

# SketchAnimator: Animate Sketch via Motion Customization of Text-to-Video Diffusion Models

Ruolin Yang<sup>1</sup>, Da Li<sup>2</sup>, Honggang Zhang<sup>1</sup> and Yi-Zhe Song<sup>2</sup>

<sup>1</sup>PRIS, School of Artificial Intelligence, Beijing University of Posts and Telecommunications, China

<sup>2</sup>SketchX, CVSSP, University of Surrey, United Kingdom

**Abstract**—Sketching is a uniquely human tool for expressing ideas and creativity. The animation of sketches infuses life into these static drawings, opening a new dimension for designers. Animating sketches is a time-consuming process that demands professional skills and extensive experience, often proving daunting for amateurs. In this paper, we propose a novel sketch animation model *SketchAnimator*, which enables adding *creative* motion to a given sketch, like “a jumping car”. Namely, given an input sketch and a reference video, we divide the sketch animation into three stages: Appearance Learning, Motion Learning and Video Prior Distillation. In stages 1 and 2, we utilize LoRA to integrate sketch appearance information and motion dynamics from the reference video into the pre-trained T2V model. In the third stage, we utilize Score Distillation Sampling (SDS) to update the parameters of the Bézier curves in each sketch frame according to the acquired motion information. Consequently, our model produces a sketch video that not only retains the original appearance of the sketch but also mirrors the dynamic movements of the reference video. We compare our method with alternative approaches and demonstrate that it generates the desired sketch video under the challenge of one-shot motion customization.

**Index Terms**—Sketch animation, diffusion process, generative model, video generation, motion extraction

## I. INTRODUCTION

Free-hand sketching is an effective tool for showcasing unique creativity through abstract expression. Dynamic sketches, in particular, have a broader range of applications (e.g., animated advertisements, educational videos and storyboarding). Users often express a desire to create imaginative concepts that defy real-world constraints using sketches, such as “a running clock”. Recently, personalized and customized video generation has been explored [1]–[3]. The objective of this task is to adapt pre-trained foundational models to generate videos depicting a specific motion concept, utilizing reference videos as guidance. However, due to the significant disparity between sketches and the data used to train large-scale generative models, research on customized and personalized dynamic sketch videos is still relatively limited. Namely, directly applying large-scale pre-trained T2V diffusion models presents challenges in terms of generalization to sketch domains. In this work, given a user sketch and a reference video sequence depicting a moving object, our framework generates a video in which the sketch is animated according to the driving sequence without changing its appearance.

We observe three key issues in previous customized sketch animation methods: (i) rely on user annotations: Su *et al.* [4] developed an interactive system where users input a motion-

indicating video sequence and mark key points on sketches to generate dynamic animations. However, this method depends heavily on accurate marking and remains challenging and time-consuming; (ii) fail to decouple motion signals apart from appearance: traditional methods [5]–[8] focus on specific domains such as faces, human figures, and related subjects. They train on datasets that gather numerous similar actions and subjects, relying on specific key-points. This limits customization in video generation. Video quality declines when objects in the driving video and source image vary in domain or shape subtly [9]; (iii) lack of creativity: current state-of-the-art method, live sketch [10], animates a single-subject sketch using motion priors from a large pre-trained T2V diffusion model. However, this model may struggle with specific motion concepts and novel subject-motion combinations unseen in real life. For example, generating a video of “A jumping tree” could challenge producing realistic and consistent results aligned with the description.

To mitigate these issues, we propose SketchAnimator, a novel approach for animating a user-provided sketch with a specific motion extracted from a reference video, *i.e.* one-shot customized sketch video generation. To avoid human annotations and subtle shape problems, we employ Low-Rank Adaptations [11] to fine-tune pre-trained T2V models on sketches and videos. This allows us to extract the motion information embedded within the video sequence and the subject-specific details inherent in the sketch in a staged manner. To enhance diversity and creativity, we combine identity features with motion patterns to generate a video prior capable of predicting the content of the target sketch animation. However, direct inference faces challenges in preserving both the structural integrity of the sketch and the fidelity of motion information. Therefore, building on prior research [12]–[14], we employ a differentiable Bézier curves rendering to achieve adaptable sketch representation. Additionally, we show how to leverage a score distillation sampling (SDS) loss [15] to guide the sketch generation while staying aligned with the customized video prior. We show superior results compared with the traditional motion transfer method and advanced pixel-based diffusion generation models. Moreover, our work allows users to animate sketches with motions that are not present in the training set of large-scale pre-trained video diffusion models.

