

Conditional Human Sketch Synthesis with Explicit Abstraction Control

Har-Yen Chen
 chendaryen@outlook.com

¹ SketchX, CVSSP
 University of Surrey, UK

Yi-Zhe Song¹
 y.song@surrey.ac.uk

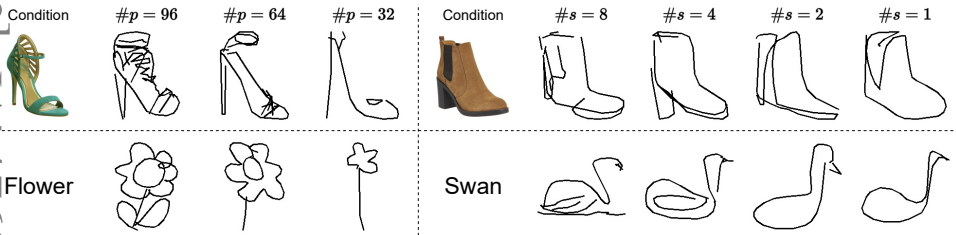


Figure 1: Our approach enables explicit control over the length of points (left) and strokes (right) in both photo-to-sketch synthesis (top) and class-conditional synthesis (bottom). This allows us to create recognizable sketches with few points (4th column) or a single stroke (9th column). Our results maintain key features at various abstraction levels, adjusting the detail amount accordingly. #p, #s: number of points and strokes in the sketch, respectively.

Abstract

This paper presents a novel free-hand sketch synthesis approach addressing explicit abstraction control in class-conditional and photo-to-sketch synthesis. Abstraction is a vital aspect of sketches, as it defines the fundamental distinction between a sketch and an image. Previous works relied on implicit control to achieve different levels of abstraction, leading to inaccurate control and synthesized sketches deviating from human sketches. To resolve this challenge, we propose two novel abstraction control mechanisms, state embeddings and the stroke token, integrated into a transformer-based latent diffusion model (LDM). These mechanisms explicitly provide the required amount of points or strokes to the model, enabling accurate point-level and stroke-level control in synthesized sketches while preserving recognizability. Outperforming state-of-the-art approaches, our method effectively generates diverse, non-rigid and human-like sketches. The proposed approach enables coherent sketch synthesis and excels in representing human habits with desired abstraction levels, highlighting the potential of sketch synthesis for real-world applications.

1 Introduction

From prehistoric cave paintings to modern digital art, sketching has been an innate human ability, encapsulating our experiences and emotions in simple strokes. Sketching encompasses two fundamental aspects: content and style. Content refers to the subject matter of the sketch. Abstraction, ranging from a single stroke to a detailed depiction, lies at the core of the style. This study aims to explore whether neural networks can emulate human sketching capabilities by conditionally synthesizing sketches with a desired level of abstraction.

The field of sketch synthesis has made significant strides since the introduction of Sketch-RNN [6], which marked the beginning of the pursuit to generate sketches in a human-like manner. Despite the progress made in the field, previous works predominantly focus on unconditional synthesis within a single category or a limited number of categories. There has been limited research on class-conditional sketch synthesis models. Photo-to-sketch synthesis, a relatively unexplored area, has witnessed only a few attempts by researchers, including works by Song *et al.* [20], Muhammad *et al.* [14], SketchLattice [17], and CLIPasso [23]. Although some of these works [14, 17, 23] propose different abstraction controlling mechanisms, they suffered from drawbacks such as incomplete sketches, inaccurate abstraction control, and rigid styles.

The limitations of past works come from the inability to learn from data how humans sketch at different abstraction levels. They need to rely on implicit support like lattice representation [17] or pretrained classifiers [14, 23]. Conversely, our explicit control approach offers an accurate and human-like solution for controlling the abstraction level in synthesized sketches. We propose two abstraction control mechanisms, state embeddings and stroke token, which integrate with transformer-based [20] latent diffusion models (LDM) [16]. These mechanisms enable accurate control over the length of points and strokes, respectively. By feeding the number of points or strokes to the diffusion model straightforwardly, our approach can adapt the sketch strategy based on explicit abstraction level constraints, generating sketches with diverse styles that emulate human traits (see Figures 5 and 6). To reduce the computational burden on the diffusion backbone, we pretrain a Variational Autoencoder (VAE) [10] to provide compressed neural representations of sketches, significantly shortening sketches and allowing more efficient leveraging of the transformer. By incorporating adaptive layer normalization and cross-attention conditioning in diffusion models, we seamlessly achieve class-conditional and photo-to-sketch synthesis.

