

# GenDoP: Auto-regressive Camera Trajectory Generation as a Director of Photography

Mengchen Zhang<sup>1,2</sup>, Tong Wu<sup>3</sup>, Jing Tan<sup>4</sup>, Ziwei Liu<sup>5</sup>, Gordon Wetzstein<sup>3</sup>, Dahua Lin<sup>2,4,6</sup>

<sup>1</sup>Zhejiang University, <sup>2</sup>Shanghai Artificial Intelligence Laboratory, <sup>3</sup>Stanford University,

<sup>4</sup>The Chinese University of Hong Kong, <sup>5</sup>Nanyang Technological University, <sup>6</sup>CPII under InnoHK

zhangmengchen@zju.edu.cn, {wutong16, gordon.wetzstein}@stanford.edu

{tj023, dhlin}@ie.cuhk.edu.hk, ziwei.liu@ntu.edu.sg

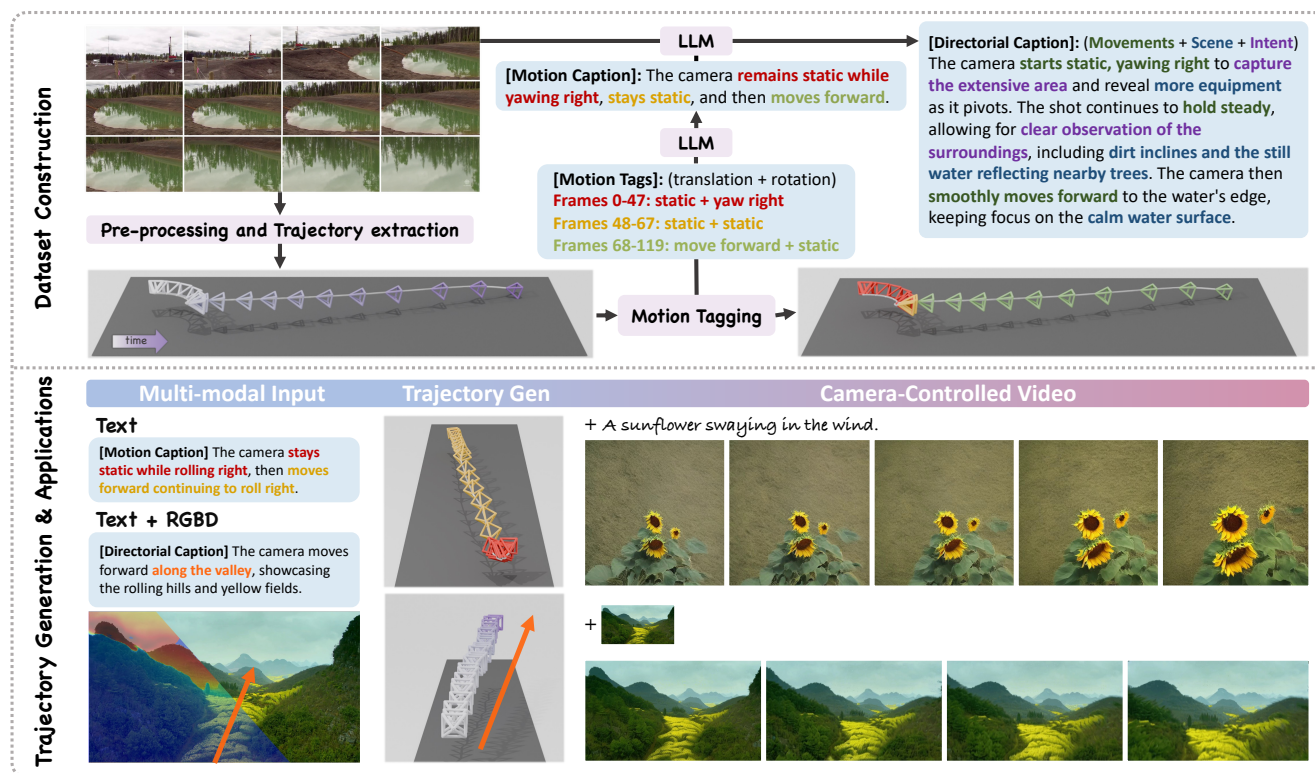


Figure 1. **Overview.** **Top:** DataDoP data construction. Given RGB video frames, we extract RGBD images and camera poses, then tag the pose sequence with different motion categories (in different colors). With LLM, we generate two types of captions from motion tags and RGBD inputs: *Motion Caption* describes the camera movements, while *Directorial Caption* describes the **camera movements** along with their **interaction with the scene** and **directorial intent**. **Bottom:** Our GenDoP method supports multi-modal inputs for trajectory creation. The generated camera sequence can be easily applied to various video generation tasks, including text-to-video (T2V) [12] and image-to-video (I2V) generation [14]. GenDoP paves the way for future advancements in *camera-controlled* video generation.

