

SVGFusion: Scalable Text-to-SVG Generation via Vector Space Diffusion

Ximing Xing¹, Juncheng Hu¹, Jing Zhang¹, Dong Xu², Qian Yu^{1*}

¹Beihang University, China

{ximingxing, hujuncheng, zhang-jing, qianyu}@buaa.edu.cn

²The University of Hong Kong, China

dongxu@cs.hku.hk

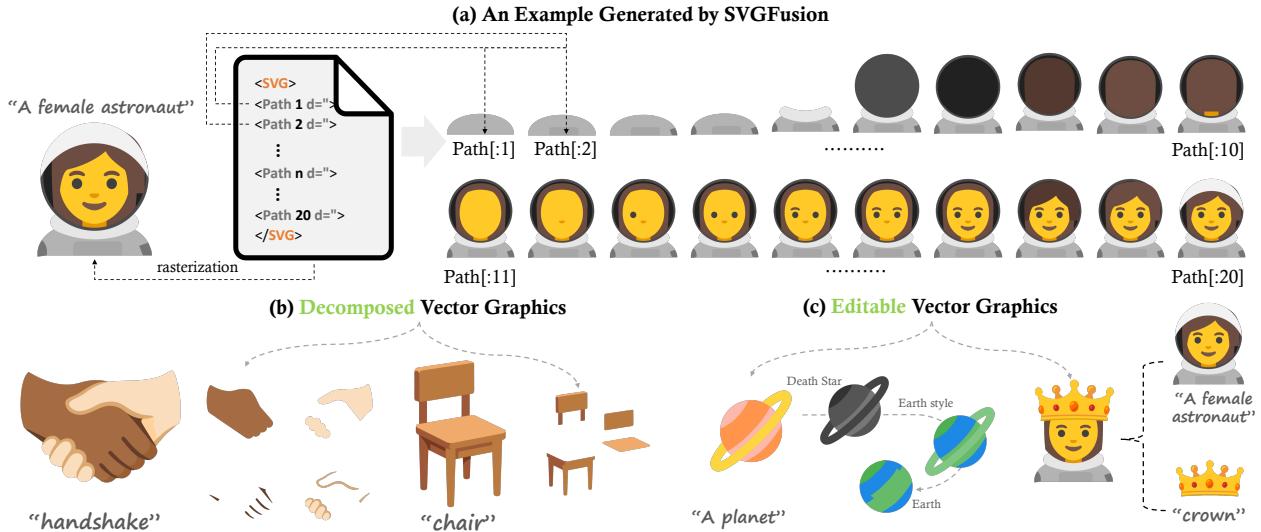


Figure 1. **Example SVGs generated by our SVGFusion.** Our proposed method, SVGFusion can generate SVGs with (a) reasonable construction, (b) a clear and systematic layering structure, and (c) highly editability.

Abstract

In this work, we introduce SVGFusion, a Text-to-SVG model capable of scaling to real-world SVG data without relying on text-based discrete language models or prolonged Score Distillation Sampling (SDS) optimization. The core idea of SVGFusion is to utilize a popular Text-to-Image framework to learn a continuous latent space for vector graphics. Specifically, SVGFusion comprises two key modules: a Vector-Pixel Fusion Variational Autoencoder (VP-VAE) and a Vector Space Diffusion Transformer (VS-DiT). The VP-VAE processes both SVG codes and their corresponding rasterizations to learn a continuous latent space, while the VS-DiT generates latent codes within this space based on the input text prompt. Building on the VP-VAE, we propose a novel rendering sequence modeling strategy which enables the learned latent space to capture the inherent creation logic of SVGs. This allows the model to generate SVGs

with higher visual quality and more logical construction, while systematically avoiding occlusion in complex graphic compositions. Additionally, the scalability of SVGFusion can be continuously enhanced by adding more VS-DiT blocks. To effectively train and evaluate SVGFusion, we construct SVGX-Dataset, a large-scale, high-quality SVG dataset that addresses the scarcity of high-quality vector data. Extensive experiments demonstrate the superiority of SVGFusion over existing SVG generation methods, establishing a new framework for SVG content creation.

1. Introduction

Scalable Vector Graphics (SVGs) are indispensable in modern digital design due to their ability to represent shapes—such as paths, curves, and polygons—using mathematical definitions instead of pixel-based raster

graphics. This vector-based approach enables SVGs to scale seamlessly across different resolutions without any loss of detail, ensuring precise quality control. Additionally, SVGs are highly efficient in terms of storage and transmission due to their compact structure. They also offer exceptional editability, allowing designers to make fine adjustments to graphic elements. Consequently, SVGs play a critical role in applications such as web design, user interfaces, and the creation of icons, logos, and emojis.

The task of SVG generation has garnered increasing attention in recent years. Existing methods can be broadly categorized into two types: optimization-based and language-model-based. Optimization-based methods, such as those described in [9, 20, 35, 46, 52, 56, 65, 68], typically initialize an SVG in vector space and then rasterize it into pixel space using differentiable rasterizers [21] for loss computation. Gradients are back-propagated to the vector space to update the SVG parameters. While these methods can produce high-quality results by leveraging priors from large pre-trained models such as CLIP [35] or Stable Diffusion [38], they are computationally expensive and often support only a limited set of differentiable SVG commands, such as Cubic Bézier curves. Moreover, the generated results are often challenging to edit, as the primitives tend to intertwine or occlude one another.

On the other hand, language-model-based methods [4, 13, 49, 58, 59, 63, 66] treat SVG generation as a sequential task, using language models like RNNs or Transformers to predict commands in an auto-regressive manner. However, these methods struggle with efficiency when the sequence of an SVG becomes excessively long and face challenges in learning the complex syntax and diverse primitive types of SVGs without extensive training data. Consequently, these methods often restrict their output primarily to black-and-white icons and single-letter shapes, significantly limiting their applicability in real-world SVG applications. Moreover, the sequential nature of these models can lead to error accumulation and limited global coherence in the generated outputs.

In this work, we propose SVGFusion, a novel model designed for text-based SVG generation. SVGFusion is capable of producing high-quality SVGs that are both decomposable and editable, as illustrated in Fig. 1. Inspired by the recent success of Latent Diffusion Models (LDMs) [38] in the Text-to-Image (T2I) task, SVGFusion adopts an LDM-like framework to learn a continuous latent space for SVG generation. It comprises two main components: 1) *Vector-Pixel Fusion Variational Autoencoder (VP-VAE)*: A transformer-based VAE that learns a unified latent space for vector graphics by combining both SVG code and its rasterized image counter-

part. 2) *Vector Space Diffusion Transformer (VS-DiT)*: A scalable diffusion backbone operating in the learned latent space, conditioned on text prompts to generate highly-editable SVGs.

Our model introduces three key designs. (1) **Vector-Pixel Fusion Encoding**. To learn a more representative latent space, features from both SVG code and its rasterization are used, leading to SVGs with higher visual quality. During the learning phase for vector graphics, the VP-VAE encoder processes both the SVG code and its rasterization as inputs, while the decoder reconstructs the SVG code. Before being fed into the encoder, the SVG code is converted into a matrix representation with geometric attributes, colors, and opacity. Additionally, a pre-trained DINOv2 [31] is employed to extract visual features from the SVG rendering. (2) **Rendering Sequence Modeling**. Existing methods often struggle to generate editable SVGs because they lack an understanding of how SVG codes correspond to their visual representations. To address this challenge, we introduce a rendering sequence modeling strategy that presents the model with an incremental accumulation of SVG codes and their corresponding renderings. This strategy enables the model to learn the design logic of human creators, leading to SVGs with more logical constructions and enhanced visual quality, as shown in Fig. 1(a). (3) **Incorporation of DiT Blocks**. Motivated by the success of DiT [32] in Text-to-Image (T2I) tasks and their scalability, we incorporate VS-DiT blocks and perform the diffusion process in the latent space specifically learned for SVGs. This design allows SVGFusion’s generative capabilities to be continuously enhanced by scaling the VS-DiT through the addition of more blocks.

To train and evaluate our model, we curated a dataset of approximately 240,000 high-quality, human-designed SVGs from different sources. Additionally, we designed an automated data pre-processing pipeline to trim the SVGs losslessly. The contributions of this work are threefold:

- **A Novel and Scalable SVG Generation Model:** We introduce SVGFusion, which adapts a popular Text-to-Image framework for Text-to-SVG generation. This model is scalable and capable of producing high-quality SVGs that are highly editable.
- **Specific Module Design and Training Strategy for SVGs:** We propose a Vector-Pixel Fusion Variational Autoencoder (VP-VAE) that learns the latent space for SVGs from both SVG codes and their corresponding rasterizations. A novel rendering sequence modeling strategy is also designed to enable the model to understand the creation logic of SVGs, thereby enhancing the visual quality and editability of the generated outputs.
- **A Large-Scale, High-Quality SVG Dataset and**

Comprehensive Evaluation: We construct *SVGX-Dataset*, a collection of approximately 240,000 high-quality, human-designed SVGs, and conduct extensive experiments to validate SVGFusion’s effectiveness in SVG generation, establishing a new benchmark for this task.

