

TransAnimate: Taming Layer Diffusion to Generate RGBA Video

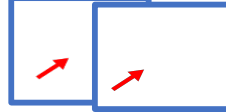
Anonymous CVPR submission

Paper ID

Image Control



Motion Control

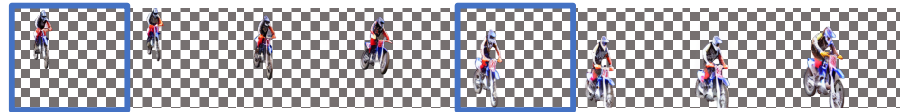


Sketch Control



Text to Video

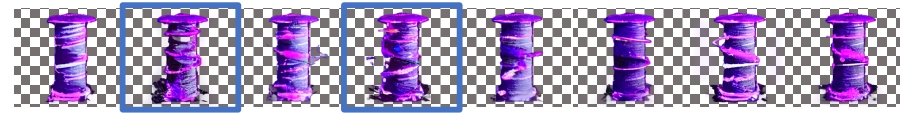
A motocross rider jumps and lands, with good motion and edge quality.



The disk shrinks moving diagonally, with good motion and edge quality.



A purple tornado swirls, with good motion and edge quality.



A horse runs on the ground, with good motion and edge quality.



Figure 1. **RGBA Video Generation with TransAnimate.** By utilizing pre-trained text-to-transparent image models, the motion-guided control mechanism, and the proposed dataset, TransAnimate enables high quality generation and effective control of video content.

Abstract

001 Text-to-video generative models have made remarkable ad-
 002 vancements in recent years. However, generating RGBA
 003 videos with alpha channels for transparency and visual ef-
 004 fects remains a significant challenge due to the scarcity
 005 of suitable datasets and the complexity of adapting exist-
 006 ing models for this purpose. To address these limitations,
 007 we present TransAnimate, an innovative framework that
 008 integrates RGBA image generation techniques with video
 009 generation modules, enabling the creation of dynamic and
 010 transparent videos. TransAnimate efficiently leverages pre-
 011 trained text-to-transparent image model weights and com-
 012 bines them with temporal models and controllability plug-
 013 ins trained on RGB videos, adapting them for controllable
 014 RGBA video generation tasks. Additionally, we introduce
 015 an interactive motion-guided control mechanism, where di-
 016 rectional arrows define movement and colors adjust scaling,
 017 offering precise and intuitive control for designing game
 018 effects. To further alleviate data scarcity, we have devel-

oped a pipeline for creating an RGBA video dataset, in-
 019 corporating high-quality game effect videos, extracted fore-
 020 ground objects, and synthetic transparent videos. Compre-
 021 hensive experiments demonstrate that TransAnimate gener-
 022 ates high-quality RGBA videos, establishing it as a practi-
 023 cal and effective tool for applications in gaming and visual
 024 effects. 025

1. Introduction

026 Text-to-video generation models have achieved remarkable
 027 progress in recent years, enabling the creation of dynamic
 028 and visually engaging content widely applied in video edit-
 029 ing, image animation, and motion customization. Further-
 030 more, methods incorporating control signals like optical
 031 flow, pose skeletons, and multimodal inputs have been pro-
 032 posed to achieve more precise video generation guidance.
 033 Alpha channels play a vital role in producing high-quality
 034 visual effects, as they allow transparent elements such as
 035 smoke, fire, and light to seamlessly integrate into complex
 036

037 scenes. This capability is particularly critical in game de-
038 velopment, where transparency effects are central to creat-
039 ing immersive and realistic experiences. However, generat-
040 ing controllable RGBA videos with special effects remains
041 highly desirable yet challenging.

042 Despite breakthroughs in video generation models,
043 RGBA video synthesis remains underdeveloped. This
044 dilemma stems from dual constraints: (1) the lack of
045 large-scale RGBA video datasets severely limits algorithmic
046 exploration, and (2) existing solutions like LayerDiffuse,
047 while capable of generating static transparent images,
048 cannot be seamlessly integrated with video generation mod-
049 els. The core challenge lies in leveraging massive RGB
050 video pretraining resources and RGBA image generation
051 models to build a data-efficient transparent video frame-
052 work—an urgent problem we address.

053 This paper presents TransAnimate, a framework that uni-
054 fies transparency modeling and video generation through
055 three synergistic innovations. First, we establish an
056 RGBA video synthesis pipeline and construct a founda-
057 tional RGBA video dataset. Second, we combine AnimateDiff’s
058 motion modeling capabilities, LayerDiffuse’s
059 transparency generation expertise, and sparse-control video
060 methods by fine-tuning adaptation layers on limited RGBA
061 video data. Third, we design motion control mechanisms
062 tailored for game visual effect artists, enabling pixel-precise
063 control through directional arrows for trajectory specifica-
064 tion and hue parameters for effect scaling.

065 While large-scale RGB video datasets contain millions
066 of samples, the scarcity of RGBA video data remains a sig-
067 nificant bottleneck. To address this limitation, we employ
068 three complementary data collection strategies: (a) curating
069 high-quality videos from game designers to capture authen-
070 tic transparency properties, (b) extracting foreground ob-
071 ject videos from instance segmentation data to enhance mo-
072 tion diversity, and (c) synthesizing controllable transparent
073 videos using parametric transformations such as translation
074 and scaling. Each strategy has distinct advantages and lim-
075 itations: (a) provides superior visual quality but lacks cat-
076 egory diversity, (b) introduces diverse motion patterns but
077 suffers from incomplete foregrounds and imperfect edges,
078 and (c) ensures precise edges and broad category cover-
079 age but is limited in motion complexity. Additionally, we
080 manually curate high-quality results from models trained on
081 these datasets and incorporate them as supplementary train-
082 ing data in an iterative refinement process. To balance these
083 trade-offs, we propose a positive trigger strategy, which as-
084 signs distinct learnable tokens to each data source. This
085 approach helps mitigate the impact of imperfect data distri-
086 butions by preventing the model from absorbing unwanted
087 features, thereby improving overall learning efficiency.

088 Architecturally, we first augment LayerDiffuse with
089 temporal modules initialized using AnimateDiff’s RGB-

pretrained weights, eliminating costly training from scratch.
Subsequent fine-tuning on RGBA video data enables effective
adaptation of motion priors from AnimateDiff. For controllable
generation, we repurpose RGB-pretrained SparseCtrl [13] with
input-layer fine-tuning—experiments show minimal parameter
adjustments suffice for high-consistency control.

