## OmniSVG: A Unified Scalable Vector Graphics Generation Model

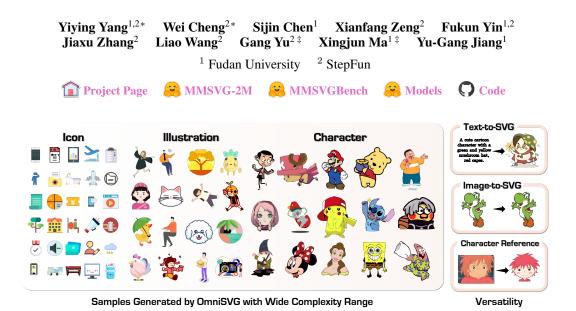


Figure 1: **OmniSVG** is capable of autoregressively generating high-quality Scalable Vector Graphs (SVG) across a wide spectrum of complexity, from simple icons to intricate anime characters. OmniSVG demonstrates remarkable versatility in generating high-quality SVGs adhering to multimodal instructions, covering tasks like Text-to-SVG, Image-to-SVG, and Character-Reference SVG, making it a powerful and flexible solution for diverse creative tasks.

## Abstract

Scalable Vector Graphics (SVG) is an important image format widely adopted in graphic design because of their resolution independence and editability. The development of autonomous SVG generation workflows is continuously drawing attention from both designers and researchers in the AIGC community. However, existing methods either produce unstructured outputs at huge computational cost or are limited to generating monochrome icons of over-simplified structures. To produce high-quality and complex SVG adhering to multi-modal instructions, we propose OmniSVG, a unified SVG generation framework that inherits knowledge from a pre-trained Vision-Language Model (VLM). By parameterizing SVG commands and coordinates into discrete token sequences, the auto-regressive nature enables us to seaminglessly adapt modern VLMs to the direct SVG generation. To further advance the development of SVG synthesis, we introduce MMSVG-2M, a multimodal dataset with two million richly annotated SVG assets, along with a standardized evaluation protocol for conditional SVG generation tasks. Extensive experiments show that OmniSVG outperforms existing methods and demonstrates its potential for integration into professional SVG design workflows.

<sup>\*</sup> Yiying Yang and Wei Cheng contributed equally to this work.

<sup>‡</sup> Corresponding Authors.

## 1 Introduction

Scalable Vector Graphics (SVG) have become a cornerstone of modern digital design because of their resolution independence, compact file size, and inherent editability. Widely adopted in professional workflows from UI/UX design to industrial CAD systems, SVG enables precise manipulation of geometric primitives (*e.g.*, Bézier curves, polygons) while maintaining high precision and consistent visual quality across varying resolutions. However, creating high-quality SVG content remains challenging for non-experts, requiring mastery of specialized tools or intricate XML syntax.

Existing methods adopt either optimization-based methods or auto-regressive approaches to generate SVG contents. The optimization-based methods [34, 12, 29] iteratively refine the SVG parameters by minimizing the differences between the input image and the raster image created by differentiable vector graphics rasterizers. Though these methods are sufficient for reconstructing SVG icons, they suffer from significant computational overhead when scaling up to more intricate samples and produce unstructured outputs with redundant anchor points, harming the editability of the reconstructed SVG samples. In contrast, auto-regressive methods build transformer models or adapt pre-trained Large Language Models (LLMs) to directly generate XML parameters [59] or codes [56, 42] representing SVGs. Benefiting from the end-to-end learning pipeline, the auto-regressive method is a more scalable approach [5] as it can learn directly from a large collection of SVG samples. However, existing auto-regressive approaches are limited to basic SVG contents [11, 24, 53] because of the limited context length and the scarcity of complex SVG data.

In this paper, we propose **OmniSVG** that harnesses native VLMs [1] for various end-to-end multimodal SVG generation tasks. By parameterizing SVG coordinates and commands into discrete tokens, OmniSVG decouples structural logic from low-level geometry, mitigating the "coordinate hallucination" problem prevalent in code-based LLMs, and produces vivid and colorful SVG results. Additionally, the next token prediction training objective enables OmniSVG to complete SVGs with diverse generation results given some partial observations. Compared to traditional auto-regressive SVG generation methods, OmniSVG is able to parameterize SVGs exceeding 30k tokens, facilitating the generation of detailed and complex SVG contents. Building upon pre-trained VLMs, our method natively integrates the ability to reason upon visual and textual instructions to synthesize editable, high-fidelity SVGs across diverse domains, from icons to intricate illustrations and anime characters.

To advance the development of SVG synthesis, we introduce MMSVG-2M, a multi-modal SVG synthesis dataset with two million richly annotated assets, encompassing icons, illustrations, and anime designs. We also establish a standardized evaluation protocol, MMSVG-Bench, for "Text-to-SVG" and "Image-to-SVG" generation. Extensive experiments show that OmniSVG can produce highly detailed and complex SVG contents, surpassing prior art both quantitatively and qualitatively.

To summarize, our key contributions include:

- We introduce OmniSVG, a family of end-to-end multimodal SVG generators that leverage native VLMs for generating complex and detailed SVGs, from simple icons to intricate anime characters.
- We present MMSVG-2M, a large-scale dataset comprising two million SVG assets, along
  with a standardized evaluation protocol for various multi-modal SVG generation tasks,
  providing a comprehensive resource for future research.
- Extensive experiments show that OmniSVG surpasses prior SVG generation methods both qualitatively and quantitatively, highlighting its potential for integration into professional SVG design workflows.

## 2 Related Works

**SVG Generation**. Early attempts to generating SVGs directly utilize architectures like RNNs [18, 41, 19, 44, 45], VAEs [4, 32, 48, 46, 51], and Transformers [4, 57] to compress SVG commands into latent representations. Meanwhile, DeepSVG [4] further parameterizes SVGs using a dual transformer architecture but struggles with geometric consistency. Recently, the advent of large language models (LLMs) [30, 64, 52, 61, 5, 6, 63, 62, 49] unleashes the potential of generating SVGs via XML code synthesis [59, 56, 42]. However, the limited context length of existing LLM-based SVG generation methods [56, 42, 59] poses significant challenges in handling complex SVGs that exceed 10k tokens.

In this paper, we explore the potential of native Vision-Language Models (VLMs) in multi-modal SVG generation. By combining pre-trained VLMs with SVG command parameterization, we validate that OmniSVG is able to follow multi-modal instructions and generate vivid and complex SVGs.

**Image Vectorization**. Recent advancements in vectorization harness diffusion models paired with differentiable rasterizers, using techniques like score distillation sampling [37, 22, 7] and specialized regularizers [29, 34] to convert raster images into SVG paths. While these methods achieve remarkable results, they face limitations such as over-smoothing, color over-saturation, and lack of editability, often producing tangled paths that fail to capture hierarchical structures inherent in professional SVG designs. In this paper, we present an end-to-end approach that follows multi-modal instructions to generate high-quality SVGs with improved path clarity and editability.

**SVG Datasets and Benchmarks**. The lack of suitable datasets for complex SVG structures presents a significant challenge. Existing datasets [11, 24, 53] primarily focus on simplified path-based SVGs or monochrome icons, overlooking the intricate layered structures and rich color semantics found in real-world designs. For example, FIGR-8-SVG [11] focuses on monochromatic icons, while StarVector [42] proposes categorized datasets, including illustrations, icons, emojis, and fonts. Therefore, existing datasets only present simple SVG samples that do not exceed 8.2k tokens, failing to capture the complexities of layered structures and rich color semantics. Benchmark evaluations, such as VGBench [70], further highlight gaps in multi-format testing and the absence of comprehensive coverage for illustrative SVGs. To this end, we introduce MMSVG-2M, a multimodal SVG synthesis dataset comprising two million richly annotated assets, including icons, illustrations, and complex anime designs. We also present a standardized evaluation protocol, MMSVG-Bench, to evaluate diverse multi-modal SVG generation tasks with varying complexity.