## II. METHOD

In this section, we begin with preliminaries. Next, we outline the details of SketchAnimator, which aim to address

three key challenges: 1) how to preserve the characteristics of the original sketch; 2) how to extract the motion pattern from the reference video; and 3) how to add motion information to the sketch without changing its appearance in a one-shot manner:

#### A. Preliminaries

**Video Diffusion Models** Similar to image diffusion models, video diffusion models (VDM) [16], [17] learn a distribution representing video data through denoising process. During training, an autoencoder map the video sample  $x$  of length  $F$  into a latent representation  $z_0 \in \mathbf{R}^{F \times h \times w \times c}$ . Next, the forward process gradually adds Gaussian noise  $\epsilon \sim \mathcal{N}(0, I)$  to the latent code to obtain a noised data, as shown as follows:

$$z_t = \alpha_t z_0 + \sigma_t \epsilon, \quad (1)$$

where  $\alpha_t, \sigma_t$  are two predetermined schedule [18] of random time step  $t \sim \mathcal{U}(0, T)$ . Video diffusion model is conditioned on text prompts  $y$  and adopt a 3D U-Net  $\epsilon_\theta$  to perform denoising. Let  $\tau_\theta(y)$  be the text encoder. The following reweighted variational bound is employed for optimization:

$$\mathcal{L} = \mathbf{E}_{z_0, y, t, \epsilon} \left[ \|\epsilon - \epsilon_W(z_t, \tau_\theta(y), t)\|_2^2 \right], \quad (2)$$

**Low-Rank Adaptation** Low-Rank Adaptation (LoRA) [11] method enables efficient adaptation of large pre-trained language models to downstream tasks. It predicts  $\Delta W$  to update the pre-trained weight matrix  $W_0 \in \mathbf{R}^{d \times k}$ . Specifically,  $\Delta W$  are decomposed by weight  $B \in \mathbf{R}^{d \times r}$  and  $A \in \mathbf{R}^{r \times k}$ . After fine-tuning,  $\Delta W$  can be merged into  $W_0$  as a plug-and-play module which altered the direction of the model's prediction as:

$$W = W_0 + \alpha \Delta W = W_0 + \alpha B A, \quad (3)$$

where  $d$  is the input dimension and  $k$  is the output dimension with  $r \ll \min(d, k)$ .

**Score-Distillation Sampling** DreamFusion [15] first proposed the score-distillation sampling (SDS) loss for optimizing pre-trained diffusion models to guided the 3D generation process. Given an arbitrary differentiable parametric function that renders images  $x$  (e.g., a NeRF),  $g_\phi$ , the gradient of the diffusion loss function with respect to the parameters  $\phi$  specified as SDS loss is given by:

$$\nabla_\phi \mathcal{L}_{SDS} = \left[ w(t)(\epsilon_W(z_t, \tau_\theta(y)), t) \frac{\partial x}{\partial \phi} \right], \quad (4)$$

where  $w(t)$  is a weighting function.

In our paper, sketch is rendering by a differentiable rasterizer [19] at stroke level following previous sketch generation works [12]–[14]. Hence, we utilize SDS loss to generate sketch videos via optimization from a customized diffusion model.

#### B. SketchAnimator

Fig. 1 illustrates our three-stage learning pipeline, which takes a user-provided static sketch in vector format and a driving video as input and produces the customized sketch video. We decouple this task into appearance learning, motion learning and video prior distillation.