Addressing game developers’ needs, we implement a vector-
chroma joint control system: Motion vectorization enables
8-directional trajectory control with velocity parameterization,
while chromatic scaling regulates effect magnitudes. This
integrated control mechanism significantly enhances flexibility
for game visual effects creation, empowering designers to
produce high-quality customizable effects with unprecedented
efficiency and precision. By unifying transparency modeling
and motion dynamics, our method achieves diverse RGBA
video generation while maintaining controllability comparable
to RGB video approaches.

Our contributions are summarized as follows:

- We propose TransAnimate, a framework for generating temporally coherent layered transparent videos and effects.
- We develop a dataset pipeline integrating game effects, segmented foregrounds, and synthetic videos, refined via iterative enhancement and a positive trigger strategy to mitigate distribution mismatches.
- We introduce novel directional-chromatic controls enabling pixel-accurate motion specification for game effects.

2. Related Work

Text-to-video generation. Recent advancements in text-to-video (T2V) generation [6, 15, 18, 21, 34, 44] have predominantly utilized diffusion models [16, 28, 32, 36], celebrated for their training stability and robust open-source ecosystems. A foundational milestone in this domain is the Video Diffusion Model [17], which extends a 2D image diffusion framework to handle video data by jointly training on image and video datasets from scratch. Building upon this, subsequent approaches leverage pre-trained image generators, such as Stable Diffusion [33], by integrating temporal layers into the existing 2D architectures and fine-tuning on expansive video datasets [1]. Among innovative methods, Align-Your-Latents [3] achieves efficient T2V conversion by aligning noise maps sampled independently for each frame, while AnimateDiff [12] incorporates a modular motion layer, enabling the generation of high-quality animations on customized image generation backbones [35]. To address temporal coherence challenges, Lumiere [2] eliminates the need for a temporal super-resolution module by directly generating videos with consistent frame rates. Further notable advancements include adopting scalable trans-

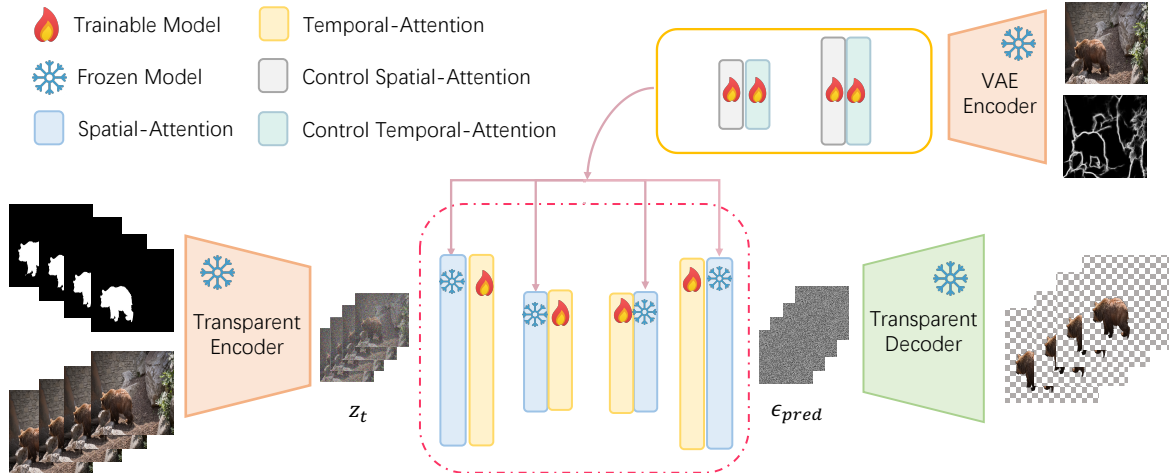


Figure 2. **Framework Overview.** TransAnimate generates transparent videos by learning motion patterns from videos. A frozen Transparent Encoder extracts features, refined by Temporal Attention and Linear Layers. Pre-trained SparseCtrl weights enable control via motion, sketches, and RGB images. A frozen Transparent Decoder reconstructs transparent frames, enhancing generation with limited RGBA data.

142 former architectures [26], leveraging spatiotemporal com-
 143 pressed latent spaces, as demonstrated by W.A.L.T.[14]
 144 and Sora[4], and utilizing discrete tokens alongside language
 145 models for video synthesis, exemplified by VideoPoet [23].
 146 Building on previous RGB video generation approaches,
 147 our method introduces text-to-RGBA video generation tai-
 148 lored for game effects design, addressing a significant gap
 149 in this domain.

150 **Controllable video generation.** Existing text-to-video
 151 models often suffer from limited control, as relying solely
 152 on text descriptions introduces ambiguity and reduces pre-
 153 cision. To address this, various methods incorporate explicit
 154 guidance signals. For instance, some approaches use depth
 155 maps or skeleton sequences to dictate scene layout or hu-
 156 man motion with greater accuracy [7, 8, 11, 20, 37, 45].
 157 Others leverage image-based control signals, which en-
 158 hance video quality and improve temporal consistency [11,
 159 27]. Despite these advances, controlling camera motion
 160 during video generation remains underexplored. Animate-
 161 Diff [12] applies LoRA [19] fine-tuning to adapt model
 162 weights for specific camera angles. Direct-a-Video [41]
 163 incorporates a camera embedder for pose control but sup-
 164 ports only basic parameters, limiting its capacity to simple
 165 motions like panning. MotionCtrl [39] extends this idea
 166 with additional input parameters to enable complex cam-
 167 era trajectories. However, its reliance on fine-tuning com-
 168 ponents of the diffusion model compromises generaliza-
 169 tion. In this work, we introduce a novel motion-guided con-
 170 trol framework specifically designed for generating trans-
 171 parent videos with predefined motion directions and scales.
 172 This method fills a critical gap by offering game effect de-
 173 signers a powerful and flexible tool for creating dynamic,
 174 high-quality content that meets the unique demands of their

workflows.