#### 3 OmniSVG Dataset

We present MMSVG-2M, a large-scale SVG dataset with two million SVG samples covering website icons, illustrations, graphic designs, anime characters, and etc (Sec. 3.1). To promote the downstream development of SVG generation methods, we also introduce MMSVG-Bench, a standardized evaluation protocol for a series of multi-modal instruction following tasks for conditional SVG generation (Sec. 3.2).

#### 3.1 MMSVG-2M

**Data Source.** With increasing visual complexity, MMSVG-2M consists of three subsets, 1) the icon subset MMSVG-Icon collected from Iconfont, 2) the illustration subset MMSVG-Illustration sourced from IconSount, and 3) the complex anime character subset MMSVG-Character both curated from Freepik and created by our data creation pipeline as shown in Fig. 2. All these websites are online platforms where users can publish and share SVGs, encompassing a broad variety of categories. Specifically, our collection of MMSVG-2M contains 1.1 million icons, 0.5 million illustrations, and 0.4 million anime characters as shown in Tab. 6.

**Data Curation.** We extract SVG samples with a comprehensive deduplication process based on filenames, SVG code, and metadata. We first fit the collected SVGs within a viewbox of  $200 \times 200$ . Then, we employ an off-the-shelf VLM, specifically BLIP-2 [28], to generate captions for the SVGs. Please find more samples from the MMSVG-2M dataset in Fig. 8, and instruction templates in Appendix A.2.

**SVG** Simplification is an essential procedure in SVG data cleansing, since the over-complicated XML grammars in the crawled SVG data will lead to ambiguities while representing basic shapes. To standardize training and evaluation, we simplify all SVG commands with atomic commands as shown in Tab. 1. Inspired by FIGR-8-SVG [11] and IconShop [57], we remove all attributes and simplify each SVG with five basic commands, including "Move To" (M), "Line To" (L), "Cubic Bézier" (C), "Elliptical Arc" (A), "ClosePath" (Z). The introduction of atomic commands further removes the ambiguities, as complex XML grammars can be approximated with the combination of several atomic commands. To efficiently produce a unified and less complex data structure, we utilize picosvg to remove grammars like "group" and "transform", and simplify the complex commands

Table 1: SVG Draw Commands. Draw commands used in this work along with their arguments and a visualization are listed. The start-position  $(x_1, y_1)$  is implicitly defined as the end-position of the preceding command.

Command	Arguments	Description	Visualization
<sop></sop>	Ø	'Start-of-Path' token.	
M (MoveTo)	$x_2, y_2$	Move the cursor to the end-point $(x_2, y_2)$ without drawing anything.	$(x_2,y_2)$
L (LineTo)	$x_2, y_2$	Draw a line to the point $(x_2, y_2)$ .	$(x_1,y_1)$ $(x_2,y_2)$
C (Cubic Bézier)	$q_{x1}, q_{y1}$ $q_{x2}, q_{y2}$ $x_2, y_2$	Draw a cubic Bézier curve with control points $(q_{x1}, q_{y1})$ , $(q_{x2}, q_{y2})$ and end-point $(x_2, y_2)$ .	$(x_1,y_1)$ $(q_{xz},q_{yz})$ $(q_{xz},q_{yz})$ $(x_2,y_2)$
A (Elliptical Arc)	$r_x, r_y \\ \varphi, f_A, f_S \\ x_2, y_2$	Draw an elliptical arc with radii $r_x$ and $r_y$ (semi-major and semi-minor axes), rotated by angle $\varphi$ to the x-axis, and end-point $(x_2, y_2)$ . $(x_2, y_2)$ .	$\begin{cases} f_1 = 1 \\ f_2 = 1 \\ p \\ (x_1, y_2) \end{cases}$
Z (ClosePath)	Ø	Close the path by moving the cursor back to the path's starting position $(x_0, y_0)$ .	$(x_0,y_0)$ $(x_1,y_1)$
F (Fill)	fill	Draw the fill attribute of the path.	Ø
<eos></eos>	Ø	'End-of-SVG' token.	

to atomic path commands. It is worth noting that atomic path commands are sufficient to represent complex SVGs shown in Fig. 1.

#### 3.2 MMSVG-Bench

To compensate for the vacancy of standardized and open evaluation for SVG generation, we introduce **MMSVG-Bench**, a comprehensive benchmark for multi-modal SVG generation. We require the corresponding benchmark to be a sufficient verification whether a model is practically useful in real-world scenarios, and avoid the excessive similarity between the benchmark input data and training data as in traditional train/test splits. Therefore, we opt to generate the benchmark inputs with GPT-40. Specifically, for Text-to-SVG task, we synthesize 150 textual prompts for each SVG subset (*i.e.* Icon and Illustration). For Image-to-SVG task, we synthesize extra 150 textual descriptions, and prompt GPT-40 to generate vector-style images with transparent backgrounds based on the above texts as the ground truth visual samples. We focus on both the visual quality and semantics of the generation results.

**Text-to-SVG** requires a model to generate SVGs from text instructions. We measure the visual quality with Frechet Inception Distance (FID) [50], aesthetic appeal with Aesthetic score [43], text-SVG alignment with CLIP score [38], and Human Preference Scores (HPS) [58].

**Image-to-SVG** evaluates a model's ability to convert images into SVGs. To quantify the distance between the input and output SVG, we calculate the cosine similarity of DinoV2 features (DinoScore) [35], Structural Similarity Index (SSIM) [54], Learned Perceptual Image Patch Similarity (LPIPS) [66], and Mean Squared Error (MSE).

Character-Reference SVG Generation evaluates a model's ability to generate novel SVGs while keeping the profile of the characters depicted in the input image. Different from image-to-SVG, the model does not reconstruct, but generates a specific character SVG for the input image (see Fig. 5). We evaluate the alignment between input character images and generated SVGs by prompting GPT-40 [21] to generate a score ranging from 1 to 10, the higher the better. [15, 23, 17]

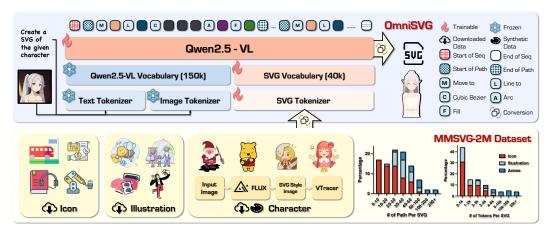


Figure 2: **Overview of OmniSVG**. OmniSVG is built on a pre-trained vision-language model Qwen2.5-VL and incorporates an SVG tokenizer. The model tokenizes both text and image inputs as prefix tokens, while the SVG tokenizer encodes vector graphics commands into a unified representation space.

### 4 OmniSVG

To support end-to-end training for multi-modal SVG generation, OmniSVG parameterizes a series of atomic SVG path commands into a sequence before feeding into a pre-trained VLM with multi-modal instructions.

**SVG Tokenizer.** As illustrated in Sec. 3, our MMSVG-2M dataset simplifies an SVG by removing all attributes and using five basic path commands (see Tab. 1). After the simplification, an SVG script G is represented as the combination of M paths,  $G = \{P_i\}_{i=1}^M$ . Here,  $P_i$  is the i-th path containing  $N_i$  commands,  $P_i = \{C_i^j\}_{j=1}^{N_i}$ , where  $C_i^j$  is the j-th command in the i-th path. Each command is represented as  $C_i^j = (U_i^j, V_i^j)$ , containing both the command type identifier  $U_i^j \in \{M, L, C, A, Z\}$  and the corresponding location argument  $V_i^j$ . To generate colored SVG contents, we assign special tokens for hex values to control the "Fill" (F) attribute, distinguishing it from the original SVG commands and coordinates. To this end, we are able to use a total six types of commands  $U_i^j \in \{M, L, C, A, Z, F\}$  to parameterize a colored SVG parameterization.