In summary, our contributions include: (i) proposing state embeddings and the stroke token to accurately and naturally control the number of points and strokes in the generated sketches, respectively. (ii) presenting a transformer-based LDM framework that accomplishes class-conditional synthesis and photo-to-sketch synthesis. (iii) being the first model that succeeds in class-conditional sketch synthesis across a wide range of categories. (iv) enhancing the real-world applicability of sketch synthesis by providing controllability of both content and abstraction level.

2 Related Work

Free-Hand Sketch Synthesis Sketch-RNN [6] initiated the study of neural representations of vector sketches by using a sequence-to-sequence Variational Autoencoder (seq2seq-VAE) [10] combined with a recurrent neural network-based autoregressive model. Addressing the issue of long sketch point sequences, Das *et al.* [2] employed Bézier curves to represent

strokes, allowing for fewer points to control smooth strokes compared to point-by-point segments. This approach leads to improved performance in synthesizing longer sketches.

Recent advancements introduced Denoising Diffusion Probabilistic Models (DDPMs) [8] for image synthesis, showcasing their potential in generating high-quality images. SketchKnitter [24] extended DDPMs to free-hand sketch synthesis, demonstrating the framework’s benefits in the sketch domain. However, a significant limitation of these studies is their models’ need for single-category training or their ability to perform only unconditional synthesis on multi-category datasets. Class-conditional sketch synthesis remains unexplored in the literature.

For instance-based synthesis, Song *et al.* [21] employed deep neural networks with shortcut cycle consistent constraints to achieve photo-to-sketch synthesis. Muhammad *et al.* [14] trained a reinforcement learning agent to automatically remove strokes from edge maps. Qi *et al.* [17] proposed a lattice representation of sketches, combining Long Short-Term Memory (LSTM) [9] with graph models to create SketchLattice, capable of generating sketches based on edge maps. Taking a different approach, CLIPasso [23] implemented a step-by-step optimization algorithm for sketch synthesis.

Sketches Abstraction Muhammad *et al.* [14] developed a reinforcement model that employs a pretrained classifier to provide recognizability as a reward, allowing the model to remove unimportant strokes and achieve various abstraction levels. SketchLattice [17] adjusts the abstraction level by controlling the density of lattice points. Another approach, CLIPasso [23], sets the number of strokes during initialization, with the optimization process guided by CLIP [18] similarity, influencing stroke shapes. However, a common issue among these methods is that the resulting sketches do not reflect human habits. It is essential to consider that people alter their sketching styles when applying different abstraction levels, a factor that should be incorporated into future models. Neither classifier, lattice, nor CLIP can express this aspect of human sketching behavior.

Diffusion Probabilistic Models In recent years, DDPMs [9, 8, 15] have gained prominence as an effective approach to image synthesis. DDPMs treat image generation as a gradual Markov denoising process, facilitating the capture of major data variations. Nonetheless, the time-consuming nature of the multi-step denoising process remains a challenge. DDIM [20] addressed this issue by proposing a non-Markov process for faster sampling and inference. Further advancements in DDPM efficiency were made by Rombach *et al.* [19] through the integration of VAEs. By capitalizing on VAEs’ dimensional reduction capabilities, the time consumption of the denoising process during training and synthesis can be significantly reduced. Although DDPMs have been applied to various sequential data types, such as speech [12], their use in sketch synthesis is relatively unexplored, with only Diff-HW [13] and SketchKnitter [24] venturing into this domain, indicating that the field remains in its infancy.

3 Method

In our method, we represent a sketch as a sequence of N points, denoted as $S = (s_1, s_2, \dots, s_N)$. Each point $s_i = (x_i, y_i, p_i)$ is defined by its normalized absolute position, i.e., $x_i, y_i \in [-1, 1]$, and the corresponding binary pen-state p_i . The pen-state is "draw" when $p_i = 0$, and "pen lift" when $p_i = 1$. To maintain a consistent shape for all sketches, we keep N constant across the dataset. We apply the Ramer–Douglas–Peucker algorithm [5] to simplify strokes. For shorter sketches, we pad them to extend their length, while for longer strokes, we truncate the sequence to the first N points.