175 **Transparent Layer Processing** Most existing ap-
 176 proaches for transparent layer generation focus primarily
 177 on image generation, closely linked to image matting tech-
 178 niques. For instance, PPMatting [5] is a neural net-
 179 work model for image matting, trained from scratch us-
 180 ing standard matting datasets. Building on advancements
 181 in foundational models, Matting Anything [24] leverages
 182 the Segment Anything Model (SAM) [22] as its back-
 183 bone for matting tasks, while VitMatte [42] utilizes a Vi-
 184 sion Transformer (ViT) in a tri-map-based matting frame-
 185 work. Expanding beyond traditional matting, recent inno-
 186 vations explore the integration of layered effects in diffu-
 187 sion models. LayerDiffuse [43] introduces the concept of
 188 "latent transparency." This technique encodes alpha chan-
 189 nel transparency directly into the latent manifold of a pre-
 190 trained latent diffusion model by modifying the Variational
 191 Autoencoder (VAE) to decode alpha channels, enabling
 192 richer transparency effects. Despite these advancements,
 193 transparent video generation remains significantly underex-
 194 plored. SAM-2 [31] integrates a Transformer architecture
 195 with streaming storage and memory mechanisms to provide
 196 coherent segmentation predictions across video sequences.
 197 However, SAM-2 is not explicitly designed for video gen-
 198 eration and struggles to create diverse, high-quality trans-
 199 parent layers essential for effects like fire, light, and smoke,
 200 which are crucial in game effect design. In this work, we
 201 introduce a novel framework tailored for game effect gen-
 202 eration, enabling text-to-RGBA video generation with sup-
 203 port for various control modalities. Our approach addresses
 204 the limitations of existing methods, providing an innova-
 205 tive solution for creating visually compelling and dynamic
 206 transparent layers in video content, specifically for gaming
 207

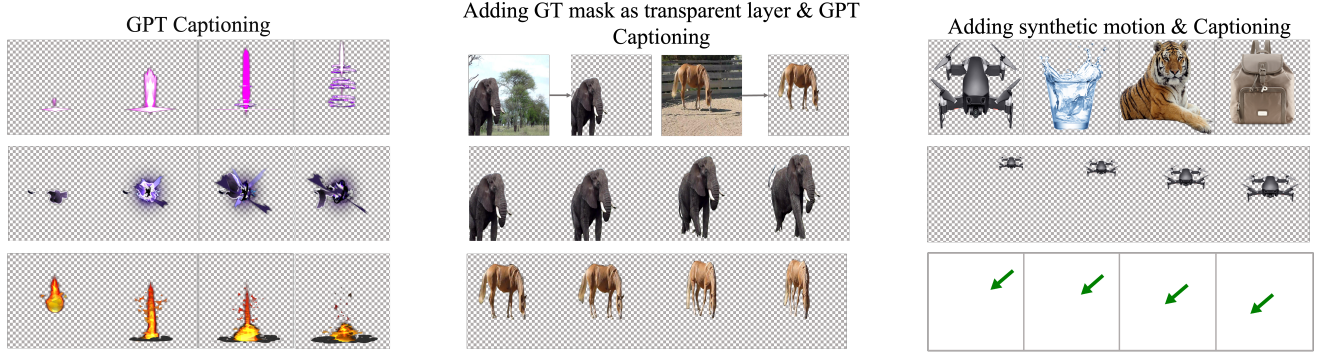


Figure 3. **Illustration of TransAnimate.** Our dataset consists of (a) Animate Dataset with 3,000 high-quality game effect videos, (b) Foreground Object Videos Dataset with 7,000 segmented videos capturing diverse motion patterns, and (c) Synthesized Transparent Motion Videos with 20,000 generated samples featuring controlled motion transformations. For synthesized dataset, from top to bottom, it represents Motion Caption, Raw Image, Synthetic, and Motion Control.

208 applications.

209 3. Method

210 Our approach enables a latent diffusion model to learn
 211 transferable motion priors from video data while leverag-
 212 ing existing text-to-transparent-image generation methods
 213 to produce transparent videos. In Section 3.1, we introduce
 214 the foundational principles of our method. Section 3.2 de-
 215 tails the dataset preparation process. Section 3.3 demon-
 216 strates how the proposed positive prompts enhance motion
 217 and edge quality. Section 3.4 elaborates on the TransAni-
 218 mate framework for transparent video generation. Finally,
 219 Section 3.5 explores how pretrained SparseCtrl weights are
 220 adapted to achieve controllable transparent video genera-
 221 tion.

222 3.1. Preliminary

223 Stable Diffusion (SD), a widely-used open-source text-to-
 224 image (T2I) model, performs the diffusion process in the
 225 latent space of a pre-trained autoencoder. This process per-
 226 turbs the encoded image $z_0 = \mathcal{E}(x_0)$ into z_t at step t by
 227 adding noise:

$$228 z_t = \sqrt{\bar{\alpha}_t} z_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \epsilon \sim \mathcal{N}(0, I), \quad (1)$$

229 where $\bar{\alpha}_t$ governs the noise strength. The denoising UNet
 230 ϵ_θ , utilizing an MSE loss, predicts noise while integrating
 231 ResNet, self-attention, and cross-attention mechanisms to
 232 incorporate text conditions effectively.

233 AnimateDiff [12] is an extension of T2I models, leverag-
 234 es video data to learn transferable motion priors, enab-
 235 ling animation generation through iterative denoising.
 236 The method consists of three core components: a domain
 237 adapter, a motion module, and a LoRA-based extension, all
 238 seamlessly integrated into the T2I framework during infer-
 239 ence.

240 To model temporal dynamics, AnimateDiff inflates the
 241 2D T2I model to handle 3D video data. Video tensors $x \in$

$\mathbb{R}^{b \times c \times f \times h \times w}$ are reshaped to allow the image layers to process
 frames independently, while a newly introduced motion
 module captures temporal dependencies across frames.

The motion module in AnimateDiff employs a temporal
 Transformer enhanced with self-attention blocks and sinu-
 soidal position encodings. Input features are reshaped into
 sequences along the temporal axis to facilitate inter-frame
 information exchange:

$$z_{out} = \text{Softmax}(QK^T/\sqrt{c}) \cdot V, \quad (2)$$

where Q, K, V are linear projections of the input features.
 This design effectively encodes motion priors, enabling the
 generation of coherent and realistic animations.