2. Related Work

2.1. Vector Graphics Generation

Scalable Vector Graphics (SVGs) are widely used in design due to their key advantages, including geometric manipulability, resolution independence, and compact file sizes. One approach to SVG generation involves training neural networks to output predefined SVG commands and attributes [4, 13, 23, 37, 49, 63]. These networks typically use architectures like RNNs [13, 37], VAEs [23], and Transformers [4, 49, 63]. However, their ability to generate intricate and diverse vector graphics is constrained by the lack of large-scale datasets for training.

Unlike raster image generation, which benefits from large-scale datasets like ImageNet [7], the vector graphics domain has limited publicly available datasets, most of which focus on narrow domains such as monochrome icons [6], emojis [10], and fonts [58]. As an alternative to directly training an SVG generation network, optimization-based approaches have been explored, where vector parameters are iteratively refined to match a target image during inference.

Li *et al.* [21] introduced DiffVG, a differentiable rasterizer that bridges vector and raster domains, enabling gradient-based SVG optimization. Subsequent works [18, 26, 46, 52, 65, 67, 69] leveraged differentiable rasterizers with vision-language models, such as CLIP [35], to enhance text-guided sketch generation [9, 56]. Meanwhile, diffusion models like DreamFusion [34] demonstrated superior generative capabilities. Recent methods, including VectorFusion [20], DiffSketcher [65], and SVGDreamer [68], integrate differentiable rasterization with diffusion models, achieving promising results in sketching and iconography. However, they still suffer from limited editability, redundant geometry, and suboptimal graphical quality. To address these issues, hybrid approaches [51, 69] incorporate geometric constraints to refine SVG paths, but remain restricted to SDS-optimized structures and path-based primitives. More recently, VecFusion [50] introduced an image-conditioned diffusion model for vector font synthesis, though its scope is limited to fonts. In this work, we present SVGFusion, a novel model for scalable SVG generation in a continuous vector space. By moving beyond discrete language models and intensive optimization, SVGFusion enables high-quality, diverse,

and editable SVG generation while overcoming key limitations of prior methods.

2.2. Diffusion Model

Denoising diffusion probabilistic models (DDPM) [8, 15, 16, 28, 38, 42–45] have demonstrated outstanding performance in generating high-quality images. The diffusion model architecture combined with the language-image pretrained model [35] shows obvious advantages in text-to-image (T2I) tasks, including GLIDE [29], Stable Diffusion [38], DALL-E 2 [36], Imagen [39] and DeepFloyd IF [47], SDXL [33]. The progress achieved by T2I diffusion models [29, 36, 38, 39] also promotes the development of a series of text-guided tasks, such as text-to-3D [34, 57, 60] and text-to-video [12, 17, 22, 41].

Recent efforts such as DreamFusion [34] explores text-to-3D generation by exploiting a Score Distillation Sampling (SDS) loss derived from a 2D text-to-image diffusion model [38, 39] instead, showing impressive results. In addition, Sora [22] based on the latent diffusion model [32] has made amazing progress in the field of video generation. Recently, the architecture of diffusion models has been shifting from U-Net [8] architectures to transformer-based architectures [3, 25, 32], narrowing the gap between image generation and language understanding tasks. In this work, we extend the diffusion transformer to the domain of vector graphics, enabling the synthesis of vector graphics. We also demonstrate the potential of the proposed method in vector design. However, the absence of a scalable foundation model for vector graphics has significantly hindered the development of this field for broader applications. To address this, we propose SVGFusion, a scalable foundation model based on vector space design.

3. Methods

Our task is to generate editable SVGs from input text prompts. As illustrated in Fig. 2, our model first trains a Vector-Pixel Fusion Variational Autoencoder (VP-VAE) to learn a latent space \mathcal{Z} for SVGs. Next, a Vector Space Diffusion Transformer (VS-DiT) is trained within this latent space to generate new latent codes conditioned on text prompts. Once trained, given an input text and a randomly sampled latent code, our model produces an SVG that is semantically aligned with the text. In the following sections, we first describe the process of converting SVG code into SVG embeddings (Sec.3.1), followed by discussions on VP-VAE (Sec.3.2) and VS-DiT (Sec. 3.3).

3.1. SVG Representation

Vector graphics consist of machine instructions composed of a series of XML elements (*e.g.* `<path>` or

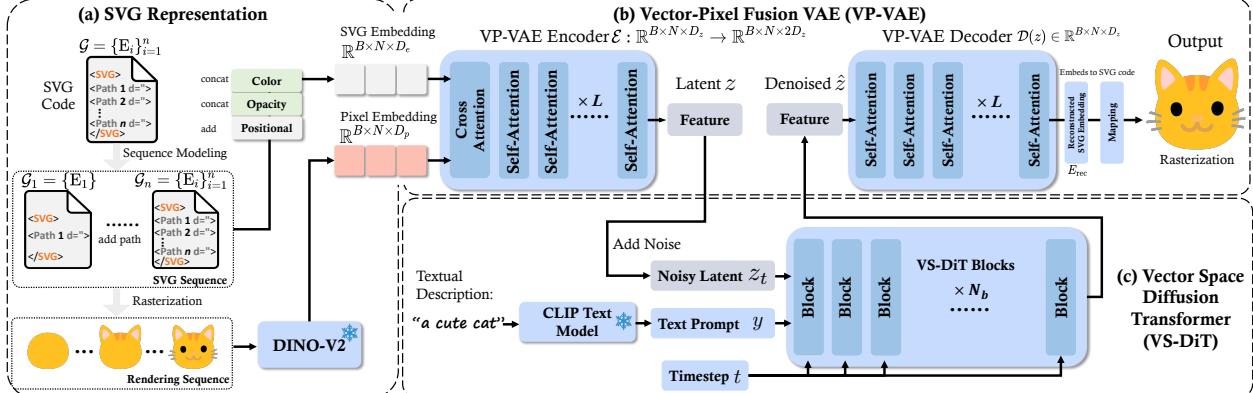


Figure 2. An Overview of SVGFusion. (a) The pipeline begins with the representation of SVGs, where XML-defined SVG code is converted into an SVG embedding (Sec. 3.1). (b) We first train a Vector-Pixel Fusion Variational Autoencoder (VP-VAE, Sec. 3.2) with a transformer-based architecture to learn a continuous latent space for SVGs by incorporating features from both SVG codes and their rendered images. (c) The Vector Space Diffusion Transformer (VS-DiT, Sec. 3.3) is then trained within the learned latent space to generate new latent codes conditioned on input text descriptions.

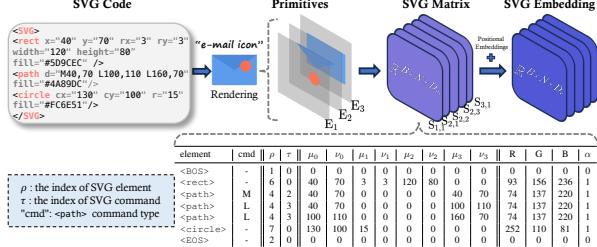


Figure 3. Illustration of the SVG embedding process. SVG code is initially converted into a matrix representation that includes geometric attributes, colors, and opacity. This matrix is subsequently mapped into a tensor via SVG embeddings.

<rect>), commands (e.g. M, C in <path> element), and attributes (e.g. d, r or fill). Inspired by prior works [4, 58], we transform these instructions into a structured, rule-based matrix representation. We formally define an SVG code as a collection of n vector primitives: $\text{SVG } \mathcal{G} = \{E_1, E_2, \dots, E_n\} = \{E_i\}_{i=1}^n$, where each E_i corresponds to an SVG element such as <path>, <rect>, <circle>, etc. Each element is defined as: $E_i = \{S_{i,j}, F_{i,j}, \alpha_{i,j}\}_{j=1}^{M_i}$, where $S_{i,j}$ is the j -th command in the i -th element, $F_{i,j} \in \{r, g, b\}$ and $\alpha_i \in [0, 1]$ correspond to the color property and the visibility of the i -th element, respectively. M_i indicates the total number of commands in E_i . Thus, the total number of primitive commands in an SVG is $N = \sum_{i=1}^n M_i$. In our work, we set $N = 1024$, determining the maximum number of commands in an SVG representation. Notably, the <path> element consists of multiple commands, while other elements typically contain only a single command.

Figure 3 illustrates the process of converting an SVG code into an SVG embedding. We begin by transforming each primitive into a vector representation, which we then organize into an SVG matrix. Specifically, each command $S_{i,j}$ is represented as $S_{i,j} =$

$(\rho, \tau, \mu_0, \nu_0, \mu_1, \nu_1, \mu_2, \nu_2, \mu_3, \nu_3)_j^i$, where ρ and τ represent the type index of an element and command, respectively. Taking the element <rect> as an example, its element index ρ is 6, and its command index τ is 0 (as it only has no <path> command), its $(\mu_i, \nu_j)_{i=0, j=0}^3$ correspond to x, y, rx, ry, width and height. Further details are provided in Table S1 in Supplementary. After converting into a matrix, we apply an embedding layer to transform ρ and τ from discrete indices into continuous representations while normalizing coordinates and colors in the matrix. Finally, we obtain the SVG embedding by adding the positional embedding.