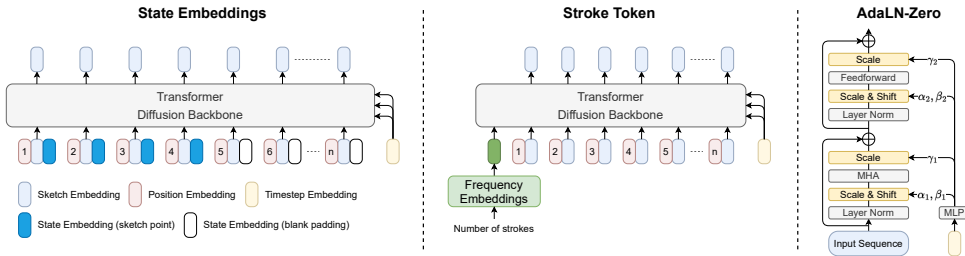


Figure 2: Conditioning mechanism architectures. Left: Our state embeddings represent the state of each sketch token in the input sequence, conveying point length information to the model. Middle: We utilize stroke embedding as an extra token to the input sequence, allowing it to interact with the sketch sequence via MHA. This interaction is critical for effective sketch synthesis with varying stroke counts. Right: Details of adaLN-Zero [16].

3.1 Latent Diffusion Framework

In this study, we employ transformer-based Latent Diffusion Models (LDM) as our synthesis framework. Contrary to traditional Diffusion Models, as used in Diff-HW [13] and SketchKnitter [24], which operate directly on sketch points, our method compresses the temporal dimension while maintaining the semantics of the sketch through a VAE, thus creating a compact low-dimensional latent space. The transformer [22] architecture naturally fits the sequential structure of sketch vectors.

By focusing on the VAE’s latent space, the diffusion backbone effectively reduces the processing of redundant information in high-dimensional spaces, which in turn significantly decreases time and memory consumption during training and inference. Let the encoder and decoder of the first-stage VAE be represented by E and D , and the diffusion backbone by ϵ_θ . LDMs are perceived as sequential denoising autoencoders $\epsilon_\theta(z_t, t), t = 1, \dots, T$, trained to estimate the noise component in z_t and generate a less noisy z_{t-1} . z_t is acquired via a diffusion process applied to $z_0 = E(S)$. This process is characterized as a Markov Chain of length T , with each step involving a minor Gaussian perturbation of the prior state. We will use $\epsilon_\theta(z_t)$ as an abbreviation for the time-dependent $\epsilon_\theta(z_t, t)$.

For conditional synthesis, $\epsilon_\theta(z_t)$ can be reformulated as $\epsilon_\theta(z_t, c)$, given the condition c . With the aid of a reweighted variational lower bound [9], the objective can be formulated as:

$$\mathcal{L} = \mathbb{E}_{z_t, \epsilon \sim \mathcal{N}(0, \mathbf{I}), t, c} [\|\epsilon - \epsilon_\theta(z_t, c)\|_2^2] \quad (1)$$

where $t \sim \mathcal{U}(1, \dots, T)$. The likelihood learning is achieved by minimizing the mean-squared error between the actual noise ϵ and the estimated noise $\epsilon_\theta(z_t, c)$.

We use AdaLN-Zero, proposed by Peebles and Xie [16], for processing timestep embeddings. AdaLN-Zero has proven to be more effective than conventional adaptive layer normalization in transformer blocks for diffusion models, as it regresses dimension-wise scaling parameters γ prior to the residual connection, in addition to regressing scale α and shift β parameters after layer normalization (see Figure 2).

3.2 Conditional Synthesis

Class-conditional synthesis By summing label embeddings with timestep embedding, AdaLN-Zero allows us to efficiently inject category information into the diffusion backbone, achiev-

ing class-conditional sketch synthesis.

Photo-to-Sketch synthesis We leverage the pretrained CLIP (Contrastive Language-Image Pretraining) [18] image encoder as the photo encoder. CLIP is trained using contrastive learning on a wide variety of image domains and linguistic concepts, effectively capturing the semantics in images. We concatenate the classification tokens from all self-attention layers in CLIP ViT-B/32 to construct a length-12 sequence. To exploit CLIP’s capabilities, we attach multi-head cross-attention layers following each multi-head self-attention (MHA) layer in the transformer architecture, which is similar to the conditioning used in Stable Diffusion [19] for text-to-image tasks. These cross-attention layers enable our model to capture the object’s semantics and establish stronger correspondence between the photo and sketch domains, resulting in coherent synthesized sketches that resemble the reference photo.

3.3 Controlling Abstraction

In sketches, the level of abstraction can be determined by two factors: point length and stroke length. Point length refers to the total number of points in a sketch, which directly impacts its complexity and abstraction. Conversely, stroke length pertains to the number of strokes that make up a sketch. While there is no definitive relationship between stroke length and point length, they generally exhibit a positive correlation. For humans, manipulating stroke length is more intuitive than adjusting point length, as strokes naturally form during sketching, whereas points are a byproduct of computer graphics. Prior research by Song *et al.* [20] and CLIPasso [23] demonstrated precise control over stroke count in generated samples, but no existing approach has managed to regulate sketch point length.