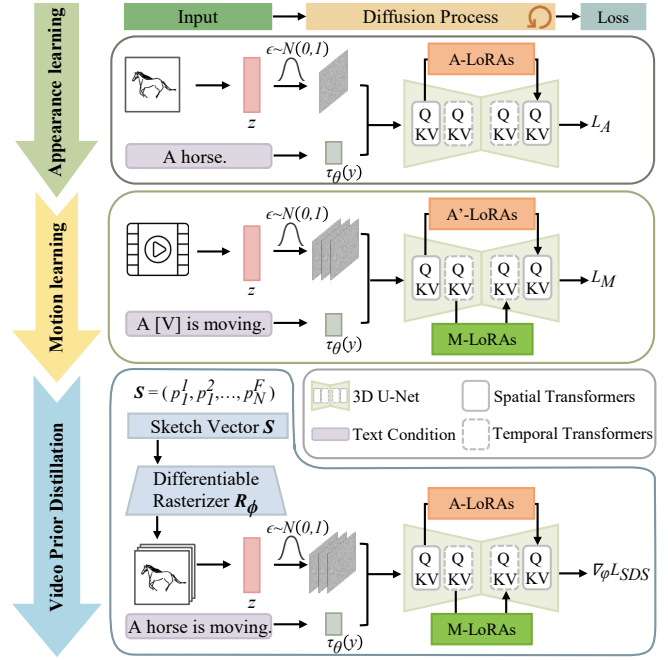


Fig. 1. Illustration of the proposed SketchAnimator pipeline. In appearance learning stage, the trainable A-LoRAs are applied to absorb the visual information of sketch. In Motion learning stage, we separate motion information using A'-LoRAs and M-LoRAs to reconstruct the video. During Video Prior Distillation stage, we combine two LoRAs and update the parameters of the Bézier curves in each sketch frame by SDS loss.

**Appearance Learning** The input sketch is firstly represented in the image format. Inspired by Dreambooth [20] and Textual Inversion [21], we learn the appearance information of sketch by optimizing the attention layers with trainable LoRAs of the base T2V model, denoted as Appearance LoRAs (A-LoRAs). Also, we find that representing the sketch as a single subject by a semantic word (e.g., “A horse.”), denoted by  $y_a$ , serves as a better textual queries than the format of “A photo of \*” or “A sketch drawing of \*.” Given that in this scenario, the input sketch image exhibits no temporal variations, we solely optimize spatial transformer layers in the 3D U-Net blocks following the objectives:

$$\mathcal{L}_{appearance} = \mathbf{E}_{z_0, y_a, t, \epsilon} \left[ \|\epsilon - \epsilon_W(z_t, \tau_\theta(y_a), t)\|_2^2 \right]. \quad (5)$$

**Motion Learning** To obtain the motion signal from the video, we decouple and reconstruct the video in both the temporal and spatial dimensions. As shown in Fig. 1, given a reference video and corresponding caption, we add LoRAs to attention blocks in spatial transformer and temporal transformer separately. For each training step, the A'-LoRAs are trained on a single frame randomly sampled from the training video to fit its appearance while ignoring its motion, based on spatial loss, which is reformulated as

$$\mathcal{L}_{spatial} = \mathbf{E}_{z_0, y_m, t, \epsilon, i \sim \mathcal{U}(0, F)} \left[ \|\epsilon - \epsilon_\theta(z_{t,i}, \tau_\theta(y_m), t)\|_2^2 \right], \quad (6)$$

where  $F$  is the number of frames of the training data and the  $z_t, i$  is the sampled frame from the latent code  $z_t$  and the  $y_m$  is the prompt depicting the video sequence.

