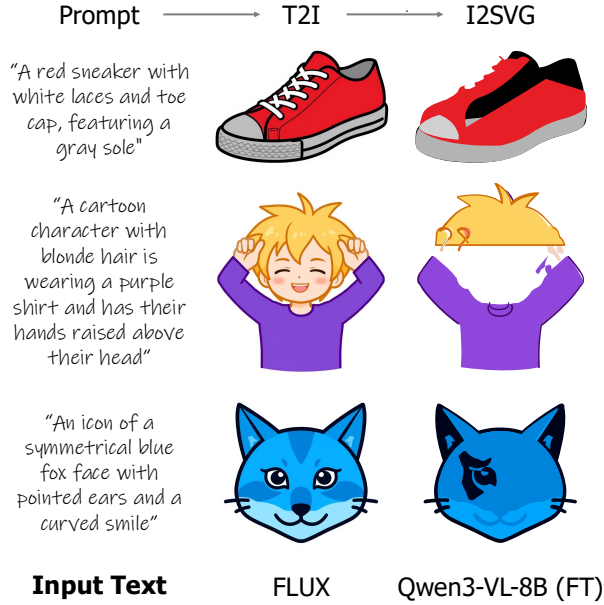
Figure 10. Additional comparisons with Flux + I2SVG pipeline.
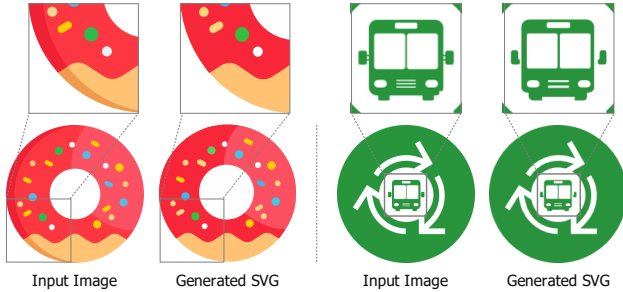


Figure 11. Failure cases.

## C. Limitation

While our model excels at SVG generation, it has limitations. When the input image contains very fine details and rich color variation, the generated SVG may miss small structures and exhibit mild color shifts (see Figure 11). A potential mitigation is a dynamic high-resolution strategy [3] that adaptively increases the number of patches fed to the vision encoder, which can improve the capture of fine details and color consistency. We plan to investigate this in future work.

## D. Details of the captioning pipeline

To enable text-to-SVG (T2SVG) training from semantically rich prompts, we rasterize each SVG and use open-source VLMs (InternVL3 [64] and Qwen2.5-VL [3]) to generate captions at three levels of detail. We then perform cross-model verification and refinement: a caption produced by

one VLM is evaluated against the rendered SVG by the other, which flags inaccuracies and suggests edits; the revised caption is subsequently adopted. The prompt templates for the three levels are provided in Figure 12, Figure 13, and Figure 14.

```
You are an expert in semantic analysis and
captioning for SVG icons and flat 2D
illustrations. Analyze the provided image and
produce a concise, descriptive title suitable
for an SVG-style icon or flat illustration.

**Guidelines**
- Write **one compact description** (a short
phrase or a brief sentence) in clear, precise
English.
- Focus on the **primary subject or action**.
Mention **dominant color(s)** and other key
attributes **only if they are essential to
the meaning**.
- Optionally add a simple **layout/background
cue** only when it is crucial (e.g., "on an
orange saucer," "within a circle").
- Be specific and accurate about the
subject's semantics.
- Output only the description text. No extra
words or formatting.

Describe this SVG-style image following the
rules above. Output only the sentence or
phrase.
```

Figure 12. Short prompt template.

```
You are an expert captioner for SVG-style
icons and flat 2D illustrations.

**Task**
Given an image, write **1-2 concise English
sentences** that **enumerate the key semantic
elements** and describe their **basic layout
and style**.

**Guidelines**
- Start with the **main subject/action**,
then briefly list important secondary
elements.
- Mention **spatial arrangement** only when
helpful (e.g., centered, above/beside, inside
a circle).
- Include **dominant colors and notable style
cues** (flat/outline/pastel/thick strokes,
etc.) when they are visually salient.
- Stay **compact but informative**; do not
describe fine-grained geometry.

**Output**
Return **only the sentences**.
```

Figure 13. Medium description template.

2

```
You are a technical analyst for vector
graphics.

**Task**
Analyze the provided image and write a
**precise, exhaustive description in 3-6
English sentences (one paragraph)**.

**Guidelines**
- **Identify all elements and semantics:**
systematically mention every visible
component and its intended meaning.
- **Shape decomposition and geometry:**
describe each element's exact shapes (e.g.,
circles, polygons, curved paths) and key
geometric traits (rounded vs. sharp corners,
symmetry, thickness, etc.).
- **Spatial relations and layout:** state
relative positions, alignment, scale,
layering, overlaps, and foreground/background
structure.
- **Styling and color per element:** name the
specific colors of each part and any fills or
shading.
- **Stylistic details:** note
outlines/strokes and their weight, flat fills,
gradients, highlights, shadows, textures, or
transparency when present.
- **Text transcription:** copy any legible
text exactly and specify its placement.

**Output format**
Use plain sentences only (no lists or
bullets). Return **only** the description
paragraph.
```

Figure 14. Detailed annotation template.