In this work, we accurately control both point length and stroke length, enabling more refined manipulation of sketch abstraction levels. The proposed state embeddings and stroke token architectures are illustrated in Figure 2.

Point length We introduce learnable state embeddings inspired by the segmentation embeddings used in BERT [9]. There are two types of state embeddings: one representing sketch points and another representing blank padding. These state embeddings, combined with the input sequence and position embeddings, enable the model to differentiate sketch points from padding.

Stroke length In order to achieve stroke length control in the sketch generation model, we first encode the desired number of strokes using frequency embeddings. The resulting stroke embedding is then concatenated to the input sequence as an additional token, referred to as the stroke token. By incorporating the stroke token, it can interact with the sketch sequence through multi-head attention within the transformer layers. This allows the model to adjust the sketch sequence based on the information encoded in the stroke token, effectively controlling the stroke length in the generated sketches.

4 Experiments

In this section, we provide evaluations in Section 4.1, comparison with SOTA [14, 17, 23] in Sections 4.2 and 4.3, and further ablative experiments in Section 4.4.

Dataset We evaluate our approach on two datasets: the *Quick, Draw!* dataset [6] and the QMUL Shoe-V2 dataset [25]. The *Quick, Draw!* dataset is a large-scale collection of hand-drawn sketches, comprising 345 categories with over 50M sketch vectors. For our experiments, we selected 75 categories to train the first-stage VAE and class-conditional models.

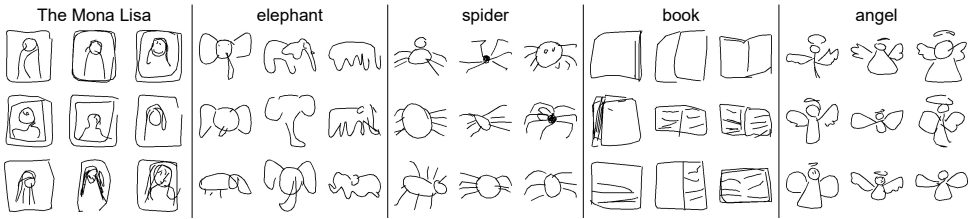


Figure 3: Class-conditional synthesis results showcasing diverse sketches across multiple categories.



Figure 4: Abstraction level controlling on class-conditional synthesis. Left: Point length controlling. Right: Stroke length controlling.

	Class/Point	Class/Stroke	Photo/Point	Photo/Stroke
Acc. (%)	99.3	96.9	99.8	90.6
Acc. ± 1 (%)	99.9	99.8	99.8	100

Table 1: Abstraction controlling accuracy. C/L: the C-conditional model with L length controlling. Acc ± 1 : Accuracy with 1 point/stroke error.

Each category consists of 50K training data and 500 testing data, resulting in a total of 3.75M training data and 37.5K testing data. The QMUL Shoe-V2 dataset contains photo and sketch vector pairs, with 5,982 training data samples and 666 testing data samples. We use the QMUL Shoe-V2 dataset to train and evaluate our photo-to-sketch models.

Implementation Details In this work, all models are trained using a single Nvidia 3080 Ti GPU. We set the number of points, N , to 192. We employ a first-stage VAE with a compression ratio of $4\times$, resulting in a sketch latent representation of length 48 for processing by the diffusion backbone. As the QMUL Shoe-V2 dataset is relatively small, we initialize the overlapping weights of our photo-to-sketch models using class-conditional models and apply Low Rank Adaptation (LORA) [14] to mitigate overfitting. Lastly, all results in this section are sampled with the full 1000 steps.

4.1 Evaluation

Qualitative evaluation Our approach achieves class-conditional sketch synthesis, as demonstrated in Figures 1 and 3, where our model generates sketches corresponding to given categories. By incorporating abstraction conditioning, we can flexibly control the abstraction level by determining either the length of points or strokes, as shown in Figure 8.

As illustrated in Figures 1, our photo-to-sketch models demonstrate the ability to capture the semantics of references with various abstraction levels. Overall, our method exhibits the capability of generating high-quality sketches, whether for category-based or instance-based synthesis, over a wide range of abstractions, showcasing its effectiveness and versatility.



Figure 5: Comparison on photo-to-sketch synthesis. Reference: human sketches in the QMUL Shoe-V2 dataset.