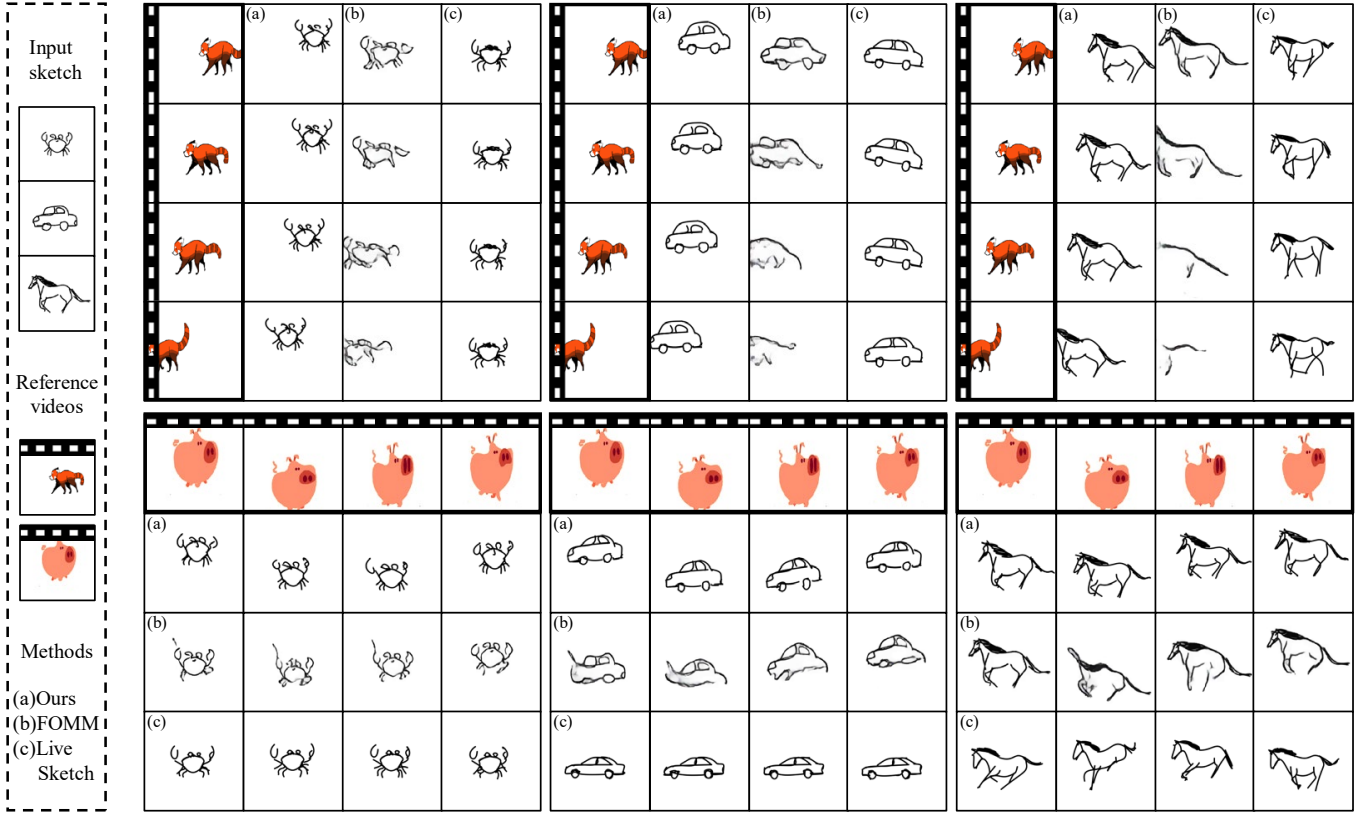


Fig. 2. Qualitative comparisons with previous methods. FOMM captures good motion information but may cause deformation in the target sketch. Sketch videos generated using Live Sketch tend to be static or exhibit slight deformations. However, our approach effectively encodes motion information while avoiding changes in appearance.