Specifically, our SVG tokenizer transforms SVG scripts  $X_s$  into an ordered SVG token sequence within the same representation space as the pre-trained VLM. Following IconShop [57], we flatten the layered structure of the SVG script by concatenating different paths into a single command sequence, where each path begins with the drawing commands followed by point coordinates. Therefore, each SVG sequence could be represented as a flattened sequence. As the generation identifier, we apply special tokens like <S0P> and <E0S> to the two ends of a SVG sequence, identifying the begining and ending of a SVG sequence. We assign special tokens for each command type, i.e.  $\{M, L, C, A, Z, F\}$ . To shorten the length of the SVG sequence, we further merge the 2D point coordinates into one token with a mapping function:  $\langle x, y \rangle \rightarrow x \times w + y$ , where w is the width of the image. The SVG sequence are then lifted into the same embedding space as the pre-trained VLM with a learnable embedding layer.

**Model Architecture.** OmniSVG adopts Qwen2.5-VL [1], an open-sourced VLM that excels at understanding intricate vision-text inputs, as its backbone (Fig. 2) to produce precise and compact SVG outputs. OmniSVG is trained to predict the SVG suffix tokens ( $x_s$ ) conditioned on the mutlimodal instruction prefix tokens ( $x_s$ ) with the standard next-token prediction objective.

$$\theta^* = \arg\max_{\theta} \prod_{i=1}^{L} P\left(x_{s,i} \mid x_{s,< i}, x_c\right) \tag{1}$$

Table 2: **Quantitative Evaluations.** Quantitative results between OmniSVG and current state-of-the-art text-to-SVG and image-to-SVG baseline methods. The bold numbers and underlined numbers represents the best and second best performance repectively. Our OmniSVG model demonstrates superior performance compared SOTA SVG generation baselines.

Evaluation Dataset	Methods	# Tokens		Text	to-SVG		Image-to-SVG			
Evaluation Dataset	Methods		FID↓	CLIP↑	Aesthetic†	HPS↑	DINO↑	SSIM↑	LPIPS↓	MSE↓
	Vectorfusion [22]	66.2k	250.77	0.240	4.76	0.237	-	-	-	-
	SVGDreamer [60]	132.0k	308.94	0.207	4.26	0.221	-	-	-	-
	Chat2SVG [56]	0.6k	190.87	0.299	4.41	0.247	-	-	-	-
	IconShop [57]	2.0k	213.28	0.288	4.55	0.244	-	-	-	-
	LIVE [34]	52.5k	-	-	-	-	0.932	0.943	0.106	0.011
MMSVG-Icon	DiffVG [29]	322.0k	-	-	-	-	0.940	0.954	0.066	0.002
	GPT-4o [21]	0.3k	-	-	-	-	0.860	0.792	0.403	0.124
	StarVector(8B) [42]	2.0k	-	-	-	-	0.895	0.881	0.231	0.059
	Vtracer	52.4k	-	-	-	-	0.993	0.966	0.039	0.002
	OmniSVG(4B)	3.8k	137.40	0.275	4.62	0.244	0.993	0.950	0.050	0.006
	OmniSVG-L(8B)	5.7k	130.56	0.276	4.60	0.242	0.922	0.893	0.235	0.040
	Vectorfusion [22]	66.1k	253.94	0.185	4.94	0.226	-	-	-	-
	SVGDreamer [60]	132.0k	419.70	0.201	4.37	0.221	-	-	-	-
	Chat2SVG [56]	1.0k	210.03	0.283	4.45	0.250	-	-	-	-
	IconShop [57]	2.6k	107.93	0.233	4.46	0.224	-	-	-	-
	LIVE [34]	52.2k	-	-	-	-	0.935	0.950	0.111	0.008
MMSVG-Illustration	DiffVG [29]	322.0k	-	-	-	-	0.945	0.955	0.065	0.001
	GPT-4o [21]	0.4k	-	-	-	-	0.875	0.854	0.373	0.077
	StarVector(8B) [42]	2.6k	-	-	-	-	0.877	0.900	0.238	0.046
	Vtracer	57.6k	_	_	-	-	0.994	0.966	0.035	0.002
	OmniSVG(4B)	5.8k	154.37	0.226	<u>4.56</u>	0.232	0.899	0.906	0.237	0.034
	OmniSVG-L(8B)	6.9k	138.42	0.231	4.51	0.232	0.905	0.907	0.231	0.031

## 5 Experiments

To validate the effectiveness of our method, we first introduce the baselines (Sec. 5.1). Then, we make quantitative comparisons with prior arts (Secs. 5.2 and 5.3) and conduct ablations (Sec. 5.4) to study the effectiveness of our design.

#### 5.1 Baselines

For the text-to-SVG task, we compare our method with language-based (LLM-based) methods, including VectorFusion [22], SVGDreamer [60], Chat2SVG [56] and IconShop [57]. For image-to-SVG task, we compare our method with baseline methods across image vectorization and Multimodal Large Language Modeling approaches, including LIVE [34], DiffVG [29], StarVector [42], Vtracer [12] and GPT-40 [21] using the official implementations with the hyperparameters proposed by the authors, and apply their pre- and post-processing code as required.

#### **5.2** Quantitative Comparisons

We compare our OmniSVG with other baseline methods on the "text-to-SVG" and "image-to-SVG" tasks in our MMSVG-Bench. In addition to the metrics mentioned in Sec. 3, we also report the average token length (# tokens) of a generated SVG sample utilizing the Qwen2.5-VL [1] tokenizer.

As shown in Tab. 2, OmniSVG demonstrates strong performance compared to state-of-the-art baselines in text-to-SVG generation, achieving superior FID scores and competitive CLIP score, aesthetic quality, and HPS. For image-to-SVG, OmniSVG also achieves competitive results with traditional vectorization methods, *i.e.* LIVE [34], DiffVG [29], and VTracer [12], but with a much shorter sequence length. When comparing to auto-regressive methods, *i.e.* GPT-4o [21] and StarVector [42], OmniSVG also achieves a superior performance across all metrics. The above results indicate that OmniSVG effectively balances the generation cost and the visual quality when generating SVGs according to multi-modal conditions.

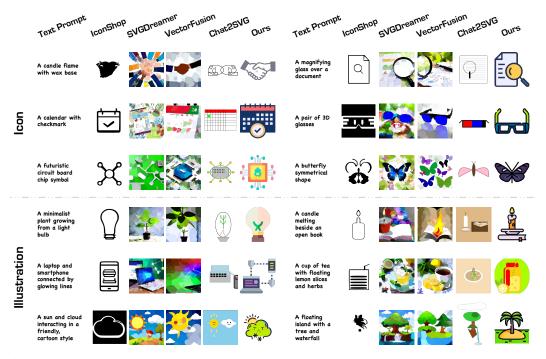


Figure 3: Qualitative Comparison with SOTA Methods on Text-to-SVG Task. We compare the propose method with SOTA Text-to-SVG methods on our evaluation benchmarks, namely Icon and Illustration.