254 3.2. New Dataset for Transparent Video Generation

255 Achieving high-quality transparent video generation
 256 presents three major challenges: (1) ensuring high-fidelity
 257 transparent layers, (2) capturing realistic and dynamic
 258 motion, and (3) maintaining sufficient diversity in video
 259 content. To address these challenges, we construct a
 260 comprehensive dataset by collecting and synthesizing three
 261 distinct types of video data, ensuring that our network
 262 learns from diverse samples with high-quality motion and
 263 transparency.

Animate Dataset. The animate dataset consists of 3,000
 professionally curated game effect videos created by expert
 designers. These videos feature complex and visually rich
 animations, such as explosions, energy transformations, and
 magical effects. Each video is carefully selected to ensure
 high transparency quality and realistic motion.

Foreground Object Videos Dataset. To enhance di-
 versity and capture a wide range of motion patterns,
 we construct a foreground object dataset by extracting
 7,000 videos from large-scale video instance segmenta-
 tion datasets, including VideoMatte240K [25], MeViS [9],
 YouTube-VOS [40], DAVIS [29], and MOSE [10]. These
 datasets cover a broad range of object categories, ensuring

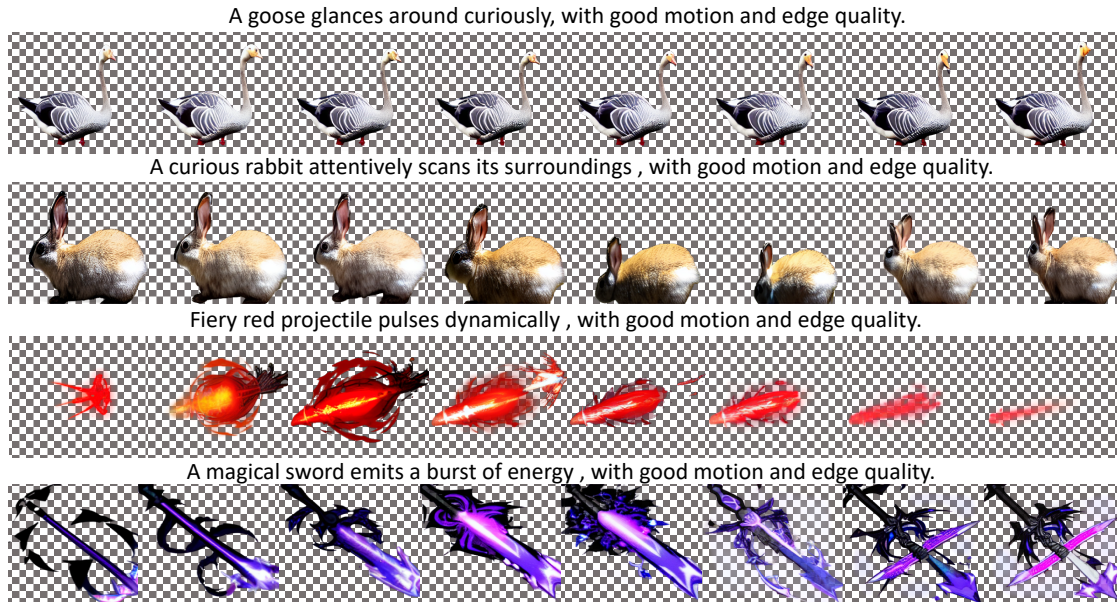


Figure 4. Text-to-RGBA video generation results of TransAnimate.

Dataset	RGBA Edge Quality	Motion Quality	Category Diversity
Game Effects Dataset	✓	✓	✗
Foreground Dataset	✗	✓	✓
Synthetic Dataset	✓	✗	✓
Iterated Dataset	✓	✓	✓

Table 1. Comparison of datasets based on three key quality aspects.

the dataset includes diverse appearances, textures, and motion trajectories. Using video instance segmentation techniques, we extract foreground objects and retain only those exhibiting significant motion dynamics.

Synthesized Transparent Motion Videos. To further improve transparency quality and motion diversity, we generate 20,000 synthetic transparent videos by applying controlled motion transformations to foreground images from predefined transparent classes. These transformations include translation, scaling, and rotation, ensuring that the dataset captures a wide variety of movement behaviors. The motion characteristics are visualized with directional arrows (for translation) and color-coded indicators (green for scaling up, blue for stability, and red for shrinking). Captions are generated based on the transparent class and applied motion transformation, ensuring a strong alignment between textual descriptions and visual content.

Data Iteration We jointly train the TransAnimate model on the aforementioned three datasets, employing a data iteration strategy to manually curate high-quality generated results as supplementary training data. This effectively mit-

igates the scarcity of RGBA data. The animate dataset ensures high-quality transparent layers and realistic motion, the foreground object dataset provides diversity in motion and appearance, and the synthesized transparent motion dataset further refines transparency and motion control. To address the limitations of different datasets, we adopt a negative trigger strategy to filter out undesirable features, significantly enhancing the performance of transparent video generation.

3.3. Positive Triggers

As shown in Table 1, the datasets we collect exhibit certain limitations in terms of quality attributes. Specifically, the Foreground Dataset demonstrates strong motion quality but suffers from poor edge quality, while the Synthetic Dataset excels in edge quality but lacks sufficient motion fidelity. To mitigate these issues, we leverage positive triggers during training, explicitly guiding the model to learn high-quality attributes from different datasets.

For instances originating from the Synthetic Dataset, we include positive triggers emphasizing good edge quality, helping the model recognize and reinforce well-defined edges. Conversely, for samples from the Foreground Dataset, we incorporate positive triggers highlighting strong motion attributes, enabling the model to better capture dynamic motion patterns. This targeted approach allows the network to develop a more comprehensive understanding of both motion and edge quality during training.