For Motion LoRAs (M-LoRAs), the loss function is following the VDM loss as (2). Therefore, the final objective of our motion learning process is the combination of temporal loss and spatial loss as follows:

$$\mathcal{L}_{motion} = \mathcal{L}_{spatial} + \mathcal{L}_{temporal}. \quad (7)$$

**Video Prior Distillation** Initially, the SVG sketch is consist of  $N$  strokes for the first frame  $P = \{p_1, \dots, p_N\} \in \mathbf{R}^{N \times 4 \times 2}$ . Each stroke is represented by a two-dimensional Bézier curve with four control point  $p_i = \{p_i^j\}_{j=1}^4 = \{(x_i, y_i)^j\}_{j=1}^4 \in \mathbf{R}^{4 \times 2}$ . Next, in accordance with the frame count of the reference video, the strokes are duplicated  $F$  times for each frame initialization, which denoted by  $S = \{P^f\}_{f=1}^F \in \mathbf{R}^{N \cdot F \times 4 \times 2}$ . Then, all the strokes are fed into a differentiable rasterizer,  $R\phi$ , to acquire raster frames. Therefore, the input sequence in the latent space is noted as:

$$z_0 = E_\theta(R_\phi(S)), \quad (8)$$

where  $E_\theta$  is the pre-trained variational autoencoder.

Subsequently, we inject the A-LoRAs of sketch  $\Delta W_A$  and M-LoRAs of reference video  $\Delta W_M$  into the pre-trained video diffusion model as expected video prior, the final weight matrix of VDM is:

$$W' = W_0 + \lambda_1 \Delta W_A + \lambda_2 \Delta W_M, \quad (9)$$

and then the Bézier curve  $S$  is updated using SDS loss until the structure, position and scale of sketch is changed as desired and we get final sketch curve  $S'$  following:

$$\nabla_\phi \mathcal{L}_{SDS} = \left[ w(t)(\epsilon_{W'}(z_t, \tau_\theta(y)), t) \frac{\partial x}{\partial \phi} \right], \quad (10)$$

$$\phi' \xleftarrow{\nabla_\phi \mathcal{L}_{SDS}} \phi, \quad (11)$$

$$S' \xleftarrow{\phi'} S, \quad (12)$$

where  $w(t)$  is a weighting function. The Eq.(11) and (12) are the brief backpropagation process for bridging the raster and vector domains by the differentiable rasterizer [19].

### III. EXPERIMENTS

#### A. Experimental Setup

**Datasets** The proposed method is evaluated on the popular motion transfer datasets MGIF dataset [22]. For each video, we pick 5 sketches from CLIPasso [12], QuickDraw [23] and SketchVOS [24]. In our experiments, all sketches are pre-processed into vector format and all the frames are pre-processed into  $256 \times 256$  pixels.

**Implementation details** We evaluate our approach primarily on Modelscope [16], a standard video diffusion model. We employ the Adam optimizer to train all the parameters on a single RTX 3090 GPU. For sketch appearance learning

TABLE I  
QUANTITATIVE EVALUATIONS AND ABLATION STUDIES

Method	Appearance Alignment ( $\uparrow$ )	Motion Alignment ( $\uparrow$ )	Temporal Consistency ( $\uparrow$ )
FOMM	0.824	0.209	0.934
Custom-A-Video	0.689	0.314	0.879
MotionDirector	0.729	0.398	0.942
DreamVideo	0.743	0.407	0.951
Live Sketch	0.948	0.460	0.980
Ours	<b>0.955</b>	<b>0.541</b>	<b>0.988</b>
Ablation results			
w/o A-LoRAs	0.947	0.512	0.995
w/o M-LoRAs	0.952	0.457	0.981

and video motion learning, the LoRAs are optimized by 500 iterations following [20]. During the video prior distillation stage, the differentiable rasterizer are updated with learning rate at  $2.0 \times 10^{-3}$ . The hyperparameters  $\lambda_1$  and  $\lambda_2$  of the appearance LoRAs and motion LoRAs are set to 0.5 and 1.

**Evaluation metrics** Assessing the quality of sketch animation poses a challenge due to the ground truth animations are not available. We evaluate our approach with the following three metrics: (1) *Appearance Alignment* measures the average cosine similarity between CLIP [25] image embeddings of all generated frames and the input sketch. (2) *Motion Alignment* computes the generated video frames and the inference prompt depicting reference video motion, in the form of X-CLIP Score [26], a metric for general video recognition. (3) *Temporal Consistency* calculates CLIP image embeddings for all generated frames and present the average cosine similarity across all pairs of consecutive frames.