Compared to previous approaches [4, 58, 63], our method supports a broader range of SVG primitives, including more elements (e.g., <circle>, <rect>) and commands (e.g. Q, A for <path>). A full list of primitives supported by our model is shown in Table S1 of Supplementary. This enhancement significantly improves the model’s ability to learn from real-world data, making the generated SVGs more structured and editable.

3.2. Vector-Pixel Fusion VAE

The Vector-Pixel Fusion VAE (VP-VAE) is designed to learn a latent space for vector graphics with a transformer-based VAE architecture, which includes an encoder and a decoder. To effectively capture both structural and visual features, we introduce Vector-Pixel Fusion Encoding. This method integrates information from the SVG code and its rendered image. Additionally, to better emulate human design logic, we design a sequence modeling strategy. This strategy enables the model to learn how vector graphics are incrementally constructed, ensuring the resultant SVGs with more reasonable construction.

Vector-Pixel Fusion Encoding. As shown in Fig. 4,

to effectively capture both the geometric structure and visual appearance of vector graphics, the VP-VAE encoder $\mathcal{E}(\cdot)$ takes two types of embeddings as the input: (1) SVG embeddings $E_{\text{emb}} \in \mathbb{R}^{N \times D_e}$. (2) Pixel embeddings $E_{\text{pix}} \in \mathbb{R}^{N \times D_p}$. The pixel embeddings are obtained by using a pretrained DINOv2 [31] to extract high-level visual features from the rendered images. To align these embeddings in the latent space, we first project them onto a common dimension D_z using separate linear layers: $E'_{\text{emb}} = W_{\text{proj}}^e E_{\text{emb}}^\top$, $E'_{\text{pix}} = W_{\text{proj}}^p E_{\text{pix}}^\top$, where $W_{\text{proj}}^e \in \mathbb{R}^{D_z \times D_e}$, $W_{\text{proj}}^p \in \mathbb{R}^{D_z \times D_p}$. The transformed embeddings are then subjected to a cross-attention layer, where the SVG embeddings act as queries (Q), while the pixel embeddings act as keys (K) and values (V). This enables the model to effectively integrate geometric and visual features, enhancing its representation of complex vector structures. Finally, the fused representation is processed through the subsequent self-attention layers and mapped to a latent distribution, from which the latent variable z is sampled.

Rendering Sequence Modeling. To enhance the construction of the generated SVGs, we introduce a Rendering Sequence Modeling strategy that enables the model to learn the creation logic of human designers. This strategy involves demonstrating to the model how an SVG is constructed from scratch by observing changes in both SVG codes and their corresponding rasterizations. During training, we represent the SVG as a progressive sequence of drawing steps. Each step incrementally adds new primitives to the SVG, simulating the progressive construction along the batch dimension. As illustrated with the cat example in Fig. 2, each SVG is processed into a batch of data of size B , where each sample corresponds to a different stage of creation, ranging from initial primitives to the complete SVG. This allows the model to observe the process of SVG creation within a single training iteration, thereby fostering a deeper understanding of the logical layers involved in SVG.

The Architecture of VP-VAE. The encoder of the VP-VAE consists of one cross-attention layer followed by L self-attention layers. Specifically, the SVG embeddings initially serve as queries to the pixel embeddings through the cross-attention layer. The tokens resulting from this interaction are then processed by the subsequent self-attention layers. The final output, latent code z , from the encoder is designed to encapsulate both visual and geometric features of the SVG.

The decoder of the VP-VAE mirrors the structure of the encoder but omits the cross-attention layer. Given a latent code z , the decoder reconstructs its corresponding SVG embeddings: $E_{\text{rec}} = \mathcal{D}(z)$, where $\mathcal{D}(\cdot)$ represents the decoder network. To obtain the final SVG representation, we apply two separate decoding processes:

1) *Coordinate Mapping:* The reconstructed embeddings

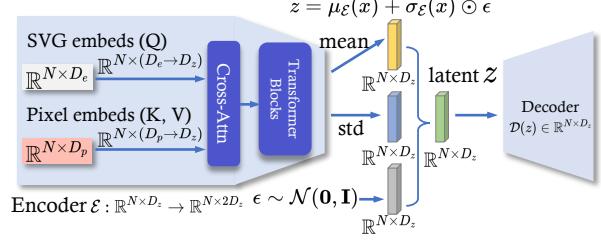


Figure 4. **Illustration of the Vector-Pixel Fusion Encoding.** The VP-VAE encoder integrates the SVG embeddings (Q) with pixel embeddings (K, V) using a cross-attention layer. After processing through L self-attention layers, the encoded features are mapped to a latent space, where the mean and standard deviation are computed for a probabilistic representation. A latent variable z is sampled using the reparameterization trick and then passed to the decoder for further processing. For clarity, the batch dimension B has been omitted.

containing normalized coordinates C_{rec} (ranging from $[-1, 1]$) are mapped back to the canvas coordinate system: $C_{\text{coord}} = (C_{\text{rec}} + 1)/2 \times V$, where V indicates the pre-defined canvas size. 2) *Element & Command Mapping:* The SVG elements and commands are recovered from their respective embeddings using a learned embedding-to-token mapping. Specifically, the reconstructed embeddings are passed through an embedding layer, which predicts the discrete SVG elements and commands, effectively reconstructing the full SVG structure. These steps ensure that the decoded SVG faithfully preserves both the structural and geometric information encoded in the latent space.

VP-VAE Objective. To ensure accurate reconstruction of vector primitives, we measure the discrepancy between the predicted primitive E_i and the ground-truth primitive \hat{E}_i using the mean squared error (MSE) loss. This loss encourages the model to generate primitives that closely match the original input. Beyond reconstruction, we impose regularization on the latent space using Kullback-Leibler (KL) divergence, which constrains the learned latent distribution to approximate a standard Gaussian prior. This prevents overfitting and ensures smooth, continuous latent representations, which are crucial for generating diverse and coherent vector structures. The final VP-VAE loss is formulated as:

$$\mathcal{L}_{\text{VAE}} = \mathcal{L}_{\text{MSE}}(E_i, \hat{E}_i) + \lambda_{\text{KL}} \mathcal{L}_{\text{KL}} \quad (1)$$

where \mathcal{L}_{KL} measures the divergence between the learned latent distribution and the prior Gaussian distribution $\mathcal{N}(0, I)$. The KL term encourages the latent space to remain compact and structured, facilitating smooth interpolation between vector primitives.

3.3. Vector Space Diffusion Transformer

The VP-VAE effectively learns a latent space tailored for vector graphics representation. Building on this founda-

tion, SVGFusion leverages DiT [32] as its core architecture and performs the diffusion process directly in the vector latent space. To facilitate interaction between textual features and vector representations, we introduce a multi-head cross-attention layer into the VS-DiT block, inspired by [5]. A detailed description of the VS-DiT architecture is provided in Supplementary Sec. C. During training, for a given input SVG, we first derive its latent representation using VP-VAE: $\mathbf{z} = \mathcal{E}(\mathcal{G})$. The diffusion process then takes place within this latent space, where the noisy latent variable \mathbf{z}_t is generated as follows: $\mathbf{z}_t = \alpha_t \mathbf{z} + \sigma_t \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ where α_t and σ_t define the noise schedule, parameterized by the diffusion time t . Following Latent Diffusion Models (LDM) [38], our VS-DiT model predicts the noise ϵ in the noisy latent representation \mathbf{z}_t , conditioned on a text prompt y . The training objective is formulated as:

$$\mathcal{L}_{\text{ldm}} = \mathbb{E}_{\mathcal{E}(\mathbf{z}), y, t, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\|\epsilon_\phi(\mathbf{z}_t, t, y) - \epsilon\|_2^2 \right] \quad (2)$$

where $t \sim \mathcal{U}(0, 1)$. We randomly set y to zero with a probability of 10% to apply classifier-free guidance [15] during training, which enhances the quality of conditional generation during inference.

4. Experiments

4.1. SVGX-Dataset

We introduce *SVGX-Dataset*, a large-scale SVG dataset designed for model training and evaluation. It comprises 240,000 high-quality *emoji* and *icon*-style SVGs sourced from Twemoji-Color-Font [54], Noto-Emoji [10], FluentUI-Emoji [27], SVG-Repo [62], and Reshot [61].

The dataset spans a diverse range of visual complexities, incorporating both Bézier curves (expressed using `<path>`) and basic geometric shapes (`<circle>`, `<rect>`, etc.). It includes high-quality, human-designed SVGs in black-and-white and color, covering various semantic categories such as people, animals, objects, and symbols. Examples are shown in Fig. S1 in the Supplementary.