During inference, we combine positive triggers for both good edge and motion quality, ensuring that the generated results exhibit improved sharpness and dynamic consistency. This strategy enhances the overall quality of

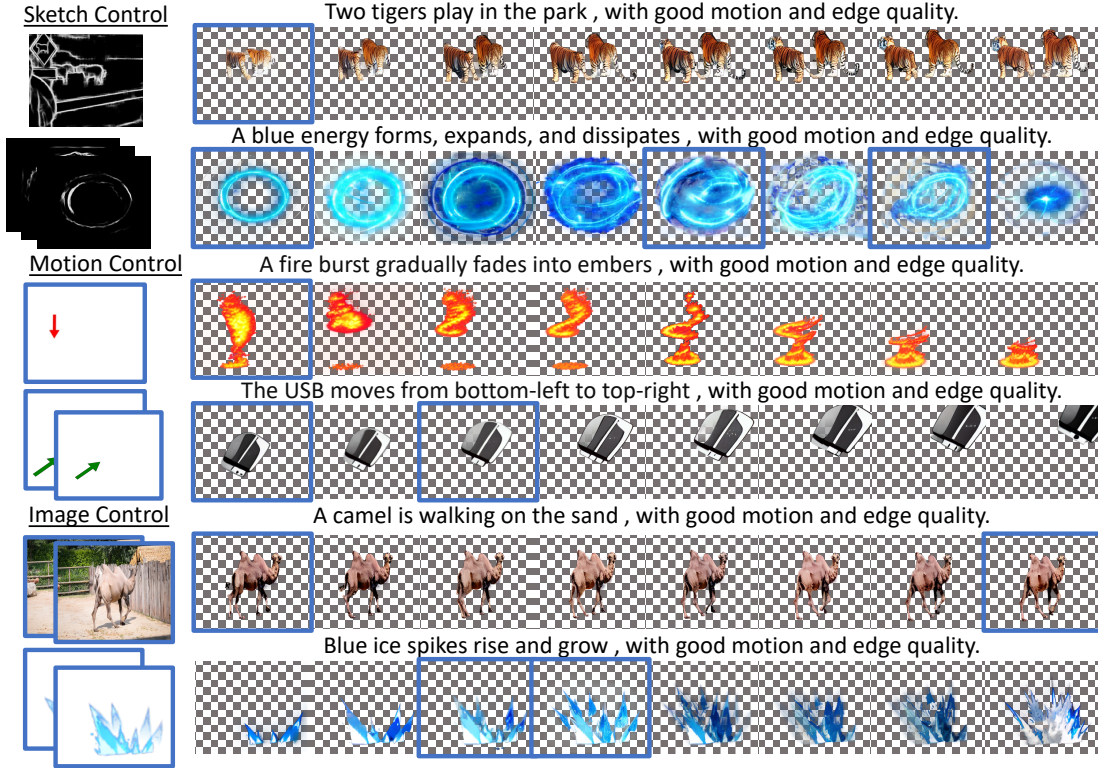


Figure 5. Conditional generation results from TransAnimate. The qualitative results are results with sketch, depth, and RGB image conditions. The input conditions are displayed on the left, while the keyframes guided by these conditions are highlighted with blue borders.

329 RGBA video generation by bridging the strengths of multi-
330 ple datasets.

331 3.4. Text2RGBA Video Generation

332 Our method introduces a novel approach to RGBA video
333 generation by extending text-to-transparent-image frame-
334 works, particularly by using LayerDiffuse [43]. Unlike
335 prior methods focused solely on RGB, we address the
336 unique requirements of generating transparent animations,
337 crucial for game effects design.

338 The framework integrates a motion module inspired
339 by AnimateDiff [12], enabling animation through iterative
340 denoising. This module is trained independently while
341 keeping the text-to-transparent-image generation network
342 frozen. By incorporating this learnable module into Lay-
343 erDiffuse, we achieve temporal animation while preserving
344 the high-quality transparency features of the original model.

345 To adapt LayerDiffuse for video generation, we ext-
346 tend its architecture to process 5D video tensors $\mathbf{x} \in$
347 $\mathbb{R}^{b \times c \times f \times h \times w}$, where b , c , f , h , and w denote the batch size,
348 channels, frames, height, and width, respectively. Follow-
349 ing a network inflation strategy similar to AnimateDiff [3],
350 we allow the image layers to process frames independently
351 by reshaping the temporal axis f into the batch axis during
352 feature extraction and restoring it afterward.

353 Conversely, the motion module reshapes spatial dimen-
354 sions (h, w) into the batch axis during temporal processing,
355 which is restored post-processing. This design ensures effi-
356 cient and independent handling of spatial and temporal di-
357 mensions.

358 The motion module captures temporal dependencies by
359 dividing the reshaped feature map along the temporal axis
360 into a sequence of vectors $\{z_1, z_2, \dots, z_f\}$, where each
361 vector corresponds to a frame. These vectors are pro-
362 cessed through self-attention layers, enabling temporal in-
363 teractions:

$$364 z_{out} = \text{Softmax} \left(\frac{QK^T}{\sqrt{c}} \right) V, \quad (3)$$

365 where $Q = W^Q z$, $K = W^K z$, and $V = W^V z$ are the
366 query, key, and value projections of the input features. This
367 mechanism ensures temporal coherence by enabling infor-
368 mation flow across frames.

369 During training, the encoder, decoder, and 3D U-Net
370 components of LayerDiffuse are frozen, and only the mo-
371 tion module is trained. The diffusion process progressively
372 adds noise to the latent representation \mathbf{x}_t over time steps
373 t . The network predicts the added noise using a denoising

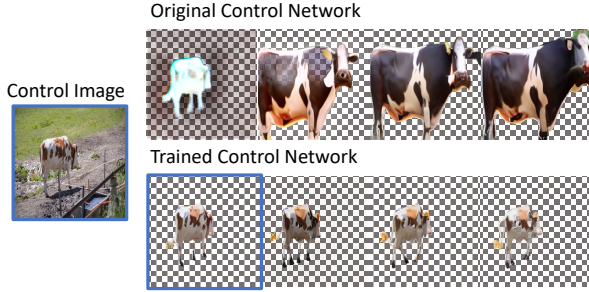


Figure 6. Qualitative comparison of SparseControl’s original weights with our trained weights.

	FVD	CLIP	RGBA Quality	Motion Quality
Original control	1400.86	19.22	16.32	18.17
Trained control	407.41	30.57	83.62	82.28

Table 2. Quantitative comparison of SparseControl’s original weights with our trained weights.