#### 5.3 Qualitative Evaluations

**Text-to-SVG task.** We compare our method with baseline approaches using seven distinct text prompts for the text-to-SVG task, as shown in Fig. 4. Optimization-based methods like SVG-Dreamer [60] and VectorFusion [22] require significant computation time due to their iterative optimization processes, which, while effective for refining SVG details, are computationally expensive. Auto-regressive methods, such as IconShop [57] and Chat2SVG [56], generate SVGs more quickly by leveraging pre-trained models but have notable limitations. IconShop produces monochrome SVGs, restricting its applicability, while Chat2SVG, though flexible, generates less detailed and semantically consistent SVGs in its first stage. Our OmniSVG consistently outperforms all baselines across various text prompts in generating high-fidelity SVGs with rich color, geometric accuracy, and the ability to handle complex visual cues.

Image-to-SVG Task. We compare our method with classical image vectorization approaches, including DiffVG [29], LIVE [34], and VLMbased methods GPT-4o [21], StarVector [42] and Vtracer [12] As shown in Fig. 4, our method outperforms these baselines in both quality and efficiency. Optimization-based methods like DiffVG and LIVE perform well on simple icons but struggle with complex images, often generating visual artifacts. The GPT-40 model, while capable of generating SVGs for complex images, is limited to icon-level outputs and cannot handle detailed illustrations. StarVector excels at simple icons but fails to produce accurate SVGs for more intricate images, highlighting its limited generalization capability. Vtracer is an image processing algorithm designed to convert raster images into SVGs. In contrast, OmniSVG effi-

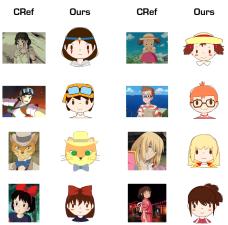


Figure 5: Generated SVG with Character-Reference (CRef) by OmniSVG.

ciently converts a wide range of images, from icons to complex illustrations and character images,

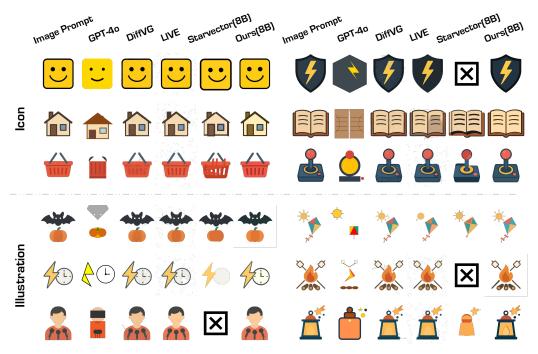


Figure 4: Qualitative Comparison with SOTA Methods on Image-to-SVG Task. We compare the propose method with SOTA Image-to-SVG methods on our evaluation benchmarks.

into high-quality, editable SVGs. This superior performance in handling diverse visual cues distinguishes OmniSVG from traditional vectorization methods. Additional visual results can be found in Fig. 12. We provide more detailed discussions with existing methods, particularly the recent works LLM4SVG [59] and StarVector [42], in the Appendix D.

**Character-Reference SVG generation task.** As shown in Fig. 5, by training on MMSVG-Character with natural character image and SVG pair data, OmniSVG is capable of generating character SVGs through image references.

#### 5.4 Ablation studies

**Effectiveness of SVG Parameterization.** We present a comprehensive comparison among different SVG parameterization strategy with the traditional non-parameterized methods for SVG representation in large language models. We ablates on the parameterization on both coordinate and color attributes of the SVG.

The results, shown in Tab. 3 and Fig. 6 demonstrate that parameterizing both coordinate and color attributes yields a better generation results under all metrics with the shortest token length. It further validates that the efficient token representation allows our method to generate complex SVGs with fewer computational resources. Additionally, qualitative results show that our method outperforms others, particularly as SVG complexity increases. The non-parameterization method fails to generate

Table 3: **Quantitative Study on SVG Parameterization.** Ablation studies on color parametrization (abbreviated as color param.) and coordinate parametrization (abbreviated as coord param.) are conducted.

Methods		Text-	to-SVG		# Tokens					
Methous	FID↓	CLIP↑	$Aesthetic \!\!\uparrow$	$HPS \!\!\uparrow$	$\text{DINO}{\uparrow}$	SSIM↑	LPIPS↓	$MSE \downarrow$	# TOKEHS	
w/o param.	218.76	0.185	3.43	0.138	0.741	0.718	0.315	0.182	18.5k	
w/o coordinate param.	193.42	0.216	3.90	0.169	0.826	0.809	0.248	0.119	10.2k	
w/o color param.	<u>167.28</u>	0.269	4.31	0.211	0.895	0.879	0.179	0.053	6.3k	
OmniSVG(4B)	145.89	0.308	4.59	0.238	0.946	0.928	0.138	0.020	4.8k	

Table 4: Ablation of the Model Size. As the model size grows, the generated samples are of higher quality.

Methods	Input	Size		Text-	to-SVG	Image-to-SVG				
Methous	Input	Size	$FID\downarrow$	CLIP↑	Aesthetic↑	HPS↑	DINO↑	SSIM↑	LPIS↓	$MSE\downarrow$
FLAN-T5-Base[10]	Text	223M	198.48	0.158	3.38	0.085	_	_	_	_
FLAN-T5-Large[10]	Text	770M	175.24	0.208	3.92	0.142	_	_	_	_
FLAN-T5-x1[10]	Text	3B	160.28	0.258	4.31	0.192	_	_	_	_
blip2-flan-t5-xl[28]	Text/Image	3.94B	152.11	0.235	4.48	0.215	0.898	0.891	0.255	0.041
OmniSVG(4B)	Text/Image	3.7B	145.89	0.308	4.59	0.238	0.946	0.928	0.138	0.020

SVGs for complex images. These findings underscore the effectiveness of our full parameterization strategy in balancing performance and resource efficiency for SVG generation tasks.

**Ablation studies on model size.** To analyze whether training a larger model benefits SVG generation, we evaluate OmniSVG base models with different sizes on the MMSVG-2M dataset in Tab. 4. We evaluate OmniSVG with base models of varying sizes on the MMSVG-2M dataset in Tab. 4 by progressively scaling up the model size. The results show that as the model size grows, we can generate SVG samples with a better quality.

Table 5: Ablation on VLM architecture.

Vision Model	Language Madel		Text-	to-SVG		Image-to-SVG				
Vision Model	Language Model	$FID\downarrow$	CLIP↑	Aesthetic <sup>↑</sup>	HPS↑	DINO↑	SSIM↑	LPIPS↓	$MSE\downarrow$	
CLIP	Qwen2.5	185.31	0.249	4.52	0.215	0.867	0.856	0.267	0.058	
VQGAN	Qwen2.5	198.74	0.234	4.49	0.203	0.839	0.828	0.295	0.071	
Qwen2.5-V	/L-3B-Instruct	145.89	0.308	4.59	0.238	0.946	0.928	0.138	0.020	
Qwen2.5-V	L-7B-Instruct	134.45	0.254	<u>4.56</u>	0.237	<u>0.914</u>	0.900	0.233	0.036	

**Ablation Studies on the VLM Architecture.** To evaluate the effectiveness of the VLM architecture, we conducted an ablation study replacing it with alternative LLM-based architectures incorporating image encoders such as CLIP ViT-B/32 [39], VQGAN [14], and Qwen2.5-VL [1].

The results in Tab. 5 show that Qwen2.5-VL consistently outperformed all alternatives under all evaluation metrics.

User Study. We extract one-tenth of the samples from the evaluation dataset and conducted a user study with 15 participants to evaluate user preferences, vividness, and the alignment between text-to-SVG and image-to-SVG. Participants are asked to assess SVGs generated by different models based on 150 text descriptions and 150 image prompts, comparing the results generated using our method and baseline models. The results in Fig. 7 show that OmniSVG is widely preferred, with higher scores for vividness and superior semantic alignment with the input conditions.