SVGs collected from the Internet often contain noise, which can hinder model performance. We identify three common issues: (1) temporary data generated by vector editing software, (2) suboptimal structural representations, and (3) unused or invisible graphical elements.

Although these elements do not affect the visual appearance, they introduce inefficiencies in training. To address this, we implement a cleaning pipeline that removes redundant elements, refines coordinate precision, and resizes canvases to 128×128 . This process significantly reduces file sizes while preserving visual fidelity. An example of the cleaning pipeline is illustrated in Fig. S3 in the Supplementary.

Method / Metric	FID \downarrow	CLIPScore \uparrow	Aesthetic \uparrow	HPS \uparrow	TimeCost \downarrow
Evolution [52]	121.43	0.193	2.124	0.115	47min23s
CLIPDraw [9]	116.65	0.249	3.980	0.135	5min10s
DiffSketcher[65]	72.30	0.310	5.156	0.242	10min22s
LIVE+VF [20]	82.22	0.310	4.517	0.253	30min01s
VectorFusion [20]	84.53	0.309	4.985	0.264	10min12s
Word-As-Img [19]	101.22	0.302	3.276	0.151	5min25s
SVGDreamer [68]	70.10	0.360	5.543	0.269	35min12s
SVG-VAE [23]	76.22	0.190	2.773	0.101	1min
DeepSVG [4]	69.22	0.212	3.019	0.114	2min
Iconshop [63]	52.22	0.251	3.474	0.140	1min03s
StrokeNUWA [49]	89.10	0.300	2.543	0.169	19s
SVGFusion-S	9.62	0.373	5.250	0.275	24s
SVGFusion-B	5.77	0.389	5.373	0.281	28s
SVGFusion-L	4.64	0.399	5.673	0.290	36s

Table 1. **Quantitative Comparison of SVGFusion vs. State-of-the-Art Text-to-SVG Methods.** Top: Optimization-based methods; Middle: language-model-based methods; Bottom: Three model variants of our proposed SVGFusion. Our method uses dpm-solver [24] for 20-step denoising.

As shown in Fig. S2, we analyze frequently occurring words in dataset names and descriptions, visualize word distributions via a word cloud, and compare data before and after cleaning. This analysis highlights the dataset’s well-structured nature, supporting effective model training. Further details on dataset collection, preprocessing, and analysis are provided in Supplementary Sec. B.

4.2. Quantitative Evaluation

We compare our proposed method with baseline methods using five quantitative indicators across three dimensions: (1) Visual quality of the generated SVGs, assessed by FID (Fréchet Inception Distance) [14]; (2) Alignment with the input text prompt, assessed by CLIP score [35], and (3) Aesthetic appeal of the generated SVGs, measured by Aesthetic score [40] and HPS (Human Preference Score) [64]. To ensure a fair comparison, we also recorded the time cost of different methods to evaluate their computational efficiency.

Comparison results are presented in Table 1. The methods are categorized into two groups: optimization-based methods (top section of Table 1) and language-model-based methods (middle section of Table 1). It is evident that our SVGFusion method surpasses other text-to-SVG methods across all evaluation metrics. This demonstrates the superiority of SVGFusion in generating vector graphics that are more closely aligned with text prompts and human preferences. Notably, compared to optimization-based methods, SVGFusion significantly reduces the time cost, enhancing its practicality and user-friendliness.

4.3. Qualitative Evaluation

Figure 5 presents a qualitative comparison between SVGFusion and existing text-to-SVG methods. The results are aligned with the quantitative results discussed in the previous section. Specifically, the

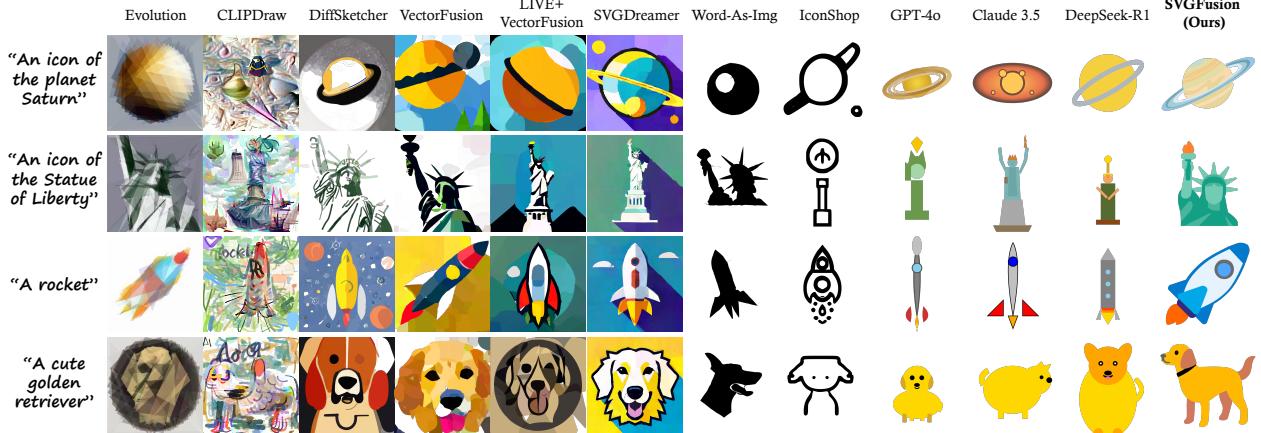


Figure 5. **Qualitative Comparison of SVGFusion and Existing Text-to-SVG Methods.** The target SVGs are in the *emoji* style. We use prompt modifiers for the optimization-based approach to encourage the appropriate style: “minimal flat 2D vector icon, emoji icon, lineal color, on a white background, trending on ArtStation.” Note that although the visual quality of results generated by optimization-based methods is high, these methods face challenges in decomposing the SVGs for further editing.

optimization-based methods, including Evolution [52], CLIPDraw [9], DiffSketcher [65], VectorFusion [20], LIVE [26]+VectorFusion [20], SVGDreamer [68], and Word-as-Img [19], use a differentiable renderer [21] to backpropagate gradients to vector parameters. Evolution [52] and CLIPDraw [9] utilize CLIP [35] as the image prior, while DiffSketcher [65], VectorFusion [20], SVGDreamer [68], and Word-as-Img [19] adopt T2I diffusion as the image prior. Despite their visual advantages, optimization-based methods often produce intertwined vector primitives, diminishing SVG editability. Language-model-based methods, such as Iconshop [63], Claude3.5-sonnet [2], GPT4o-latest [1] and DeepSeek-R1 [11], rely on language models, and can generate decoupled vector primitives but overly simplistic content.

It is worth noting that although optimization-based methods may produce more realistic or artistic visual effects, they rely on an LDM [38] sample as the target for optimization and require a differentiable rasterizer as the medium for this process. Additionally, they depend on differentiable vector primitives as the underlying representation for SVGs. As a result, these methods can only use `<path>` primitives described by Bézier curves. This leads to the need for a large amount of staggered overlapping primitives to closely fit the LDM sample, even for relatively simple shapes. Consequently, even simple regular shapes such as rectangular cannot be described using the corresponding basic shape primitives, thus losing the advantage of SVG’s editability, making it difficult to use in real-world scenarios.

As illustrated in Fig. 1(a) and Fig. 6, our SVGFusion can generate SVGs using only the necessary primitives, such as `<circle>` and `<rect>`, ensuring a compact and structured representation. It allows for flexibility in the design process: one can either start by sketching the general shape and then add local details, or begin with a

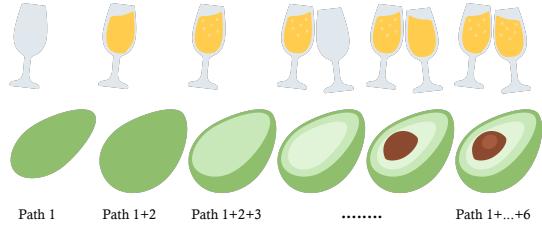


Figure 6. **Path Rendering Sequence.** SVGFusion is designed to align with human logic in SVG creation. The top diagram illustrates the left-to-right sequence of object placement, while the bottom diagram depicts the drawing order of an object from simple to complex.

specific part and gradually add elements to complete the SVG. This coincides with the process of creating SVGs by human designers.