374 model ϵ_θ , minimizing the following objective:

$$375 \mathcal{L} = \mathbb{E}_{\mathcal{E}(x_0^{1:f}), y, \epsilon^{1:f} \sim \mathcal{N}(0, I), t} \left[\|\epsilon - \epsilon_\theta(z_t^{1:f}, t, \tau_\theta(y))\|_2^2 \right]. \quad (4)$$

376 At inference, the pre-trained latent transparency decoder
377 of LayerDiffusion, $\mathcal{D}(\cdot, \cdot)$, reconstructs the video’s RGB and
378 alpha channels. Given the adjusted latent representation x_a ,
379 the decoder outputs the transparent video:

$$380 [\hat{I}_c \hat{I}_\alpha] = \mathcal{D}(\hat{I}, x_a), \quad (5)$$

381 where \hat{I}_c and \hat{I}_α represent the reconstructed color and alpha
382 channels, respectively. This process ensures high-quality
383 RGBA video outputs suitable for game effects.

384 3.5. Controllable RGBA Video Generation

385 This section demonstrates how to apply SparseCtrl [11],
386 which was trained on RGB videos, to RGBA controllable
387 generation, incorporating three types of controls: motion
388 control, RGB image control, and sketch control. Notably,
389 we propose utilizing only the RGB channels of RGBA
390 videos as control inputs, which is more practical for down-
391 stream tasks. We reuse the pre-trained weights of Spar-
392 seCtrl and keep them frozen, while fine-tuning the Adapter
393 layer on RGBA videos to align with TransAnimate Unet.
394 This strategy effectively leverages the prior knowledge of
395 the pre-trained RGB model, enabling RGBA video genera-
396 tion and control even with limited RGBA training data.

397 **Motion-Guided Video Generation.** Motion guidance
398 is crucial for intuitive game effect design, enabling creators
399 to precisely control the direction and scale of animations.
400 Existing methods often lack support for such fine-grained
401 functionality. Our motion control mechanism allows users

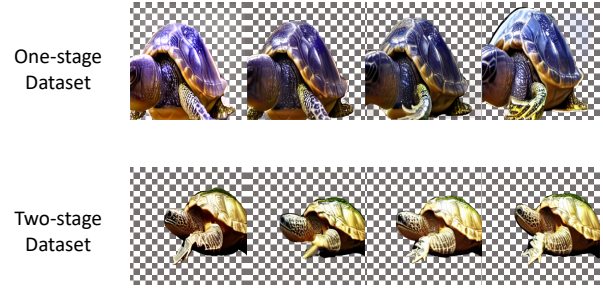


Figure 7. Qualitative comparison of results obtained using one-stage and two-stage training datasets.

	FVD	CLIP	RGBA Quality	Motion Quality
One-stage	485.63	28.55	67.05	69.17
Two-stage	429.82	29.32	76.32	78.79

Table 3. Quantitative comparison of results obtained using one-stage and two-stage training datasets.

to define the generation direction using arrows and indi- 402
cate the effect scale using colors. For example, a transpar- 403
ent video effect can be guided to move from left to right 404
while gradually shrinking, providing game designers with 405
enhanced control and adaptability for their creative work- 406
flows. 407

Sketch-to-Video Generation. Sketches offer an intu- 408
itive and accessible way for non-professional users to guide 409
T2V generation. SparseCtrl [11] allows users to input 410
sketches to shape video content. A single sketch can define 411
the overall layout, while multiple sketches—such as 412
for the first, last, and key intermediate frames—can guide 413
coarse motion and transitions. This functionality is espe- 414
cially suited for tasks like storyboarding, enabling creators 415
to visualize and iterate on video concepts effortlessly. 416

Image Animation and Transition. SparseCtrl [11] unifies 417
various video generation tasks, including video predic- 418
tion, animation, and interpolation, under a single frame- 419
work leveraging RGB image conditions. Image animation 420
generates videos based on the first frame, while transitions 421
are guided by both the first and last frames. Video predic- 422
tion uses initial frames to extrapolate motion, and interpo- 423
lation creates smooth transitions between sparsely provided 424
keyframes. By unifying these tasks, SparseCtrl broadens 425
the applicability of video generation methods, making them 426
versatile tools for diverse creative and practical scenarios. 427

428 4. Experiments

429 4.1. Experiment Setting

Setup. Our method leverages the pre-trained VAE encoder, 430
decoder, and 3D-UNet architecture from LayerDiffuse, with 431

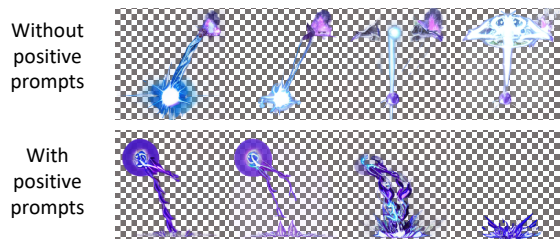


Figure 8. Qualitative comparison of results with and without the proposed positive triggers.

432 the motion module adapted from AnimateDiff. The training
 433 resolution is set to 256×256 for 16 frames. We train the
 434 model for a total of 3,000 iterations with a batch size of 16,
 435 utilizing two NVIDIA A100 80GB GPUs. The learning rate
 436 is set to 1×10^{-5} , ensuring stable and effective optimization.
 437 **Evaluation Metrics.** We evaluate RGBA and motion quality
 438 using CLIP Score [30], FVD [38], and a user study. In
 439 the user study, 15 participants assess 50 generated videos
 440 based on two key aspects: (1) the quality of RGBA edges and
 441 (2) the realism and smoothness of motion.

442 4.2. Qualitative Results

443 Figure 4 presents qualitative results. These examples il-
 444 lustrate the model’s ability to produce high-quality trans-
 445 parency effects while preserving fine details and struc-
 446 tural consistency. Furthermore, the results demonstrate the
 447 model’s strong generalization across various content types,
 448 reinforcing its robustness and adaptability.

449 4.3. Qualitative Controllable Results

450 Our approach seamlessly integrates with SparseControl to
 451 facilitate multi-modal video generation, enabling precise
 452 control over object structure, movement, and appearance.
 453 Users can provide sketches, motion trajectories, or refer-
 454 ence images to guide the generation process. Figure 5 il-
 455 lustrates qualitative results, demonstrating the model’s ca-
 456 pability to generate transparent images with high fidelity
 457 while maintaining user-defined control inputs. These results
 458 highlight the model’s ability to produce visually compelling
 459 outputs and generalize effectively across diverse content.