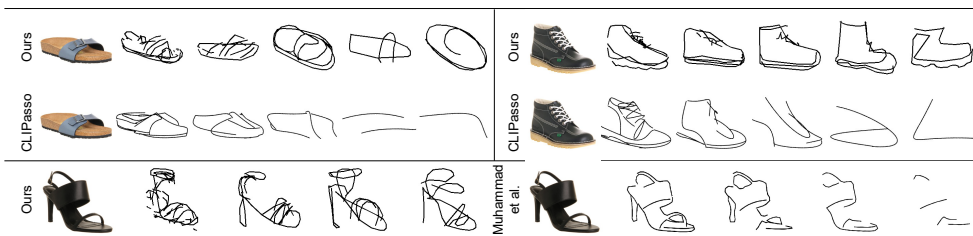


Figure 6: Stroke length controlling comparison. The sketches in the top two rows using 16, 8, 4, 2, and 1 strokes from left to right. We use 32, 16, 8, and 4 strokes in the bottom row.

Quantitative Evaluation To demonstrate that our models can accurately control abstraction levels, we provide the accuracy of the length of points and strokes of our models in Table 1. As shown in the table, our models exhibit high precision for both class-conditional and photo-to-sketch synthesis, and for both point length and stroke length. It can be observed that almost all the results fall within a slight error, further emphasizing the effectiveness of proposed state embeddings and the stroke token in explicitly controlling sketch length.

4.2 Qualitative Comparison

Photo-to-Sketch synthesis In Figure 5, we compare our results with SketchLattice [17] and CLIPasso [23]. Due to the absence of SketchLattice’s implementation, we use the results directly from their paper, hence some results are missing. CLIPasso [23] presents results with close geometric shapes to the images; however, the style is significantly different from the reference sketches drawn by humans. This discrepancy arises because their optimization-based algorithm has no opportunity to learn from real sketch vectors. The CLIP-similarity minimization in CLIPasso leads to results that appear monotonous and lack variety. Instead, our model successfully captures the way humans tend to sketch: less precise and ambiguous yet still presenting the semantics of the object. Our results exhibit different styles, as people have various drawing habits. Our model generates diverse and non-rigid sketches, which effectively capture human sketching characteristics.

Controlling abstraction In Figure 6, we compare our results with two stroke-length controlling works: Muhammad *et al.* [24] and CLIPasso [23]. Since the implementation of Muhammad *et al.* is not available, we use the results directly from their paper. Muhammad

Model	CLIP score \uparrow	Acc@3 \uparrow	Acc@10 \uparrow	Speed \downarrow
CLIPasso #s = 1	45.8	0.03	0.14	90
CLIPasso #s = 4	55.7	0.08	0.21	91
CLIPasso #s = 16	67.9	0.31	0.61	94
Our #s = 1	46.4 / 52.8	0.06 / 0.10	0.18 / 0.29	18 / 20
Our #s = 4	47.8 / 54.0	0.07 / 0.08	0.21 / 0.23	-
Our #s = 16	46.6 / 52.5	0.06 / 0.10	0.15 / 0.26	-
Our #p = 16	43.4 / 47.3	0.06 / 0.06	0.13 / 0.18	17 / 25
Our #p = 32	47.3 / 54.8	0.07 / 0.14	0.17 / 0.28	-
Our #p = 64	50.1 / 56.7	0.10 / 0.13	0.25 / 0.38	-
Human sketch	53.3	0.07	0.20	-

Table 2: Quantitative comparison on photo-to-sketch synthesis. For our models, the second figure is the best result from 10 samples with the highest CLIP score for each object. Acc@k: top-k accuracy of FG-SBIR. Speed: inference time in seconds.

et al. [14] increase abstraction by removing strokes, which can make the sketches appear incomplete. Additionally, the remaining strokes cannot be adjusted accordingly.

CLIPasso [13] can use different arrangements for different abstraction levels. However, the fixed number of control points in Bézier curves limits the expressiveness of strokes, so when there are too few strokes, CLIPasso cannot express the sketches effectively. In contrast, our model successfully generates recognizable sketches even with only one stroke. In the QMUL Shoe V2 dataset, a sketch has at most 15 strokes. Therefore, when using too many strokes (i.e., out of the training domain), our results will present some noisy strokes, as shown in the leftmost image in the bottom row. Despite this limitation, our method is flexible and suitable for generating sketches with different levels of abstraction while maintaining recognizability.