Comparison with Large Language Model. As illustrated in Fig. 5, we also compare our proposed SVGFusion with existing state-of-the-art approaches that directly generate SVGs using Large Language Models (LLMs). The results indicate that the performance of Claude3.5-sonnet [2], GPT4o-latest [1], and DeepSeek-R1 [11] in SVG generation is not particularly outstanding. In most cases, they can only use simple shapes to roughly assemble objects, but the positioning of each element lacks harmony. As a result, the overall shapes are overly simplistic, and the visual effects are less satisfactory. Regarding color design, these LLMs struggle to apply colors accurately to each part of the SVG, leading to color schemes that are often neither harmonious nor reasonable. In terms of semantic expression, the SVG code generated by LLMs is too simple to fully capture the meaning conveyed by the input text description. In contrast, our proposed SVGFusion method produces more balanced and harmonious results in terms of shape selection, color matching, and semantic representation. In Supplementary Sec. E, we provide more comparison of

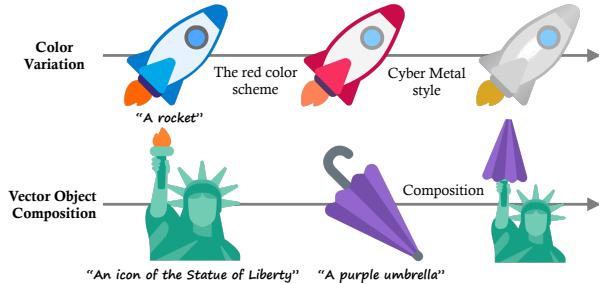


Figure 7. The Editability of SVGFusion Results. The SVGs generated by SVGFusion exhibit a clear hierarchical structure, thereby facilitating straightforward edits, such as color modifications (top example) or recomposition into new graphics (bottom example).

Model	Layer N_b	Hidden size d	Heads	Gflops
VS-DiT S	12	384	6	1.4
VS-DiT B	12	768	12	5.6
VS-DiT L	24	1024	16	19.9

Table 2. Details of VS-DiT Models. We follow model configurations for the Small (S), Base (B), and Large (L) variants.

our method with language model-based methods.

Editbility of SVGFusion Results. Figure 7 shows the editability of the SVG generated by our SVGFusion. Since the SVGs we generate have a clean and concise structure, we can easily edit the properties of the primitives, such as their color attributes. For instance, the rocket we generated can be changed from blue to red, or even transformed into a cyber-metal style, simply by adjusting the color attributes. Furthermore, the capability of editbility empowers users to efficiently reuse synthesized vector elements and create new vector compositions. As illustrated in the second example of Fig. 7, our method composes a new vector graphic by replacing the torch with an umbrella.

4.4. Ablation Study

Model Size. We apply a sequence of N_b DiT blocks, each operating at the hidden dimension size d . Following DiT [32], we use standard Transformer configs that jointly scale N_b , d and attention heads. Specifically, we use three configs: VS-DiT S, VS-DiT B, VS-DiT L. They cover a wide range of model sizes and flop allocations, from 1.4 to 19.9 Gflops, allowing us to gauge scaling performance. The details of the configs are provided in Table 2.

Ablation on Vector-Pixel Fusion Encoding. In SVGFusion, we introduce the Visual-Pixel Fusion Autoencoder (VP-VAE) to learn a continuous latent space for SVGs. This module processes both SVG codes and their corresponding rendered images as inputs, aligning the SVG codes with visual features during the encoding phase. We evaluate the impact of integrating visual features by comparing results from model variants with

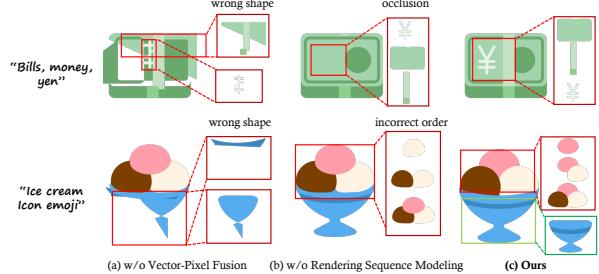


Figure 8. Effects of VP-VAE and Rendering Sequence Modeling. (a) vs. (c) demonstrates that VP-VAE cannot accurately reconstruct or generate shapes without the Vector-Pixel Fusion (incorporating DINOv2 visual prior [31]). (b) vs. (c) indicates that employing Rendering Sequence Modeling results in more reasonable SVG outcomes, due to the order of primitives being better aligned with human creation logic.

and without this feature. As illustrated in Fig. 8, the variant incorporating Vector-Pixel Fusion Encoding demonstrates superior visual quality in the generated SVGs.

Furthermore, we conducted experiments to validate the effectiveness of the feature extractor, DINOv2, employed in our approach. Detailed results and analyses of these experiments are available in the Supplementary Sec. D.

Ablation on Rendering Sequence Modeling. The Rendering Sequence Modeling strategy in SVGFusion is designed to enhance the construction of generated SVGs, making them more editable. In Fig. 8, we evaluate the impact of employing this strategy compared to scenarios where it is omitted. The results clearly show that SVGs generated with Rendering Sequence Modeling exhibit higher visual quality and improved structural integrity.

In the first example, issues arise when shapes with lighter colors are occluded by shapes with darker colors, leading to a poorly represented money icon. In the second example, the creation order of the scoops is critical as their relational positioning greatly influences the visual coherence. These examples underscore the effectiveness of our Rendering Sequence Modeling strategy.

5. Conclusion & Discussion

In this work, we have introduced SVGFusion, a novel Text-to-SVG model that effectively scales to real-world vector graphics without relying on text-based discrete language models or prolonged optimization. By learning a continuous latent space through VP-VAE and generating structured SVG representations with VS-DiT, our approach captures the underlying creation logic of vector graphics, enabling high-quality and semantically coherent generation. To support this, we have curated a dataset called *SVGX-Dataset*. Extensive experiments validate the superiority of our SVGFusion over existing

methods, demonstrating its potential as a scalable and efficient framework for SVG content creation.

Limitations. Due to the limited availability of publicly accessible data, the dataset collected in this work covers a restricted range of vector graphic types. Incorporating a broader diversity of vector graphics would not only expand the dataset size but also further explore the model’s scalability potential.

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. 7
- [2] Anthropic. Claude 3.5 sonnet. <https://www.anthropic.com/news/clause-3-5-sonnet>, 2024. 7
- [3] Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth words: A vit backbone for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22669–22679, 2023. 3
- [4] Alexandre Carlier, Martin Danelljan, Alexandre Alahi, and Radu Timofte. Deepsvg: A hierarchical generative network for vector graphics animation. *Advances in Neural Information Processing Systems (NeurIPS)*, 33: 16351–16361, 2020. 2, 3, 4, 6, 14, 16
- [5] Junsong Chen, Jincheng YU, Chongjian GE, Lewei Yao, Enze Xie, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart-\$\alpha\$: Fast training of diffusion transformer for photorealistic text-to-image synthesis. In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024. 6
- [6] Louis Clouâtre and Marc Demers. Figr: Few-shot image generation with reptile. *arXiv preprint arXiv:1901.02199*, 2019. 3
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, 2009. 3
- [8] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems (NeurIPS)*, 34:8780–8794, 2021. 3
- [9] Kevin Frans, Lisa Soros, and Olaf Witkowski. CLIP-Draw: Exploring text-to-drawing synthesis through language-image encoders. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2, 3, 6, 7
- [10] Google. Noto emoji fonts. <https://github.com/googlefonts/noto-emoji>, 2014. 3, 6, 12
- [11] 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. 7
- [12] Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024. 3
- [13] David Ha and Douglas Eck. A neural representation of sketch drawings. In *International Conference on Learning Representations (ICLR)*, 2018. 2, 3
- [14] 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. *Advances in neural information processing systems (NeurIPS)*, 30, 2017. 6
- [15] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 3, 6
- [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 6840–6851, 2020. 3
- [17] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:8633–8646, 2022. 3
- [18] Teng Hu, Ran Yi, Baihong Qian, Jiangning Zhang, Paul L Rosin, and Yu-Kun Lai. Supersvg: Superpixel-based scalable vector graphics synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 24892–24901, 2024. 3
- [19] Shir Iluz, Yael Vinker, Amir Hertz, Daniel Berio, Daniel Cohen-Or, and Ariel Shamir. Word-as-image for semantic typography. *ACM Transactions on Graphics (TOG)*, 42(4), 2023. 6, 7
- [20] Ajay Jain, Amber Xie, and Pieter Abbeel. Vectorfusion: Text-to-svg by abstracting pixel-based diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 3, 6, 7
- [21] Tzu-Mao Li, Michal Lukáč, Gharbi Michaël, and Jonathan Ragan-Kelley. Differentiable vector graphics rasterization for editing and learning. *ACM Transactions on Graphics (TOG)*, 39(6):193:1–193:15, 2020. 2, 3, 7
- [22] Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao, Ruoxi Chen, Zhengqing Yuan, Yue Huang, Hanchi Sun, Jianfeng Gao, et al. Sora: A review on background, technology, limitations, and opportunities of large vision models. *arXiv preprint arXiv:2402.17177*, 2024. 3
- [23] Raphael Gontijo Lopes, David Ha, Douglas Eck, and Jonathon Shlens. A learned representation for scalable vector graphics. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 3, 6
- [24] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:5775–5787, 2022. 6, 12