460 4.4. Ablative Study

461 **Effectiveness of Control Training.** We evaluate the effec-
 462 tiveness of our control method by comparing it with the pre-
 463 trained weights provided by SparseControl. As shown in
 464 Figure 6 and Table 2, directly applying RGB pre-trained
 465 weights fails to preserve structural consistency in control-
 466 lable RGBA video generation. In contrast, our approach
 467 successfully adapts RGB-based control to RGBA videos,
 468 achieving superior alignment with the control image while
 469 maintaining transparency effects. These results highlight

	FVD	CLIP	RGBA Quality	Motion Quality
w/o Positive Prompts	458.37	28.95	71.58	73.31
w Positive Prompts	429.827	29.32	76.32	78.79

Table 4. Quantitative comparison of results with and without the proposed positive triggers.

the ability of pre-trained RGB control weights to general-
 ize effectively to RGBA video generation, addressing the
 challenges posed by the scarcity of RGBA training data.

Two-Stage Training. We examine the impact of a two-
 stage training strategy on generation quality. As illustrated
 in Figure 7, this approach enhances dataset diversity and
 scalability without compromising visual fidelity, resulting
 in more stable and high-quality video generation. Addi-
 tionally, the quantitative results in Table 3 confirm that the
 two-stage training strategy significantly improves the final
 performance of the generated RGBA videos, demonstrating
 its effectiveness in refining output quality.

Positive Triggers. We analyze the contribution of the
 proposed positive prompts strategy in enhancing generation
 quality, particularly in motion coherence and RGBA edge
 fidelity. As shown in Figure 8, this technique enriches
 dataset diversity while preserving structural integrity, lead-
 ing to more consistent and visually refined outputs. Fur-
 thermore, Table 4 provides quantitative evidence that incor-
 porating positive prompts significantly enhances the overall
 quality of generated RGBA videos, reinforcing its effective-
 ness as a control mechanism.

5. Conclusion

In this paper, we address the underexplored challenge
 of RGBA video generation by introducing TransAnimate,
 a novel framework that combines transparency modeling
 with motion dynamics to generate high-quality, layered,
 transparent videos with temporal coherence. To tackle
 the scarcity of RGBA datasets, we propose a three-step
 data creation strategy, leveraging high-quality game effect
 videos, extracted foreground objects, and synthetic trans-
 parent videos with controlled motion dynamics. We also
 introduce a motion-guided control mechanism that enables
 precise adjustments to motion direction and scaling, tai-
 lored to game effect creation. By utilizing RGB channels
 as control inputs, our approach effectively leverages pre-
 trained RGB models, enabling robust RGBA video gener-
 ation and control with limited RGBA data. Our contribu-
 tions bridge critical gaps in RGBA video generation, pro-
 viding a unified framework with controllability compar-
 able to RGB methods while extending capabilities to layered
 transparency. This work opens new possibilities for immer-
 sive and customizable content creation in gaming and visual
 effects.

514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570**References**

- [1] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1728–1738, 2021. 2
- [2] Omer Bar-Tal, Hila Chefer, Omer Tov, Charles Herrmann, Roni Paiss, Shiran Zada, Ariel Ephrat, Junhwa Hur, Yuanzhen Li, Tomer Michaeli, et al. Lumiere: A space-time diffusion model for video generation. *arXiv preprint arXiv:2401.12945*, 2024. 2
- [3] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22563–22575, 2023. 2, 6
- [4] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. 2024. 3
- [5] Guowei Chen, Yi Liu, Jian Wang, Juncai Peng, Yuying Hao, Lutao Chu, Shiyu Tang, Zewu Wu, Zeyu Chen, Zhiliang Yu, et al. Pp-matting: high-accuracy natural image matting. *arXiv preprint arXiv:2204.09433*, 2022. 3
- [6] Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, Chao Weng, and Ying Shan. Videocrafter1: Open diffusion models for high-quality video generation, 2023. 2
- [7] Weifeng Chen, Jie Wu, Pan Xie, Hefeng Wu, Jiashi Li, Xin Xia, Xuefeng Xiao, and Liang Lin. Control-a-video: Controllable text-to-video generation with diffusion models. *arXiv preprint arXiv:2305.13840*, 2023. 3
- [8] Zhimin Chen, Longlong Jing, Yingwei Li, and Bing Li. Bridging the domain gap: Self-supervised 3d scene understanding with foundation models. *Advances in Neural Information Processing Systems*, 36:79226–79239, 2023. 3
- [9] Henghui Ding, Chang Liu, Shuting He, Xudong Jiang, and Chen Change Loy. Mevis: A large-scale benchmark for video segmentation with motion expressions. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2694–2703, 2023. 4
- [10] Henghui Ding, Chang Liu, Shuting He, Xudong Jiang, Philip HS Torr, and Song Bai. Mose: A new dataset for video object segmentation in complex scenes. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 20224–20234, 2023. 4
- [11] Yuwei Guo, Ceyuan Yang, Anyi Rao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Sparsectrl: Adding sparse controls to text-to-video diffusion models. *arXiv preprint arXiv:2311.16933*, 2023. 3, 7
- [12] Yuwei Guo, Ceyuan Yang, Anyi Rao, Yaohui Wang, Yu Qiao, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023. 2, 3, 4, 6
- [13] Yuwei Guo, Ceyuan Yang, Anyi Rao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Sparsectrl: Adding sparse controls to text-to-video diffusion models. In *European Conference on Computer Vision*, pages 330–348. Springer, 2024. 2
- [14] Agrim Gupta, Lijun Yu, Kihyuk Sohn, Xiuye Gu, Meera Hahn, Li Fei-Fei, Irfan Essa, Lu Jiang, and José Lezama. Photorealistic video generation with diffusion models. *arXiv preprint arXiv:2312.06662*, 2023. 3
- [15] Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion models for high-fidelity video generation with arbitrary lengths. *arXiv preprint arXiv:2211.13221*, 2022. 2
- [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 2
- [17] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *arXiv:2204.03458*, 2022. 2
- [18] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022. 2
- [19] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 3
- [20] Li Hu, Xin Gao, Peng Zhang, Ke Sun, Bang Zhang, and Liefeng Bo. Animate anyone: Consistent and controllable image-to-video synthesis for character animation. *arXiv preprint arXiv:2311.17117*, 2023. 3
- [21] Johanna Karras, Aleksander Holynski, Ting-Chun Wang, and Ira Kemelmacher-Shlizerman. Dreampose: Fashion video synthesis with stable diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22680–22690, 2023. 2
- [22] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 3
- [23] Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Rachel Hornung, Hartwig Adam, Hassan Akbari, Yair Alon, Vighnesh Birodkar, et al. Videopoet: A large language model for zero-shot video generation. *arXiv preprint arXiv:2312.14125*, 2023. 3
- [24] Jiachen Li, Jitesh Jain, and Humphrey Shi. Matting anything. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1775–1785, 2024. 3
- [25] Shanchuan Lin, Andrey Ryabtsev, Soumyadip Sengupta, Brian L Curless, Steven M Seitz, and Ira Kemelmacher-Shlizerman. Real-time high-resolution background matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8762–8771, 2021. 4
- [26] Xin Ma, Yaohui Wang, Gengyun Jia, Xinyuan Chen, Ziwei Liu, Yuan-Fang Li, Cunjian Chen, and Yu Qiao. Latte: Latent diffusion transformer for video generation. *arXiv preprint arXiv:2401.03048*, 2024. 3