#### 6 Conclusions

**Conclusions.** We introduce OmniSVG, a unified framework for multimodal SVG generation that leverages pre-trained Vision-Language Models (VLMs). By parameterizing SVG com-



Figure 6: Qualitative Study on Parametrization.

mands and coordinates as discrete tokens, OmniSVG efficiently decouples structural logic from geometry, addressing issues like "coordinate hallucination" while maintaining design expressiveness. Our method outperforms existing approaches in both quality and efficiency, offering high-quality, editable SVG across various design domains. Additionally, we proposed MMSVG-2M, a large-scale multimodal dataset with two million annotated SVG assets and a standardized evaluation protocol. Extensive experiments show that OmniSVG surpasses prior SVG generation methods in various

conditional generation tasks, highlighting its potential for integration into professional SVG design workflows.

Limitations and Future Work. During inference, OmniSVG generates tens of thousands of tokens for complex samples, which inevitably leads to a considerable generation time. OmniSVG is only bounded by vector style image prompt and fails on natural images. As for future work, recent endeavors on multi-token prediction [15, 2] and KV-cache compression [68, 3] provide a promising way to save the generation cost. Additionally, the autoregressive nature of OmniSVG also unlocks future

Figure 7: User Study of OmniSVG and baselines.

Method	Preference ↑	<b>Vividity</b> ↑	Alignment↑
Vectorfusion [22]	35	58	76
SVGDreamer [60]	41	65	79
Chat2SVG [56]	55	61	86
IconShop [57]	79	57	75
GPT-4o [21]	38	54	80
StarVector(8B) [42]	37	81	68
DiffVG [29]	88	76	96
LIVE [34]	86	70	95
OmniSVG	96	88	98

opportunities for in-context learning [67, 69, 47], chain-of-thought reasoning [55, 16], and multi-turn interleaved generation [20, 31], thereby providing a more precise user control.

## Acknowledgements

This work is in part supported by National Key R&D Program of China (Grant No. 2022ZD0160103), National Natural Science Foundation of China (Grant No. 62276067), and National Natural Science Foundation of China (Grant No. 62472104).

The computations in this research were performed using the CFFF platform of Fudan University.

#### References

- [1] 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.
- [2] Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, and Tri Dao. Medusa: Simple framework for accelerating llm generation with multiple decoding heads. *Retrieved December*, 2023.
- [3] Zefan Cai, Yichi Zhang, Bofei Gao, Yuliang Liu, Tianyu Liu, Keming Lu, Wayne Xiong, Yue Dong, Baobao Chang, Junjie Hu, et al. Pyramidkv: Dynamic kv cache compression based on pyramidal information funneling. arXiv preprint arXiv:2406.02069, 2024.
- [4] Alexandre Carlier, Martin Danelljan, Alexandre Alahi, and Radu Timofte. Deepsvg: A hierarchical generative network for vector graphics animation. *NeurIPS*, 2020.
- [5] Sijin Chen, Xin Chen, Anqi Pang, Xianfang Zeng, Wei Cheng, Yijun Fu, Fukun Yin, Billzb Wang, Jingyi Yu, Gang Yu, et al. Meshxl: Neural coordinate field for generative 3d foundation models. *NeurIPS*, 2024.
- [6] Sijin Chen, Xin Chen, Chi Zhang, Mingsheng Li, Gang Yu, Hao Fei, Hongyuan Zhu, Jiayuan Fan, and Tao Chen. Ll3da: Visual interactive instruction tuning for omni-3d understanding reasoning and planning. In CVPR, 2024.
- [7] Zehao Chen and Rong Pan. Svgbuilder: Component-based colored svg generation with text-guided autoregressive transformers. arXiv preprint arXiv:2412.10488, 2024.
- [8] Wei Cheng, Ruixiang Chen, Siming Fan, Wanqi Yin, Keyu Chen, Zhongang Cai, Jingbo Wang, Yang Gao, Zhengming Yu, Zhengyu Lin, et al. Dna-rendering: A diverse neural actor repository for high-fidelity human-centric rendering. In *ICCV*, 2023.
- [9] Wei Cheng, Su Xu, Jingtan Piao, Chen Qian, Wayne Wu, Kwan-Yee Lin, and Hongsheng Li. Generalizable neural performer: Learning robust radiance fields for human novel view synthesis. *arXiv preprint arXiv:2204.11798*, 2022.
- [10] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *JMLR*, 2024.
- [11] Louis Clouâtre and Marc Demers. Figr: Few-shot image generation with reptile. *arXiv preprint* arXiv:1901.02199, 2019.
- [12] Vision Cortex. Vtracer. https://www.visioncortex.org/vtracer-docs, 2023.
- [13] Nyanko Devs. Danbooru2023: A large-scale crowdsourced and tagged anime illustration dataset. Hugging Face, 2023.
- [14] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In CVPR, 2021.
- [15] Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière, David Lopez-Paz, and Gabriel Synnaeve. Better & faster large language models via multi-token prediction. *arXiv preprint arXiv:2404.19737*, 2024.
- [16] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. arXiv preprint arXiv:2501.12948, 2025.
- [17] Han Guo, Songlin Yang, Tarushii Goel, Eric P Xing, Tri Dao, and Yoon Kim. Log-linear attention. *arXiv* preprint arXiv:2506.04761, 2025.
- [18] David Ha and Douglas Eck. A neural representation of sketch drawings. In ICLR, 2018.
- [19] Teng Hu, Ran Yi, Baihong Qian, Jiangning Zhang, Paul L Rosin, and Yu-Kun Lai. Supersvg: Superpixel-based scalable vector graphics synthesis. In *CVPR*, 2024.
- [20] Minbin Huang, Yanxin Long, Xinchi Deng, Ruihang Chu, Jiangfeng Xiong, Xiaodan Liang, Hong Cheng, Qinglin Lu, and Wei Liu. Dialoggen: Multi-modal interactive dialogue system for multi-turn text-to-image generation. arXiv preprint arXiv:2403.08857, 2024.
- [21] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. arXiv preprint arXiv:2410.21276, 2024.

- [22] Ajay Jain, Amber Xie, and Pieter Abbeel. Vectorfusion: Text-to-svg by abstracting pixel-based diffusion models. In CVPR, 2023.
- [23] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International conference on machine learning*, pages 5156–5165. PMLR, 2020.
- [24] Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia Li, Chenghao Mou, Carlos Muñoz Ferrandis, Yacine Jernite, Margaret Mitchell, Sean Hughes, Thomas Wolf, et al. The stack: 3 tb of permissively licensed source code. *arXiv preprint arXiv:2211.15533*, 2022.
- [25] Kozea. Cairosvg. https://cairosvg.org/, 2023.
- [26] Black Forest Labs. Flux. https://github.com/black-forest-labs/flux, 2024.
- [27] Black Forest Labs. Flux.1-redux-dev. https://huggingface.co/black-forest-labs/FLUX. 1-Redux-dev, 2024.
- [28] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, 2023.
- [29] Tzu-Mao Li, Michal Lukáč, Gharbi Michaël, and Jonathan Ragan-Kelley. Differentiable vector graphics rasterization for editing and learning. SIGGRAPH Asia, 2020.
- [30] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In NeurIPS, 2023.
- [31] Ziyu Liu, Tao Chu, Yuhang Zang, Xilin Wei, Xiaoyi Dong, Pan Zhang, Zijian Liang, Yuanjun Xiong, Yu Qiao, Dahua Lin, et al. Mmdu: A multi-turn multi-image dialog understanding benchmark and instruction-tuning dataset for lylms. arXiv preprint arXiv:2406.11833, 2024.
- [32] Raphael Gontijo Lopes, David Ha, Douglas Eck, and Jonathon Shlens. A learned representation for scalable vector graphics. In CVPR, 2019.
- [33] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017.
- [34] Xu Ma, Yuqian Zhou, Xingqian Xu, Bin Sun, Valerii Filev, Nikita Orlov, Yun Fu, and Humphrey Shi. Towards layer-wise image vectorization. In *CVPR*, 2022.
- [35] 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.
- [36] Dongwei Pan, Long Zhuo, Jingtan Piao, Huiwen Luo, Wei Cheng, Yuxin Wang, Siming Fan, Shengqi Liu, Lei Yang, Bo Dai, et al. Renderme-360: a large digital asset library and benchmarks towards high-fidelity head avatars. *NeurIPS*, 2023.
- [37] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In ICLR, 2023.
- [38] 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.
- [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.
- [40] Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations toward training trillion parameter models. In SC20: International Conference for High Performance Computing, Networking, Storage and Analysis. IEEE, 2020.
- [41] Pradyumna Reddy, Michael Gharbi, Michal Lukac, and Niloy J Mitra. Im2vec: Synthesizing vector graphics without vector supervision. In CVPR, 2021.
- [42] 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.