- [25] Nanye Ma, Mark Goldstein, Michael S Albergo, Nicholas M Boffi, Eric Vanden-Eijnden, and Saining Xie. Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformers. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2024. 3
- [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 (CVPR)*, pages 16314–16323, 2022. 3, 7
- [27] Microsoft. Fluent emoji. <https://github.com/microsoft/fluentui-emoji>, 2021. 6, 12
- [28] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International conference on machine learning (ICLR)*, pages 8162–8171, 2021. 3
- [29] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. GLIDE: Towards photorealistic image generation and editing with text-guided diffusion models. In *Proceedings of the 39th International Conference on Machine Learning (ICML)*, pages 16784–16804, 2022. 3
- [30] OpenAI. Introducing chatgpt. <https://openai.com/index/chatgpt/>, 2023. 12
- [31] Maxime Oquab, Timothée Darctet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOV2: Learning robust visual features without supervision. *Transactions on Machine Learning Research (TMLR)*, 2024. 2, 5, 8, 14
- [32] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4195–4205, 2023. 2, 3, 6, 8
- [33] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image synthesis. In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024. 3
- [34] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *The Eleventh International Conference on Learning Representations (ICLR)*, 2023. 3
- [35] 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 (ICML)*, pages 8748–8763. PMLR, 2021. 2, 3, 6, 7
- [36] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 3
- [37] Pradyumna Reddy, Michael Gharbi, Michal Lukac, and Niloy J Mitra. Im2vec: Synthesizing vector graphics without vector supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7342–7351, 2021. 3
- [38] 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 (CVPR)*, pages 10684–10695, 2022. 2, 3, 6, 7
- [39] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 36479–36494, 2022. 3
- [40] Christoph Schuhmann. Improved aesthetic predictor. <https://github.com/christophschuhmann/improved-aesthetic-predictor>, 2022. 6
- [41] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-a-video: Text-to-video generation without text-video data. In *The Eleventh International Conference on Learning Representations (ICLR)*, 2023. 3
- [42] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 2256–2265, 2015. 3
- [43] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations (ICLR)*, 2021. 12
- [44] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [45] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations (ICLR)*, 2021. 3
- [46] Yiren Song, Xuning Shao, Kang Chen, Weidong Zhang, Zhongliang Jing, and Minzhe Li. Clipvg: Text-guided image manipulation using differentiable vector graphics. In *Proceedings of the Conference on Artificial Intelligence (AAAI)*, 2023. 2, 3
- [47] StabilityAI. If by deepfloyd lab at stabilityai. <https://github.com/deep-floyd/IF>, 2023. 3

- [48] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomput.*, 568(C), 2024. 12
- [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. 2, 3, 6, 16
- [50] Vikas Thamizharasan, Difan Liu, Shantanu Agarwal, Matthew Fisher, Michaël Gharbi, Oliver Wang, Alec Jacobson, and Evangelos Kalogerakis. Vecfusion: Vector font generation with diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7943–7952, 2024. 3
- [51] Vikas Thamizharasan, Difan Liu, Matthew Fisher, Nanxuan Zhao, Evangelos Kalogerakis, and Michal Lukac. Nivel: Neural implicit vector layers for text-to-vector generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4589–4597, 2024. 3
- [52] 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*, pages 275–291. Springer, 2022. 2, 3, 6, 7
- [53] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *ArXiv*, abs/2302.13971, 2023. 12
- [54] Twitter. Twitter color emoji svginot font. <https://github.com/13rac1/twemoji-color-font>, 2016. 6, 12
- [55] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*. Curran Associates, Inc., 2017. 12
- [56] Yael Vinker, Ehsan Pajouheshgar, Jessica Y Bo, Roman 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, 3
- [57] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A. Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12619–12629, 2023. 3
- [58] Yizhi Wang and Zhouhui Lian. Deepvecfont: Synthesizing high-quality vector fonts via dual-modality learning. *ACM Transactions on Graphics (TOG)*, 40(6), 2021. 2, 3, 4, 14, 16
- [59] 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 (CVPR)*, pages 18320–18328, 2023. 2, 16
- [60] 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 distillation. In *Thirty-seventh Conference on Neural Information Processing Systems (NeurIPS)*, 2023. 3
- [61] ReShot website. Reshot - free icons & illustrations. design freely with instant downloads and commercial licenses. <https://www.reshot.com/>, . 6, 12
- [62] SVGRepo website. Open-licensed svg vector and icons. <https://www.svgrepo.com/>, . 6, 12
- [63] Ronghuan Wu, Wanchao Su, Kede Ma, and Jing Liao. Iconshop: Text-based vector icon synthesis with autoregressive transformers. *arXiv preprint arXiv:2304.14400*, 2023. 2, 3, 4, 6, 7, 14, 16
- [64] Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score: Better aligning text-to-image models with human preference. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2096–2105, 2023. 6
- [65] Ximing Xing, Chuang Wang, Haitao Zhou, Jing Zhang, Qian Yu, and Dong Xu. Diffsketcher: Text guided vector sketch synthesis through latent diffusion models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. 2, 3, 6, 7
- [66] 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. 2
- [67] Ximing Xing, Qian Yu, Chuang Wang, Haitao Zhou, Jing Zhang, and Dong Xu. Svdreamer++: Advancing editability and diversity in text-guided svg generation. *arXiv preprint arXiv:2411.17832*, 2024. 3
- [68] Ximing Xing, Haitao Zhou, Chuang Wang, Jing Zhang, Dong Xu, and Qian Yu. Svdreamer: Text guided svg generation with diffusion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4546–4555, 2024. 2, 3, 6, 7
- [69] Peiyang Zhang, Nanxuan Zhao, and Jing Liao. Text-to-vector generation with neural path representation. *ACM Trans. Graph.*, 43(4), 2024. 3, 14

SVGFusion: Scalable Text-to-SVG Generation via Vector Space Diffusion

Supplementary Material

Overview

This supplementary material provides additional details and analyses related to **SVGFusion**, organized as follows:

- Appendix A: *Implementation Details*. We describe the implementation details of SVGFusion, including hyperparameter settings and training configurations.
- Appendix B: *Dataset Details*. We provide an in-depth overview of our newly introduced dataset, covering data representation, preprocessing, and cleaning procedures.
- Appendix C: *VS-DiT Architecture*. We present details of the Vector Space Diffusion Transformer (VS-DiT) architecture, explaining its components and design choices.
- Appendix D: *Ablation Studies*. We conduct ablation experiments to analyze the impact of different feature extractors on the performance of SVGFusion.
- Appendix E: *Comparison with Language Model-Based Methods*. We discuss the key differences between diffusion-based and language model-based SVG generation approaches, highlighting how SVGFusion leverages an image-based paradigm for improved SVG synthesis.
- Appendix F: *SVG Primitive Support*. We demonstrate that SVGFusion supports a broader range of SVG primitives than existing methods, enabling it to better capture the structural and visual diversity of real-world vector graphics.
- Appendix G: *Additional Results*. We provide qualitative results, including visualizations of the diffusion process and an extended set of generated SVGs, showcasing SVGFusion’s ability to black-and-white icon style vector graphics.

A. Implementation Details

To ensure consistency across all SVG data, we adopted relative positional coordinates. The model parameters are initialized randomly and optimized using the AdamW optimizer (with $\beta_1 = 0.9$, $\beta_2 = 0.999$) at an initial learning rate of 3×10^{-4} . The learning rate is warmed up over the first 2000 steps and then decayed to 1.5×10^{-5} following a cosine schedule. Additionally, we applied weight decay of 0.1 for regularization and constrain gradients by clipping their norms to a maximum value of 2.0. Besides, we normalized the input SVG embeddings into the $[-1, 1]$ range to stabilize the training process. To further enhance the model’s ability



Figure S1. Samples from our collected SVG dataset. The dataset contains a diverse set of SVG illustrations, including outlined and filled designs, covering themes such as nature, objects, symbols, animals, and food.

to capture sequential dependencies, we integrate rotary position embeddings (RoPE) [48], a technique widely used in advanced large language models [30, 53]. RoPE effectively encodes positional relationships within sequences, allowing the model to better understand the temporal and structural progression of SVG creation. We utilized Transformers [55] as the fundamental building block for VP-VAE. Both the encoders and decoders are based on the Transformer architecture, consisting of 4 layers, a feed-forward dimension of 512. To investigate scaling trends, we trained VS-DiT models at three different sizes: 0.16B, 0.37B, and 0.76B parameters. We trained VP-VAE for 1000k steps using a total batch size of 128 across 2 A800 GPUs, requiring approximately two day. Subsequently, leveraging the trained VP-VAE, we trained VS-DiT for 500k steps with a batch size of 256 on 8 A800 GPUs, taking approximately four days. In the sampling phase, we used DDIM [43] to de-noise 100 steps; dpm-solver [24] also supports 20-step denoising for more efficient sampling.

B. SVG Dataset Collection, Representation, Cleaning and Analysis

SVG Data Collection. We collected a large-scale SVG dataset for model training and evaluation, which includes 240k high-quality SVGs with *emoji* and *icon* style. The dataset consists of emoji and icon SVGs sourced from various sources, including Twemoji-Color-Font [54], Noto-Emoji [10], FluentUI-Emoji [27], SVG-Repo [62] and Reshot [61]. Twemoji-Color-Font, Noto-Emoji, and FluentUI-Emoji each contribute approximately 4,000 SVGs. An additional 30,000 SVGs are sourced from Reshot, while the majority—200,000 SVGs—are contributed by SVGRepo.