### B. Main results

Table I and Fig. 2 present comparisons with cutting-edge methods. Our method achieves the best results with a better trade-off between the visual quality and motion alignments among all methods. In particular, FOMM [5] leads to sketch deformation in details, especially when the object in the reference video differs significantly from the sketch in both shape and semantics. On the contrary, our approach disentangles a cleaner motion signal, leading to videos characterized by abundant variety. Live sketch [10] is capable of generating high-resolution images. However, the preservation of original structure is far from accurate and the animation results are closed to static as shown in 8th row. In contrast, our method is able to generate more creative sketch video (e.g., combining “car” with “jumping”). State-of-the-art customized video generation methods [1]–[3] perform worse in both appearance preservation and motion translation correctness.

### C. Ablation study

We conduct an ablation study to demonstrate the necessity of each component. Specifically, we test two settings on the generator and report the results in the Table I: (1) w/o Appearance LoRAs (A-LoRAs), we observe a decrease in Appearance Alignment scores and the structure of the sketch undergoes deformation shown in the third row of Fig. 3. This suggests the

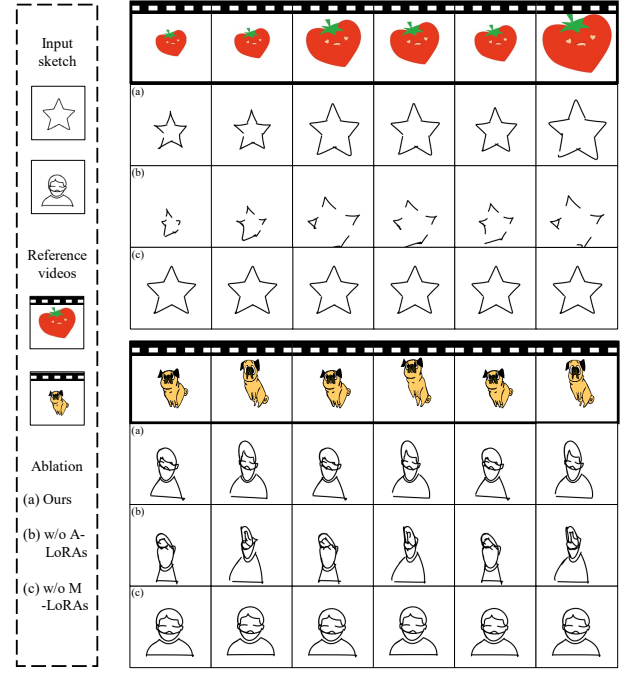


Fig. 3. Visualization of ablation study. SketchAnimator faithfully retains both the identity of the subject and the pattern of its motion.

effectiveness of learning sketch information. (2) w/o Motion LoRAs (M-LoRAs), we notice a significant decline in motion alignment scores and the sketch has transitioned to a static state illustrated in the fourth row of Fig. 3. This indicates that the M-LoRAs play a crucial role in predicting motion corresponding to the reference video.

### D. Discussion

We further illustrate the results of advanced Image-to-Video generation models in Fig. 4 including Gen2 [27], DynamiCrafter [28], SVD [29], MotionDirector [1], Custom-A-Video [3] and DreamVideo [2]. All these methods show redundant information, and commonly fail to produce a sketch. In contrast, our method has significant advantages in these areas.

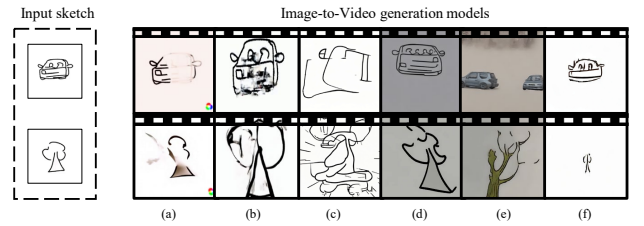


Fig. 4. The results of Image-to-Video generation models. (a) Gen2, (b) SVD, (c) DynamiCrafter, (d) MotionDirector, (e) Custom-A-Video and (f) DreamVideo.