- [43] Christoph Schuhmann. Improved aesthetic predictor. https://github.com/christophschuhmann/improved-aesthetic-predictor, 2022.
- [44] I-Chao Shen and Bing-Yu Chen. Clipgen: A deep generative model for clipart vectorization and synthesis. *TVCG*, 2022.
- [45] Yiren Song, Xuning Shao, Kang Chen, Weidong Zhang, Zhongliang Jing, and Minzhe Li. Clipvg: Text-guided image manipulation using differentiable vector graphics. In AAAI, 2023.
- [46] Hao Su, Xuefeng Liu, Jianwei Niu, Jiahe Cui, Ji Wan, Xinghao Wu, and Nana Wang. Marvel: Raster gray-level manga vectorization via primitive-wise deep reinforcement learning. *TCSVT*, 2023.
- [47] Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiying Yu, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative multimodal models are in-context learners. In CVPR, 2024
- [48] Zecheng Tang, Chenfei Wu, Zekai Zhang, Mingheng Ni, Shengming Yin, Yu Liu, Zhengyuan Yang, Lijuan Wang, Zicheng Liu, Juntao Li, et al. Strokenuwa: Tokenizing strokes for vector graphic synthesis. arXiv preprint arXiv:2401.17093, 2024.
- [49] Zecheng Tang, Chenfei Wu, Zekai Zhang, Mingheng Ni, Shengming Yin, Yu Liu, Zhengyuan Yang, Lijuan Wang, Zicheng Liu, Juntao Li, et al. Strokenuwa: Tokenizing strokes for vector graphic synthesis. arXiv preprint arXiv:2401.17093, 2024.
- [50] Lucas Theis, Aäron van den Oord, and Matthias Bethge. A note on the evaluation of generative models. arXiv preprint arXiv:1511.01844, 2015.
- [51] Yingtao Tian and David Ha. Modern evolution strategies for creativity: Fitting concrete images and abstract concepts. In *Artificial Intelligence in Music, Sound, Art and Design*, 2022.
- [52] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. arXiv preprint arXiv:2409.12191, 2024.
- [53] Yizhi Wang and Zhouhui Lian. Deepvecfont: synthesizing high-quality vector fonts via dual-modality learning. TOG, 2021.
- [54] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. TIP, 2004.
- [55] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *NeurIPS*, 2022.
- [56] Ronghuan Wu, Wanchao Su, and Jing Liao. Chat2svg: Vector graphics generation with large language models and image diffusion models. *arXiv preprint arXiv:2411.16602*, 2024.
- [57] Ronghuan Wu, Wanchao Su, Kede Ma, and Jing Liao. Iconshop: Text-guided vector icon synthesis with autoregressive transformers. TOG, 2023.
- [58] Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score: Better aligning text-to-image models with human preference. In *ICCV*, 2023.
- [59] Ximing Xing, Juncheng Hu, Guotao Liang, Jing Zhang, Dong Xu, and Qian Yu. Empowering llms to understand and generate complex vector graphics. *arXiv* preprint arXiv:2412.11102, 2024.
- [60] Ximing Xing, Haitao Zhou, Chuang Wang, Jing Zhang, Dong Xu, and Qian Yu. Svgdreamer: Text guided svg generation with diffusion model. In CVPR, 2024.
- [61] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. arXiv preprint arXiv:2412.15115, 2024.
- [62] Yiying Yang, Fukun Yin, Wen Liu, Jiayuan Fan, Xin Chen, Gang Yu, and Tao Chen. Pm-inr: Prior-rich multi-modal implicit large-scale scene neural representation. In AAAI, 2024.
- [63] Fukun Yin, Xin Chen, Chi Zhang, Biao Jiang, Zibo Zhao, Wen Liu, Gang Yu, and Tao Chen. Shapegpt: 3d shape generation with a unified multi-modal language model. TMM, 2025.
- [64] Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*, 2024.

- [65] Zhengming Yu, Wei Cheng, Xian Liu, Wayne Wu, and Kwan-Yee Lin. Monohuman: Animatable human neural field from monocular video. In *CVPR*, 2023.
- [66] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.
- [67] Yuanhan Zhang, Kaiyang Zhou, and Ziwei Liu. What makes good examples for visual in-context learning? NeurIPS, 2023.
- [68] Xiabin Zhou, Wenbin Wang, Minyan Zeng, Jiaxian Guo, Xuebo Liu, Li Shen, Min Zhang, and Liang Ding. Dynamickv: Task-aware adaptive kv cache compression for long context llms. *arXiv preprint arXiv:2412.14838*, 2024.
- [69] Yucheng Zhou, Xiang Li, Qianning Wang, and Jianbing Shen. Visual in-context learning for large vision-language models. *arXiv* preprint arXiv:2402.11574, 2024.
- [70] Bocheng Zou, Mu Cai, Jianrui Zhang, and Yong Jae Lee. Vgbench: A comprehensive benchmark of vector graphics understanding and generation for large language models. In EMNLP, 2024.

## **NeurIPS Paper Checklist**

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

## IMPORTANT, please:

- Delete this instruction block, but keep the section heading "NeurIPS Paper Checklist",
- · Keep the checklist subsection headings, questions/answers and guidelines below.
- Do not modify the questions and only use the provided macros for your answers.

#### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The main claims made in the abstract and introduction accurately reflect the paper's contributions and scope.

### Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
  contributions made in the paper and important assumptions and limitations. A No or
  NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Limitations are discussed.

#### Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

#### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The paper does not include theoretical results.

#### Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

## 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: In the experiment section, we give detailed information about the experimental setup, evaluated models and evaluation metrics.

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We provide access to the data and evaluation code for reproduction.

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
  to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).

 Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Experimental settings are indicated in the experiment section.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental
  material.

#### 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: Error bars are not reported.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Experiments compute resources are indicated in the supplementart material.

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.

- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

#### 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: Our research conforms to the NeurIPS Code of Ethics.

#### Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

#### 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: Broader impacts is discussed in the abstract and introduction. We aim to provide valuable tools to the community for developing more powerful video understanding models.

#### Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

#### 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [Yes]
Justification: We

#### Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

#### 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We have cited all papers. The license and copyright information related to data from existing datasets and benchmarks are discussed in the Supplementary Materials.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- · For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We introduce the details of the dataset, and more details are discussed in the Supplementary materials.

Guidelines: Details of the dataset and code, including the license and limitations, are discussed in the Supplementary Materials.

- The answer NA means that the paper does not release new assets.
- · Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- · At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

## 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [Yes]

Justification: The labeling process is discussed in Section 3. The annotation guidelines are provided in the Supplementary Materials.

#### Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: There are no potential risks.

#### Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

## 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

## **Appendix**

## A Additional Details of MMSVG-2M dataset

#### A.1 Samples of MMSVG-2M Dataset

We visualize samples of our MMSVG-2M dataset in Fig. 8. In our MMSVG-2M dataset, 55% of the SVG samples belongs to the MMSVG-Icon, 25% belongs to the MMSVG-Illustration, and the rest 20% belongs to the MMSVG-Character. Among the SVG samples within the MMSVG-Character category, half of them comes from Freepik, while another half is generated by our data creation pipeline. We also collect image-SVG pairs for the character-reference SVG generation tasks during the generation process.

Table 6: **Data Statistics for MMSVG-2M.** Our MMSVG-2M consists of 1.1 million SVG icons, 0.5 million SVG illustrations, and 0.4 million SVG anime characters.

Dataset	Train	Val	Total	Source	Token Length
MMSVG-Icon	990k	110k	1,100k	Iconfont	$2.2k \pm 0.9k$
MMSVG-Illustration MMSVG-Character	450k 350k	50k 50k	500k 400k	IconScout Freepik & generated	$8.1 ext{k} \pm 3.3 ext{k}  28 ext{k} \pm 7.3 ext{k}$

#### A.2 SVG-Image-Text Pairs Construction

Our *MMSVG-2M* dataset comprises two million SVG samples with the corresponding rasterized images. We generate captions on the rasterized images with BLIP-2 [28], thereby providing textual descriptions that enable us to fine-tune our model to follow these instructions. We use CairoSVG [25] for rasterization and remove samples that produced completely white images.

**Annotation.** We employ an off-the-shelf VLM, specifically BLIP-2 [28], to generate SVG captions with the prompt below. To reduce hallucinations, we drop the samples with CLIP scores less than 30. We also visualize the distribution annotated keywords of MMSVG-2M dataset in Fig. 10 with word cloud format. And the instruction template for annotation is shown in Tab. 7.

**Instruction templates.** MMSVGBench provides three tasks, including text-to-SVG task, image-to-SVG task and character-reference SVG generation task. Each task needs different instruction templates. For the text and image conditioning SVG generation, we provide the input text or image with VLM architecture. For character-reference SVG generation, we provide the natural character

#### **Instructions for Different Tasks**

- **Employed BLIP2 for SVG Captioning:** You are a helpful assistant. Your task is to describe this image in a single sentence, including the object, its color, and its overall arrangement. For example: "Yellow cheers with glasses of alcohol drinks." / "Heart emojis represent love on Valentine's Day."
- **Text-to-SVG:** You are a helpful SVG Generation assistant, designed to generate SVG. We provide the text description as input, generate SVG based on the text.
- **Image-to-SVG:** You are a helpful SVG Generation assistant, designed to generate SVG. We provide an image as input, generate SVG for this image.
- **Character-Reference SVG Generation:** You are a helpful SVG Generation assistant, designed to generate SVG. We provide a natural image as input, please generate the simplified character SVG based on the reference input image.

Table 7: Instructions for Different Tasks. Instructions including annotation, text-to-SVG, image-to-SVG and character-reference SVG generation.

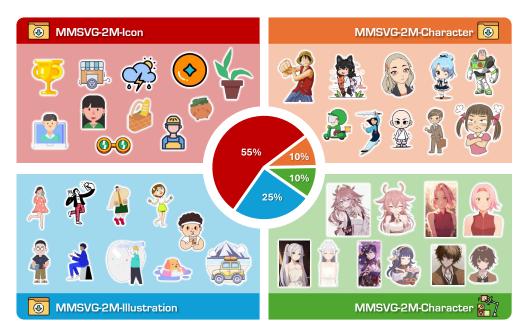


Figure 8: Samples from MMSVG-2M Dataset. The proposed MMSVG-2M dataset can be separated into three subset, namely Icon, Illustration and Character. Samples from Icon, Illustration and part of Character subsets are downloaded from Internet. Another part of Character subset is generated by our data creation pipeline, which can provide image and SVG pairs for image prompting task.

reference image and the original image with the VLM architecture. The list of instruction templates for different tasks are shown in Tab. 7.

#### A.3 Character-SVG Pairs Construction

As illustrated in the Fig. 6, part of our proposed MMSVG-2M-Character subset is constructed using a generative pipeline. As shown in the pipeline diagram in Fig. 2, we employ a FLUX [26]-based generative model enhanced with a vector-style LoRA to enable the generation of SVG-style data. For image-based conditioning, we adopt FLUX-Redux [27], which injects image features via a SigLIP encoder and projects them into image embeddings. These embeddings are then concatenated with the text tokens as conditioning inputs for FLUX [26]. However, in practice, the original Redux [27] conditioning proves to be overly strong. To address this, we adopt a community-implemented variant of Redux that downsamples the image embeddings in 2D space. As observed in our experiments shown in Fig. 9, a downsampling factor between  $2 \times$  and  $3 \times$  yields the most reasonable SVG-style character references. Finally, we employ VTracer [12] to perform near-instant vectorization of the generated images. To construct the MMSVG-2M-Character subset, we first filter 103k character instances from the Danbooru [13] dataset and apply the aforementioned pipeline with motion and expression keywords like previous works [8, 9, 36, 65]. We compare the raw FLUX [26] outputs and their vectorized counterparts, retaining only those samples with PSNR and SSIM scores above a certain threshold as valid data.

## **B** Additional Details

#### B.1 Scaling Up

To study the effectiveness of scaling up multimodal SVG generation, we scale up OmniSVG from 4B to 8B parameters. We present training perplexity in Fig. 11, where both models are trained from scratch on 250 billion tokens. We show that, as the size of the model grows, the model achieves a lower validation perplexity, indicating a higher probability of producing the validation data.

## **B.2** Implementation Details

We train our models in bfloat16 with the ZeRO-2 strategy [40] for memory-efficient training. We also adopt the AdamW [33] optimizer with a learning rate decaying from  $3\times10^{-4}$  to  $3\times10^{-6}$  and a weight decay of 0.1 to train our model. In practice, we load the pre-trained weights from the Qwen2.5-VL [1] model and initialize the SVG embeddings from scratch. Without further specification, we generate SVGs with the top-k and top-p sampling strategy with k=50 and p=0.95 for diversity.

## C Additional Results

As list in full comparisons in Tab. 2, including all the baselines mentioned in Sec. 5. For the text-to-SVG task, we compare our method with language-based (LLM-based) methods, including VectorFusion [22], SVGDreamer [60], Chat2SVG [56] and IconShop [57]. For image-to-SVG task, we compare our method with baseline methods across image vectorization and Multimodal Large Language Modeling ap-



Figure 10: Word Cloud Visualization of Label Distribution in the MMSVG-2M Dataset. The size of each label corresponds to its frequency of occurrence. The larger the label, the more frequently it appears in the dataset.

proaches, including LIVE [34], DiffVG [29], StarVector [42] and GPT-40 [21] using the official implementations with the hyperparameters proposed by the authors, and apply their pre- and postprocessing code as required. Specifically, for the text-to-SVG task, the optimization-based method SVGDreamer excels in enhancing editability by employing a semantic-driven image vectorization process that effectively separates foreground objects from the background, while failing to handle complex scenes. Another optimization-based work, VectorFusion, stands out for generating SVGexportable vector graphics without relying on large captioned datasets. However, Vectorfusion is also unable to handle complex scenarios and diverse styles. The significant problem with these optimization-based works is that the optimization time is too long. Generating an SVG usually takes more than ten minutes, which is too expensive. For the LLM-based method, Chat2SVG integrates Large Language Models (LLMs) with image diffusion models to create semantically rich SVG templates. However, Chat2SVG still needs to optimize the output SVG script from LLM, which introduces increased computational complexity and poses challenges during model training. In comparison, IconShop utilizes a transformer-based architecture to autoregressively model SVG path sequences, demonstrating exceptional performance in simplified icon SVGs, which offers effective solutions for text-to-SVG generation. It can only generate black simple Icon SVGs.