- 628 [27] Haomiao Ni, Changhao Shi, Kai Li, Sharon X Huang, and
629 Martin Renqiang Min. Conditional image-to-video gener-
630 ation with latent flow diffusion models. In *Proceedings of*
631 *the IEEE/CVF Conference on Computer Vision and Pattern*
632 *Recognition*, pages 18444–18455, 2023. 3
- 633 [28] William Peebles and Saining Xie. Scalable diffusion models
634 with transformers. In *Proceedings of the IEEE/CVF Inter-*
635 *national Conference on Computer Vision*, pages 4195–4205,
636 2023. 2
- 637 [29] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Ar-
638 beláez, Alex Sorkine-Hornung, and Luc Van Gool. The 2017
639 davis challenge on video object segmentation. *arXiv preprint*
640 *arXiv:1704.00675*, 2017. 4
- 641 [30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
642 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
643 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning
644 transferable visual models from natural language supervi-
645 sion. In *International conference on machine learning*, pages
646 8748–8763. PMLR, 2021. 8
- 647 [31] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang
648 Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman
649 Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junt-
650 ing Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-
651 Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feicht-
652 enhofer. Sam 2: Segment anything in images and videos.
653 *arXiv preprint arXiv:2408.00714*, 2024. 3
- 654 [32] Xuqian Ren, Wenjia Wang, Dingding Cai, Tuuli Tuominen,
655 Juho Kannala, and Esa Rahtu. Mushroom: Multi-sensor hy-
656 brid room dataset for joint 3d reconstruction and novel view
657 synthesis. In *Proceedings of the IEEE/CVF Winter Confer-*
658 *ence on Applications of Computer Vision*, pages 4508–4517,
659 2024. 2
- 660 [33] Robin Rombach, Andreas Blattmann, Dominik Lorenz,
661 Patrick Esser, and Björn Ommer. High-resolution image
662 synthesis with latent diffusion models. In *Proceedings of*
663 *the IEEE/CVF Conference on Computer Vision and Pattern*
664 *Recognition*, pages 10684–10695, 2022. 2
- 665 [34] Ludan Ruan, Yiyang Ma, Huan Yang, Huiguo He, Bei Liu,
666 Jianlong Fu, Nicholas Jing Yuan, Qin Jin, and Baining
667 Guo. Mm-diffusion: Learning multi-modal diffusion mod-
668 els for joint audio and video generation. In *Proceedings of*
669 *the IEEE/CVF Conference on Computer Vision and Pattern*
670 *Recognition*, pages 10219–10228, 2023. 2
- 671 [35] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch,
672 Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine
673 tuning text-to-image diffusion models for subject-driven
674 generation. In *Proceedings of the IEEE/CVF Conference*
675 *on Computer Vision and Pattern Recognition*, pages 22500–
676 22510, 2023. 2
- 677 [36] Jiaming Song, Chenlin Meng, and Stefano Ermon.
678 Denoising diffusion implicit models. *arXiv preprint*
679 *arXiv:2010.02502*, 2020. 2
- 680 [37] Matias Turkulainen, Xuqian Ren, Iaroslav Melekhov, Otto
681 Seiskari, Esa Rahtu, and Juho Kannala. Dn-splatter: Depth
682 and normal priors for gaussian splatting and meshing. *arXiv*
683 *preprint arXiv:2403.17822*, 2024. 3
- 684 [38] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach,
Raphaël Marinier, Marcin Michalski, and Sylvain Gelly. Fvd: A new metric for video generation. 2019. 8
- [39] Zhouxia Wang, Ziyang Yuan, Xintao Wang, Tianshui Chen, Menghan Xia, Ping Luo, and Ying Shan. Motionctrl: A unified and flexible motion controller for video generation. *arXiv preprint arXiv:2312.03641*, 2023. 3
- [40] Ning Xu, Linjie Yang, Yuchen Fan, Dingcheng Yue, Yuchen Liang, Jianchao Yang, and Thomas Huang. Youtube-vos: A large-scale video object segmentation benchmark. *arXiv preprint arXiv:1809.03327*, 2018. 4
- [41] Shiyuan Yang, Liang Hou, Haibin Huang, Chongyang Ma, Pengfei Wan, Di Zhang, Xiaodong Chen, and Jing Liao. Direct-a-video: Customized video generation with user-directed camera movement and object motion. *arXiv preprint arXiv:2402.03162*, 2024. 3
- [42] Jingfeng Yao, Xinggang Wang, Shusheng Yang, and Baoyuan Wang. Vitmatte: Boosting image matting with pre-trained plain vision transformers. *Information Fusion*, 103: 102091, 2024. 3
- [43] Lvmin Zhang and Maneesh Agrawala. Transparent image layer diffusion using latent transparency. *arXiv preprint arXiv:2402.17113*, 2024. 3, 6
- [44] Shiwei Zhang, Jiayu Wang, Yingya Zhang, Kang Zhao, Hangjie Yuan, Zhiwu Qin, Xiang Wang, Deli Zhao, and Jingren Zhou. I2vgen-xl: High-quality image-to-video synthesis via cascaded diffusion models. *arXiv preprint arXiv:2311.04145*, 2023. 2
- [45] Yabo Zhang, Yuxiang Wei, Dongsheng Jiang, Xiaopeng Zhang, Wangmeng Zuo, and Qi Tian. Controlvideo: Training-free controllable text-to-video generation. *arXiv preprint arXiv:2305.13077*, 2023. 3