SVG Data Representation In terms of visual complexity, our collected dataset includes primitives that vary from simple to intricate. Regarding primitive representation, the dataset includes both Bézier curves ex-

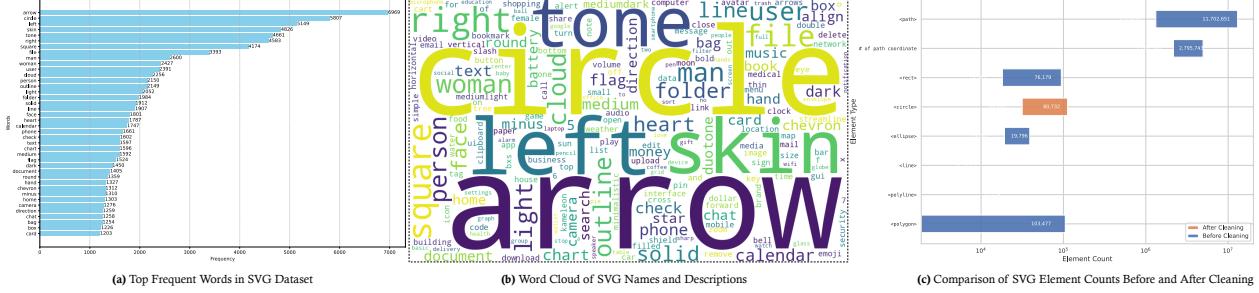


Figure S2. Overview of our collected SVG dataset analysis. (a) displays the most frequently appearing words in the SVG dataset's names and descriptions. Common words such as “arrow”, “circle”, and “left” suggest a predominance of directional and geometric elements in the dataset. (b) This word cloud visualizes the most common words found in the dataset. Larger words indicate higher frequency, highlighting key themes such as geometric shapes (“circle”, “square”), directional terms (“left”, “right”), and common UI elements (“file”, “folder”). (c) compares the number of different SVG elements (e.g., `<path>`, `<rect>`, `<circle>`) before and after data cleaning. The reduction in element count suggests significant optimization, potentially removing redundant or unnecessary elements.

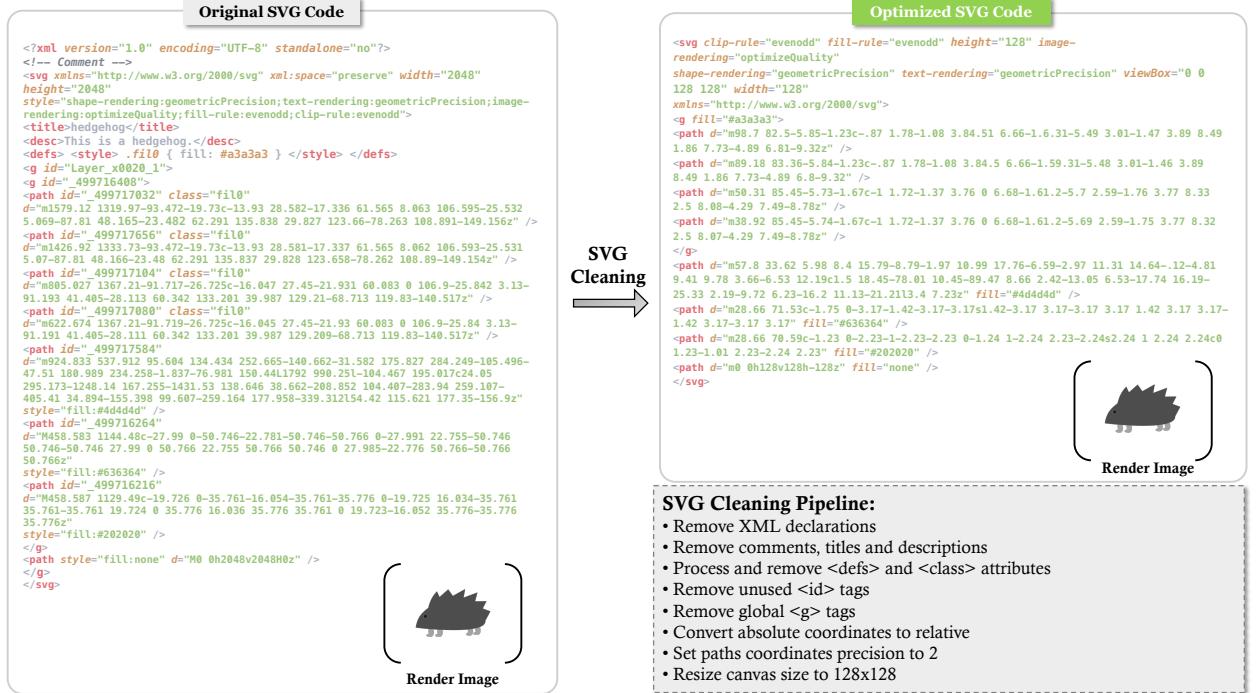


Figure S3. Overview of the SVG cleaning process. The pipeline illustrates the transformation of an original SVG file into a preprocessed version through an SVG cleaning pipeline. The pipeline involves several steps: removing XML declarations, comments, titles, and descriptions; processing and eliminating `<defs>` and `<class>` attributes; removing unused `<id>` tags and global `<g>` tags; converting absolute coordinates to relative; setting path coordinate precision to 2; and resizing the canvas to 128×128 pixels. The cleaned SVG results in a more optimized and compact representation while preserving the visual integrity of the image.

pressed using `<path>` and basic shapes represented by `<circle>` and `<rect>`, etc. As illustrated in Fig. S1, in terms of color representation, the dataset contains both black-and-white graphics and vibrant, harmonious color images. Additionally, in terms of semantics, it covers a wide range of subjects, including people, animals, objects, and symbols. Such a dataset enables us to design a model capable of leveraging prior knowledge of human-created SVGs to synthesize high-quality SVG

outputs

SVG Data Cleaning Pipeline. However, SVG files collected from the web often contain noise, such as redundant data stemming from different designers' workflows, and using it directly for learning can hinder the model's representation capabilities. As illustrated in Fig. S3, about half of the data in an SVG file is redundant for visual rendering, including (1) temporary data from vector editing applications, (2) non-optimal struc-

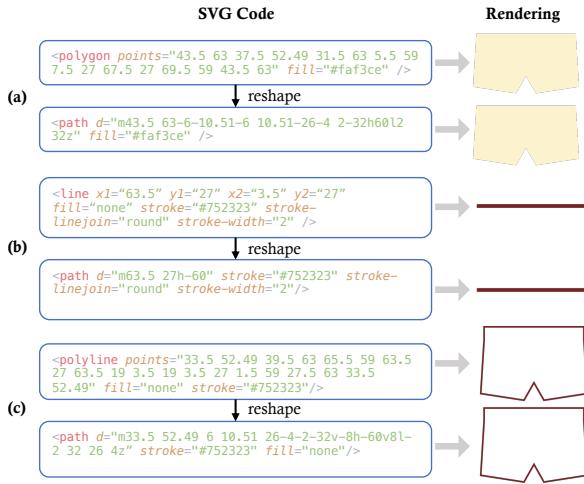
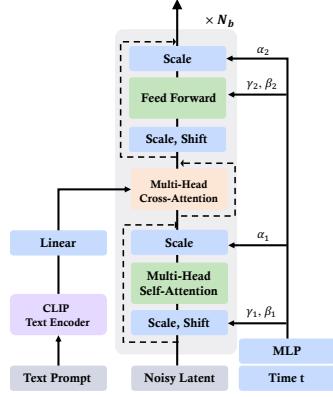


Figure S4. **Illustration of SVG primitive reshaping.** We show that (a) `<polygon>`, (b) `<line>`, and (c) `<polyline>` can be losslessly transformed into `<path>`.



tural representations, and (3) unused or invisible graphic elements. These components are superfluous in finalized, published SVG files and can be safely removed to improve file efficiency. We propose a preprocessing pipeline that losslessly reduces the size of SVG files by eliminating unnecessary XML declarations, comments, titles, and descriptions, as well as removing unused SVG tags. Additionally, we convert absolute coordinates to relative ones, set the precision of path coordinates to two decimal places, and resize the canvas to 128×128 .

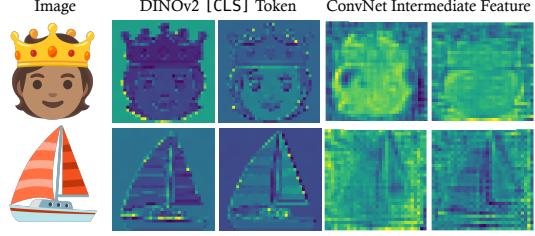


Figure S6. **Comparison of Feature Extractors.** We employ pretrained DINOv2 [31] as our feature extractor, as it provides a more precise spatial-geometric representation of objects without interference from background noise. In contrast, previous methods [58, 69] relied on ConvNet-based feature extraction, which introduced additional background noise, ultimately degrading the quality of geometric representations.