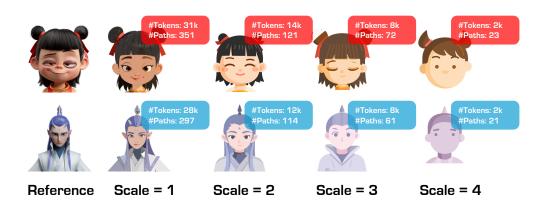


Figure 9: Image Prompting Dataset Creation of MMSVG-2M Character. By utilizing FLUX-Redux and SVG vectorization tools, image prompting data pairs can be generated. We adpot FLUX-Redux downsampling scale with 2, 3 in practice by trading-off the character similarity and complexity of generated SVG.

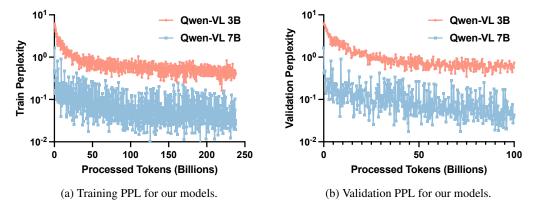


Figure 11: **Training and Validation Perplexity (PPL) for OmniSVG Models.** We train all the models from scratch on 250 billion tokens. We observe that the performance grows with model sizes.

For the image-to-SVG task, we compare our method with the image vectorization methods. LIVE allows progressive and efficient generation of SVGs, optimizing closed vector paths under raster image supervision with shape complexity control. However, LIVE needs to optimize for a long time when generating complex SVGs. DiffVG enables end-to-end differentiability in vector graphics rasterization, improving optimization through anti-aliasing and gradient-based methods while also is computationally expensive due to the complexity of the forward-backward rasterization process. Recently, the Multimodal Large Language Model (MLLM) based method StarVector leverages the visual understanding to apply accurate SVG primitive to the LLM architecture, which also can generate SVGs from both text and image inputs. However, it still fails to generate complex SVGs. Since Starvector [42] has not yet opened up its text-to-SVG model weights, our MMSVGBench does not evaluate Starvector's text-to-SVG capabilities. MMSVG-Bench also evaluates our methods with VLM methods, GPT-40, to conduct a comprehensive assessment. We compare our method with these baselines on our MMSVG-2M dataset, from simple MMSVG-Icon datset, a bit complex MMSVG-illustration dataset, to the very complex MMSVG-Character dataset.

#### D More details of the baselines

## D.1 Text-to-SVG Task

**SVGDreamer** [60] uses a semantic-driven image vectorization (SIVE) process to separate foreground objects and background, improving editability. The SIVE process utilizes attention-based primitive control and an attention-mask loss function to manipulate individual elements effectively. To address issues in existing text-to-SVG generation methods, the proposed Vectorized Particle-based Score Distillation (VPSD) approach models SVGs as distributions of control points and colors, improving shape, color diversity, and convergence speed.

**VectorFusion** [22] leverages a text-conditioned diffusion model trained on pixel representations to generate SVG exportable vector graphics without needing large captioned SVG datasets. By optimizing a differentiable vector graphics rasterizer, it distills semantic knowledge from a pretrained diffusion model and uses Score Distillation Sampling to generate an SVG consistent with a caption. Experiments show that VectorFusion improves both quality and fidelity, offering a variety of styles such as pixel art and sketches.

**Chat2SVG** [56] proposes a hybrid framework that combines the strengths of Large Language Models (LLMs) and image diffusion models for text-to-SVG generation. The approach first uses an LLM to create semantically meaningful SVG templates from basic geometric primitives. A dual-stage optimization pipeline, guided by image diffusion models, refines paths in latent space and adjusts point coordinates to enhance geometric complexity.

**IconShop** [57] uses a transformer-based architecture to encode path commands and learn to model SVG path sequences autoregressively. It has shown excellent results in simplified icon scenarios and provides a good solution to Text-to-SVG generation by extending the FIGR-8-SVG dataset with

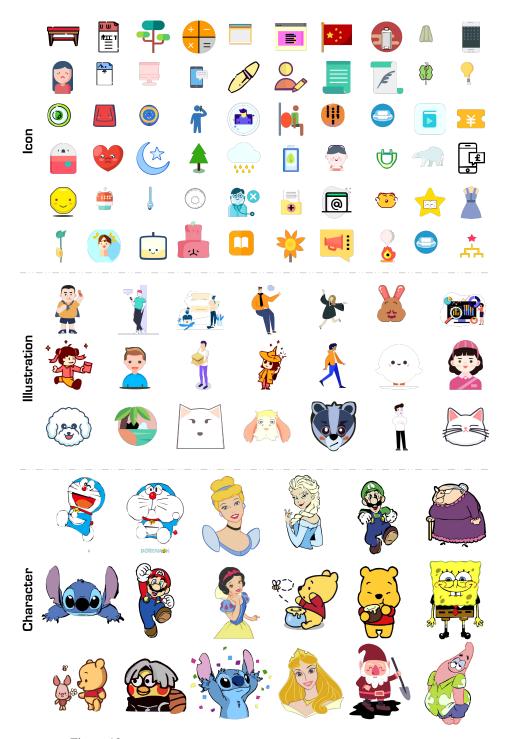


Figure 12: Illustration of the SVG Generation Capabilities of OmniSVG.

captions. We have access to their dataset and original splits and have trained our model on that data using a pre-trained checkpoint (trained on OmniVG dataset). We have extracted the results from IconShop and included them here to compare our method.

**LLM4SVG** [59] is a framework that leverages Large Language Models (LLMs) to understand and generate Scalable Vector Graphics (SVGs). It employs a structured SVG encoding approach, utilizing learnable semantic tokens to accurately represent SVG components and their properties. This design enables LLMs to produce SVGs that are both semantically aligned with textual descriptions and visually coherent. However, LLM4SVG also has a maximum token length of 2048, limiting its ability to generate highly complex SVGs that require longer sequences.

## D.2 Image-to-SVG Task

**LIVE** (Layer-wise Image Vectorization) [34] is a method for progressively generating SVGs that closely fit a given raster image by recursively adding and optimizing closed vector paths. Using a differentiable renderer (based on DiffVG [29]), LIVE enables direct optimization of paths under raster image supervision while controlling shape complexity by adjusting the number of path segments. It introduces component-wise path initialization, identifying key visual components to ensure efficient topology extraction and minimize redundant shapes.

**DiffVG** [29] is a landmark in vector graphics research, pioneering deep learning-based methods with the first differentiable vector graphics rasterization pipeline. By leveraging a combination of antialiasing techniques and gradient-based optimization, DiffVG ensures differentiability. Unlike methods relying on non-differentiable curve-to-mesh conversions, DiffVG employs a forward-backward rasterization process, where the forward pass generates antialiased images and the backward pass computes gradients with respect to vector graphic parameters.

**StarVector** [42] works directly in the SVG code space, leveraging visual understanding to apply accurate SVG primitives. StarVector employs a transformer-based architecture that integrates an image encoder with a language model, enabling it to process visual inputs and produce precise SVG code. StarVector effectively handles diverse SVG types, including icons, logos, and complex diagrams, demonstrating robust generalization across various vectorization tasks. However, with a 16k token context window, StarVector may struggle to process highly complex SVGs that require longer sequences.

**Vtracer** [12] is an image processing algorithm designed to convert raster images into SVGs. The algorithm follows a three-step pipeline, which involves the hierarchical clustering of images for vectorization. Initially, the pixels are transformed into paths, which are subsequently simplified into polygons. In the final step, these polygons are smoothed and approximated using a Bezier curve fitter.