4.3 Quantitative Comparison

We performed a quantitative comparison between our model and CLIPasso using the CLIP score and zero-shot CLIP fine-grained sketch-based image retrieval (FG-SBIR) as metrics. To ensure a fair comparison, we used the CLIP ViT-L/14 model, which was not employed in either of the methods, to calculate the similarity between generated sketches and corresponding object photos. The evaluation was conducted on 100 objects from the QMUL Shoe-V2 testing dataset. Table 2 presents the CLIP scores, retrieval accuracies, and inference times for both approaches.

CLIPasso’s performance is heavily impacted by the number of strokes, with its CLIP score and retrieval accuracy dropping substantially as fewer strokes are used. Notice that the Bézier curve used in CLIPasso makes the number of strokes directly affect the overall length. On the other hand, our model performs consistently regardless of the number of strokes used. Our point-length controlling model demonstrates a correlation between the number of points used and both the CLIP score and retrieval accuracy, supporting the claim that point count is a more direct influence on abstraction than stroke count.

Moreover, our model’s FG-SBIR accuracy is close to the baseline provided by human sketches. When taking the highest score of 10 samples for each object, our model’s performance closely resembles the baseline CLIP score, suggesting its ability to adapt its strategy to more closely align with human sketching patterns across varied points and strokes.

4.4 Ablation Study

In our approach, we take the point sequence of sketches as input, causing strokes is a more challenging concept than points. This prompts us to conduct an ablation study on stroke length controlling, as shown in Table 3. We compare our stroke token approach with two other mechanisms: (1) summing up with timestep embedding and using adaLN-Zero, and (2) forming a length-one sequence and applying additional cross-attention layers. The cross-attention conditioning has the most extra parameters but performs the worst. In contrast, both stroke token and adaLN-Zero are parameter-efficient. The performance of the proposed stroke token is slightly better in terms of both accuracy and FID [17]. This highlights the beneficial interaction that MHA provides for the sketch sequence and stroke token for enhanced control. The stroke token not only constrains the synthesis process unilaterally but can also obtain information from the sketch sequence. Beside accurate control, the lower FID indicates that this information exchange, which adaLN-zero can not offer, might be helpful for generating high-quality sketches.

Type	Acc. (%) \uparrow	Acc. ± 1 (%) \uparrow	FID \downarrow	Extra Params \downarrow
stroke token	96.9	99.8	6.69	886K
adaLN-Zero	96.2	99.8	7.02	1.5M
cross-attention	8.7	25.4	7.55	30M

Table 3: Ablative results on conditioning type of stroke controlling on class-conditional synthesis. The proposed stroke token demonstrates superior control over stroke length while using fewer extra parameters.

5 Limitation

Despite demonstrating promising results, our method encounters some limitations. First, we have not thoroughly studied the first-stage VAE, which is functional but restricts our method’s overall performance. For example, the current VAE may struggle to capture fine-grained details in more intricate sketches. Developing a more compact and semantically rich sketch latent representation might lead to enhanced outcomes and efficiency with longer, complex sketches. Second, due to a limited dataset, our photo-to-sketch synthesis is constrained to shoes. As collecting free-hand sketches is challenging, exploring alternative learning approaches, such as few-shot meta-learning or leveraging synthetic data, could be a valuable solution to expand this task to a broader domain.

6 Conclusion

In this work, we have for the first time successfully developed a neural network-based method to explicitly control the length of points and strokes in synthesized free-hand sketches, demonstrating diverse natural styles and abstraction degrees akin to human sketches. The introduced LDM framework accomplishes both class-conditional synthesis and photo-to-sketch synthesis. To the best of our knowledge, this is the first instance of a class-conditional sketch synthesis model across a broad range of categories. By providing more flexible control, our approach enhances the real-world applicability of sketch synthesis.

References

- [1] Hareesh Bahuleyan, Lili Mou, Olga Vechtomova, and Pascal Poupart. Variational attention for sequence-to-sequence models. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1672–1682, Santa Fe, New Mexico, USA, August 2018. Association for Computational Linguistics. URL <https://aclanthology.org/C18-1142>.
- [2] Ayan Das, Yongxin Yang, Timothy Hospedales, Tao Xiang, and Yi-Zhe Song. Béziersketch: A generative model for scalable vector sketches, 2020. URL <https://arxiv.org/abs/2007.02190>.
- [3] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- [4] Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat GANs on image synthesis. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=AAWuCVzaVt>.
- [5] David H. Douglas and Thomas K. Peucker. Algorithms for the reduction of the number of points required to represent a digitized line or its caricature. *Cartographica: The International Journal for Geographic Information and Geovisualization*, 10:112–122, 1973.
- [6] David Ha and Douglas Eck. A neural representation of sketch drawings, 2017. URL <https://arxiv.org/abs/1704.03477>.
- [7] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/8ald694707eb0fef65871369074926d-Paper.pdf.
- [8] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/4c5bcfec8584af0d967f1ab10179ca4b-Paper.pdf>.
- [9] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.