## IV. CONCLUSION

In this paper, we explore dynamic sketch generation customized by a reference video and construct a multi-stage generation framework for better solving the challenges mentioned above. Extensive experiments demonstrate that our approach achieved a high-quality and better motion alignment sketch animation generation compared to previous methods. This

research opens up new possibilities for generating highly customized and creative animations tailored to individual users.

## REFERENCES

- [1] R. Zhao, Y. Gu, J. Z. Wu, D. J. Zhang, J. Liu, W. Wu, J. Keppo, and M. Z. Shou, "Motiondirector: Motion customization of text-to-video diffusion models," *arXiv preprint arXiv:2310.08465*, 2023.
- [2] Y. Wei, S. Zhang, Z. Qing, H. Yuan, Z. Liu, Y. Liu, Y. Zhang, J. Zhou, and H. Shan, "Dreamvideo: Composing your dream videos with customized subject and motion," in *Proc. IEEE/CVF Conf. Comput. Vis. and Pattern Recognit.*, 2024.
- [3] Y. Ren, Y. Zhou, J. Yang, J. Shi, D. Liu, F. Liu, M. Kwon, and A. Srivastava, "Customize-a-video: One-shot motion customization of text-to-video diffusion models," *arXiv preprint arXiv:2402.14780*, 2024.
- [4] Q. Su, X. Bai, H. Fu, C.-L. Tai, and J. Wang, "Live sketch: Video-driven dynamic deformation of static drawings," in *Proc. 2018 CHI Conf. Human Factors Comput. Syst.*, pp. 1–12, 2018.
- [5] A. Siarohin, S. Lathuilière, S. Tulyakov, E. Ricci, and N. Sebe, "First order motion model for image animation," in *Proc. Adv. Neural Inf. Proces. Syst.*, vol. 32, 2019.
- [6] A. Siarohin, O. J. Woodford, J. Ren, M. Chai, and S. Tulyakov, "Motion representations for articulated animation," in *Proc. IEEE/CVF Conf. Comput. Vis. and Pattern Recognit.*, pp. 13653–13662, 2021.
- [7] J. Zhao and H. Zhang, "Thin-plate spline motion model for image animation," in *Proc. IEEE/CVF Conf. Comput. Vis. and Pattern Recognit.*, pp. 3657–3666, 2022.
- [8] F.-T. Hong, L. Zhang, L. Shen, and D. Xu, "Depth-aware generative adversarial network for talking head video generation," in *Proc. IEEE/CVF Conf. Comput. Vis. and Pattern Recognit.*, pp. 3397–3406, 2022.
- [9] B. Xu, B. Wang, J. Deng, J. Tao, T. Ge, Y. Jiang, W. Li, and L. Duan, "Motion and appearance adaptation for cross-domain motion transfer," in *Proc. Eur. Conf. Comput. Vis.*, pp. 529–545, Springer, 2022.
- [10] R. Gal, Y. Vinker, Y. Alaluf, A. Bermano, D. Cohen-Or, A. Shamir, and G. Chechik, "Breathing life into sketches using text-to-video priors," in *Proc. IEEE/CVF Conf. Comput. Vis. and Pattern Recognit.*, pp. 4325–4336, 2024.
- [11] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," *arXiv preprint arXiv:2106.09685*, 2021.
- [12] Y. Vinker, E. Pajouheshgar, J. Y. Bo, R. C. Bachmann, A. H. Bermano, D. Cohen-Or, A. Zamir, and A. Shamir, "Clipasso: Semantically-aware object sketching," *ACM Transact. Graph. (TOG)*, vol. 41, no. 4, pp. 1–11, 2022.
- [13] Y. Vinker, Y. Alaluf, D. Cohen-Or, and A. Shamir, "Clipascene: Scene sketching with different types and levels of abstraction," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, pp. 4146–4156, 2023.