SVG Dataset Analysis. As illustrated in Fig. S2, our dataset analysis focuses on three key aspects: (a) the 40 most frequently occurring words in the names of the SVG dataset, (b) a word cloud representing the most common words found in SVG names and descriptions, and (c) a comparative analysis of the dataset before and after the cleaning process.

C. VS-DiT Block Architecture

As illustrated in Fig. S5, it is positioned between the self-attention layer and the feed-forward layer so that the model can flexibly interact with the text embedding extracted from the language model. Our architecture supports flexible stacking of N_b VS-DiT blocks, enabling models with varied parameter scales that allow SVGFusion to adaptively expand with increasing data volume, thus enhancing scalability with data growth.

D. Additional Ablation Studies

Comparison of Feature Extractors. As shown in Fig. S6, we utilize pretrained DINOv2 [31] as our feature extractor to capture precise geometric features of objects at the pixel level. In contrast, [58, 69] trained ConvNet for feature extraction, which was less effective at pixel-level characterization and introduced significant background noise.

E. Comparison with Language Model-based Methods

Language model-based methods such as DeepSVG [4] and IconShop [63] treat SVGs as sequences of tokens and predict outcomes in an autoregressive manner. However, the sequential nature of autoregressive approaches inherently limits them. As illustrated in the first row of Fig. S7, the prediction of each token heavily depends on the accuracy of the preceding tokens. Consequently, errors can propagate throughout the sequence,

SVG Primitives (Element/Command)	Argument	Explanation	Example
<code><circle></code>	r, cx, cy	The <code><circle></code> element is used to create a circle with center at (cx, cy) and radius r .	
<code><ellipse></code>	rx, ry, cx, cy	The <code><ellipse></code> element is used to create an ellipse with center at (cx, cy) , and radii rx and ry .	
<code><rect></code>	rx, ry, cx, cy	The <code><rect></code> element is used to create a rectangle, optionally with rounded corners if rx and ry are specified. The center is at (cx, cy) .	
<code><path> Move To (M)</code>	μ_3, ν_3	Moves the cursor to the specified point (μ_3, ν_3) .	
<code><path> Line To (L)</code>	μ_3, ν_3	Draws a line segment from the current point to (μ_3, ν_3) .	
<code><path> Cubic Bézier (C)</code>	$\mu_1, \nu_1, \mu_2, \nu_2, \mu_3, \nu_3$	Draws a cubic Bézier curve with control points (μ_1, ν_1) , (μ_2, ν_2) , and endpoint (μ_3, ν_3) .	
<code><path> Quadratic Bézier (Q)</code>	$\mu_1, \nu_1, \mu_3, \nu_3$	Draws a quadratic Bézier curve with control points (μ_1, ν_1) and endpoint (μ_3, ν_3) .	
<code><path> Elliptical Arc (A)</code>	$rx, ry, rotate,$ $LargeArcFlag,$ $SweepFlag, \mu_3, \nu_3$	Draws an elliptical arc from the current point to (μ_3, ν_3) . The ellipse has radii rx, ry , rotated by $rotate$ degrees. $LargeArcFlag$ and $SweepFlag$ control the arc direction.	
<code><path> Close Path (Z)</code>	\emptyset	Closes the path by moving the cursor back to the path's starting position (μ_0, ν_0) .	
<code><SOS></code>	\emptyset	Special token indicating the start of an SVG sequence.	N/A
<code><EOS></code>	\emptyset	Special token indicating the end of an SVG sequence.	N/A

Table S1. **SVG Primitives Supported by Our SVGFusion.** Compared to existing works, SVGFusion supports a broader range of elements (shown in the first three rows) and commands (shown in the 7th and 8th rows), which are highlighted in bold.

especially when initial predictions are suboptimal, ultimately leading to degraded performance.

Instead, our SVGFusion employs a diffusion-based approach to synthesize shapes in parallel rather than sequentially. As shown in the second row of Fig. S7, the process begins with random, chaotic strokes, and the SVG elements, such as vector coordinates, are progressively denoised at each timestep. This method overcomes the inherent limitations of autoregressive techniques by mitigating error accumulation. Over successive iterations, the correct content gradually emerges while preserving the visual coherence of the generated SVG.

F. More Primitive Types

SVG code comprises a suite of primitives and syntax rules. For example, the element-level primitive `<rect>` defines a rectangle shape with specified position, width, and height, and can be represented as `<rect x="10" y="20" width="50" height="80"/>`. How-

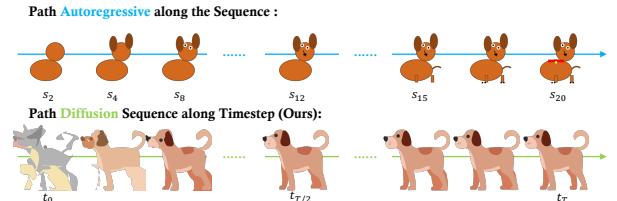


Figure S7. **Comparison of autoregressive language model prediction and diffusion-based generation of SVGFusion.** This comparison highlights the differences between two approaches for SVG generation: autoregressive language models and diffusion-based methods. Autoregressive models generate vector graphics sequentially, predicting one token at a time based on prior outputs, which can lead to error accumulation and limited global coherence. In contrast, SVGFusion leverages diffusion-based generation, which operates in a continuous vector space, enabling more holistic and globally consistent SVG synthesis while mitigating the limitations of discrete token-based generation.

ever, given the multitude of SVG primitive types, it is challenging for Language Model-based methods to model these effectively due to the need for a complex

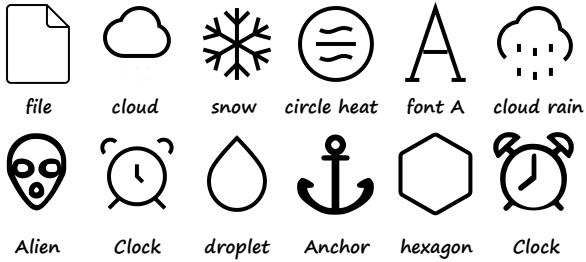


Figure S8. A Set of black-and-white SVG Icons Generated Using **SVGFusion**, illustrating various symbols representing weather, time, shapes, and miscellaneous objects.

data structure and extensive data. Consequently, some studies [4, 49, 58, 59, 63] simplify this task by focusing solely on one element-level primitive, `<path>`, and three command-level primitives: “Move To” (M), “Line To” (L), and “Cubic Bézier” (C). Such simplification is lossy and thereby those methods fail to accurately represent real SVGs as designed by humans.

To address these limitations, we have expanded the SVG representation in this work to support a broader range of SVG elements and commands. Specifically, as detailed in Table S1, our model now handles additional element-level commands such as `<circle>`, `<rect>`, and `<ellipse>`. Furthermore, as demonstrated in Figure S4, elements like `<line>`, `<polygon>`, and `<polyline>` can be losslessly converted into `<path>`. These elements are converted to `<path>` during the data preprocessing phase for consistency and efficiency.

Moreover, we have extended the path command set to include M (Move To), L (Line To), Q (Quadratic Bézier Curve), C (Cubic Bézier Curve), A (Arc To), and Z (Close Path). Other `<path>` commands, such as H (Horizontal Line To) and V (Vertical Line To), are losslessly represented by L, while S (Smooth Cubic Bézier Curve) and T (Smooth Quadratic Bézier Curve) are simplified versions of C and Q, respectively. These modifications streamline the process without sacrificing the accuracy of the representations.

G. Additional Qualitative Results

Icon. As shown in Fig. S8, the proposed method is also adept at synthesizing simple black-and-white icon-style SVGs, which are typically composed of a limited number of primitives and are defined by their minimalist yet expressive design characteristics.

Editability. Figure S9 demonstrates the high editability of SVGs produced by **SVGFusion**. The clean and systematic structure of these SVGs facilitates efficient modification of primitive properties, including color attributes, enabling seamless customization.

Conceptual Combination. The inherent editability of

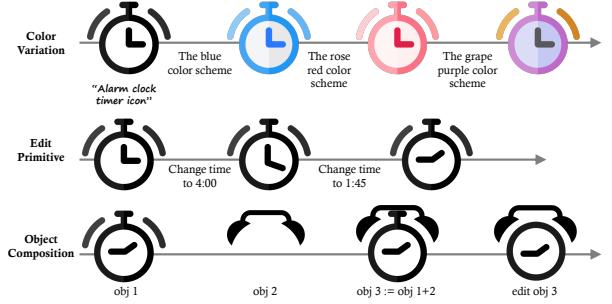


Figure S9. **The editability of SVGFusion results.** The SVGs generated by **SVGFusion** have a clear hierarchical structure, making them easy to edit, such as changing the colors (the top example) or recomposing into a new graphic (the bottom example).

SVGFusion outputs enables users to efficiently repurpose synthesized vector elements and construct novel vector compositions. Notably, SVGs generated by **SVGFusion** can be seamlessly recombined to form entirely new designs, showcasing the flexibility and reusability of the framework.