- [10] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *CoRR*, abs/2106.09685, 2021. URL <https://arxiv.org/abs/2106.09685>.
- [11] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2013. URL <https://arxiv.org/abs/1312.6114>.
- [12] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Dif-fwave: A versatile diffusion model for audio synthesis. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=a-xFK8Ymz5J>.
- [13] Troy Luhman and Eric Luhman. Diffusion models for handwriting generation. *CoRR*, abs/2011.06704, 2020. URL <https://arxiv.org/abs/2011.06704>.
- [14] Umar Riaz Muhammad, Yongxin Yang, Yi-Zhe Song, Tao Xiang, and Timothy M. Hospedales. Learning deep sketch abstraction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [15] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models, 2021. URL <https://openreview.net/forum?id=-NEXDKk8gZ>.
- [16] William Peebles and Saining Xie. Scalable diffusion models with transformers. *arXiv preprint arXiv:2212.09748*, 2022.
- [17] Yonggang Qi, Guoyao Su, Pinaki Nath Chowdhury, Mingkan Li, and Yi-Zhe Song. Sketchlattice: Latticed representation for sketch manipulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 953–961, October 2021.
- [18] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. *CoRR*, abs/2103.00020, 2021. URL <https://arxiv.org/abs/2103.00020>.
- [19] 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, June 2022.
- [20] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv:2010.02502*, October 2020. URL <https://arxiv.org/abs/2010.02502>.
- [21] Jifei Song, Kaiyue Pang, Yi-Zhe Song, Tao Xiang, and Timothy M. Hospedales. Learning to sketch with shortcut cycle consistency. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

- [22] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>.
- [23] 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, 2022.
- [24] Qiang Wang, Haoge Deng, Yonggang Qi, Da Li, and Yi-Zhe Song. Sketchknitter: Vectorized sketch generation with diffusion models. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=4eJ43EN2g6l>.
- [25] Qian Yu, Feng Liu, Yi-Zhe Song, Tao Xiang, Timothy M. Hospedales, and Chen-Change Loy. Sketch me that shoe. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.

Supplementary Material

A First Stage VAE

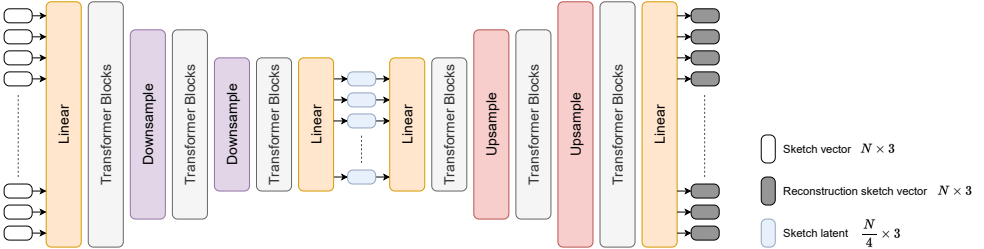


Figure 7: Architecture of first-stage VAE.

In this section, we demonstrate the architecture and training objective of the first-stage VAE, designed to compress the sketch vector’s length dimension. As Figure 7 illustrates, the transformer-based VAE’s encoder is composed of a series of transformer blocks interspersed with linear downsampling layers. Mirroring this architecture, the decoder is constructed of a stack of transformer blocks interspersed with linear upsampling layers. Utilizing a compression ratio of 4, the sketch latent representation has a reduced dimension of $\frac{N}{4} \times 3$, where the original dimension of the sketch vector is $N \times 3$.

Our losses include L_2 loss for absolute positions (x_i, y_i) , binary cross-entropy loss for binary pen states p_i , and Kullback-Leibler (KL) loss for latent z regularization. Denote $(\hat{x}_i, \hat{y}_i, \hat{p}_i)$ as the reconstruction of (x_i, y_i, p_i) from the VAE. The combination objective is as follows:

$$\mathcal{L}_{abs} = \mathbb{E} [\| (x_i, y_i) - (\hat{x}_i, \hat{y}_i) \|_2^2] \quad (2)$$

$$\mathcal{L}_{state} = -\mathbb{E} [(p_i \log(\hat{p}_i) + (1 - p_i) \log(1 - \hat{p}_i))] \quad (3)$$

$$\mathcal{L}_{KL} = \frac{1}{2} \sum (1 + \log(\sigma_z^2) - \mu_z^2 - \sigma_z^2) \quad (4)$$

$$\mathcal{L}_{VAE} = \mathcal{L}_{abs} + \mathcal{L}_{state} + \lambda \mathcal{L}_{KL} \quad (5)$$

The low-weighted coefficient λ is set to 10^{-6} in this work.