- [14] X. Xing, C. Wang, H. Zhou, J. Zhang, Q. Yu, and D. Xu, "Diffsketcher: Text guided vector sketch synthesis through latent diffusion models," in *Proc. Adv. Neural Inf. Proces. Syst.*, vol. 36, 2024.
- [15] B. Poole, A. Jain, J. T. Barron, and B. Mildenhall, "Dreamfusion: Text-to-3d using 2d diffusion," *arXiv preprint arXiv:2209.14988*, 2022.
- [16] J. Wang, H. Yuan, D. Chen, Y. Zhang, X. Wang, and S. Zhang, "Modelscope text-to-video technical report," *arXiv preprint arXiv:2308.06571*, 2023.
- [17] H. Chen, M. Xia, Y. He, Y. Zhang, X. Cun, S. Yang, J. Xing, Y. Liu, Q. Chen, X. Wang, *et al.*, "Videocrafter1: Open diffusion models for high-quality video generation," *arXiv preprint arXiv:2310.19512*, 2023.
- [18] P. Dhariwal and A. Nichol, "Diffusion models beat gans on image synthesis," in *Proc. Adv. Neural Inf. Proces. Syst.*, vol. 34, pp. 8780–8794, 2021.
- [19] T.-M. Li, M. Lukáč, M. Gharbi, and J. Ragan-Kelley, "Differentiable vector graphics rasterization for editing and learning," *ACM Transact. Graph. (TOG)*, vol. 39, no. 6, pp. 1–15, 2020.
- [20] N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein, and K. Aberman, "Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation," in *Proc. Adv. Neural Inf. Proces. Syst.*, pp. 22500–22510, 2023.
- [21] R. Gal, Y. Alaluf, Y. Atzmon, O. Patashnik, A. H. Bermano, G. Chechik, and D. Cohen-or, "An image is worth one word: Personalizing text-to-image generation using textual inversion," in *Proc. Eleventh Int. Conf. Learning Represent.*, 2023.
- [22] A. Siarohin, S. Lathuilière, S. Tulyakov, E. Ricci, and N. Sebe, "Animating arbitrary objects via deep motion transfer," in *Proc. IEEE/CVF Conf. Comput. Vis. and Pattern Recognit.*, pp. 2377–2386, 2019.
- [23] S. Cheema, S. Gulwani, and J. LaViola, "Quickdraw: improving drawing experience for geometric diagrams," in *Proc. SIGCHI Conf. Human Factors Comput. Syst.*, pp. 1037–1064, 2012.
- [24] R. Yang, D. Li, C. Hu, T. Hospedales, H. Zhang, and Y.-Z. Song, "Sketch-based video object segmentation: Benchmark and analysis," in *Proc. 34th British Machine Vis. Conf.*, 2023.
- [25] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, *et al.*, "Learning transferable visual models from natural language supervision," in *Proc. Int. Conf. machine learning*, pp. 8748–8763, 2021.
- [26] B. Ni, H. Peng, M. Chen, S. Zhang, G. Meng, J. Fu, S. Xiang, and H. Ling, "Expanding language-image pretrained models for general video recognition," in *Proc. Eur. Conf. Comput. Vis.*, pp. 1–18, 2022.
- [27] "Runway. gen-2: Text driven video generation." <https://research.runwayml.com/gen2>, 2023.
- [28] J. Xing, M. Xia, Y. Zhang, H. Chen, X. Wang, T.-T. Wong, and Y. Shan, "Dynamicrafter: Animating open-domain images with video diffusion priors," *arXiv preprint arXiv:2310.12190*, 2023.
- [29] A. Blattmann, T. Dockhorn, S. Kulal, D. Mendelevitch, M. Kilian, D. Lorenz, Y. Levi, Z. English, V. Voleti, A. Letts, *et al.*, "Stable video diffusion: Scaling latent video diffusion models to large datasets," *arXiv preprint arXiv:2311.15127*, 2023.