B Additional Qualatative Results

We provide additional length controlling results of class-conditional models and photo-to-sketch models in Figure 8 and 9, respectively.

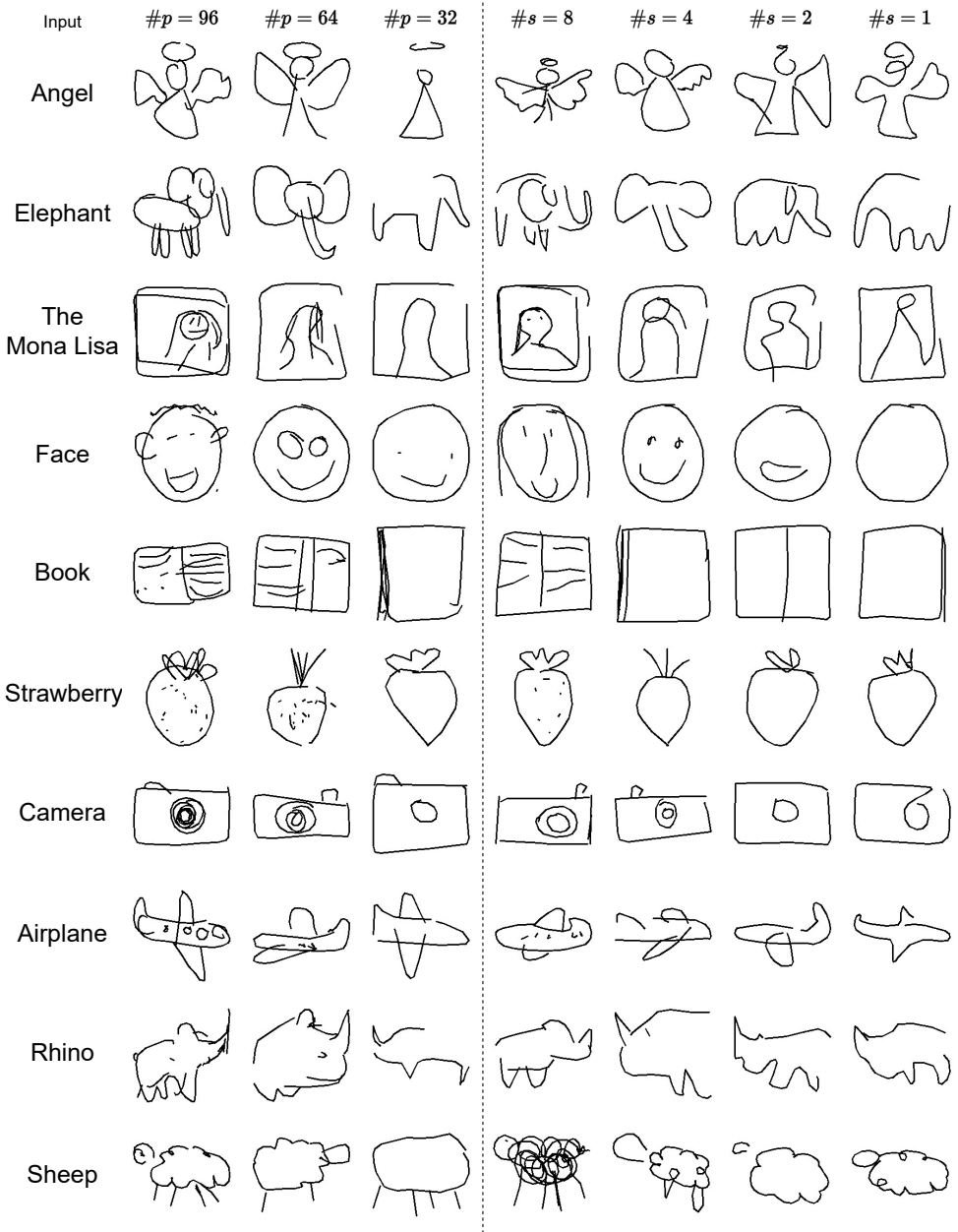


Figure 8: Additional length control results of class-conditional models.



Figure 9: Additional length control results of photo-to-sketch models.