## Abstract

Camera trajectory design plays a crucial role in video production, serving as a fundamental tool for conveying directorial intent and enhancing visual storytelling. In cinematography, Directors of Photography meticulously craft camera movements to achieve expressive and intentional framing. However, existing methods for camera trajec-

tory generation remain limited: Traditional approaches rely on geometric optimization or handcrafted procedural systems, while recent learning-based methods often inherit structural biases or lack textual alignment, constraining creative synthesis. In this work, we introduce an auto-regressive model inspired by the expertise of Directors of Photography to generate artistic and expressive camera trajectories. We first introduce **DataDoP**,

*a large-scale multi-modal dataset containing 29K real-world shots with free-moving camera trajectories, depth maps, and detailed captions in specific movements, interaction with the scene, and directorial intent. Thanks to the comprehensive and diverse database, we further train an auto-regressive, decoder-only Transformer for high-quality, context-aware camera movement generation based on text guidance and RGBD inputs, named **GenDoP**. Extensive experiments demonstrate that compared to existing methods, GenDoP offers better controllability, finer-grained trajectory adjustments, and higher motion stability. We believe our approach establishes a new standard for learning-based cinematography, paving the way for future advancements in camera control and filmmaking. Our project website: <https://kszpxxzm.github.io/GenDoP/>.*

## 1. Introduction

In video production, the camera serves as the window of observation, playing a crucial role in presenting scene content, conveying the director’s intent, and achieving visual effects. In recent years, video generation technology has advanced [1, 3, 23, 46], and several cutting-edge studies have explored camera-controlled video generation [12, 14, 29, 36]. However, these works often rely on predefined, simplistic camera trajectories to demonstrate their results. The generation of artistic, expressive, and intentional camera movements remains largely unexplored.

Trajectory generation has been a long-standing problem. Traditional approaches include optimization-based camera motion planning [4, 10, 27] and learning-based camera control [5, 8, 17, 20]. However, these techniques demand geometric modeling or cost function engineering for each motion, which limits creative synthesis. Meanwhile, oversimplified procedural systems impede precise text control. Recent advances in diffusion-based camera trajectory generation [7, 21, 24] have expanded creative possibilities for text-driven cinematography. However, CCD [21] and E.T. [7] inherit structural biases from human-centric tracking datasets, constraining camera movements to oversimplified character-relative motion patterns. Director3D [24] introduces object/scene-centric 3D trajectories from multi-view datasets [42, 47], but the lack of trajectory-level captions limits text-to-motion alignment. As a result, the generated paths are driven by geometric plausibility rather than directorial intent. These dataset constraints hinder the creation of artistically coherent free-moving trajectories that interpret creative vision without relying on specific subjects.

In this work, we tackle the problems above with several key designs. First, we introduce **DataDoP** Dataset, a multi-modal, free-moving camera motion dataset extracted from real video clips, which includes accurate camera trajectories

extracted by state-of-the-art and scene compositions. We extract camera trajectories and corresponding depth maps using MonST3R [43], and employ GPT-4o to generate comprehensive descriptions of the camera trajectories and scene focus, capturing both motion dynamics and directorial intent. DataDoP comprises over 29K shots, totaling 11M frames, with corresponding camera trajectories and diverse textual descriptions. Furthermore, given the inherently sequential nature of camera trajectories, we propose **GenDoP**, which treats camera parameters as discrete tokens and leverages an auto-regressive model for camera trajectory generation. Our model incorporates multi-modal condition as inputs, including fine-grained textual descriptions and optionally RGBD information from the first frame, to produce stable, complex, and instruction-aligned camera movements.

We conduct rigorous human validation to ensure the dataset quality. Extensive experiments confirm that GenDoP outperforms state-of-the-art methods [7, 21, 24] across fine-grained textual controllability, motion stability, and complexity, while exhibiting enhanced robustness. As AI-driven video creation evolves, multi-modal camera trajectory generation emerges as a timely and crucial direction. We believe that this work paves the way for future advancements in camera-controlled video generation and a wide range of trajectory-related downstream applications.

## 2. Related Work

**Camera trajectory datasets.** While existing datasets [7, 21, 26, 42, 44, 47] document camera trajectories, their cinematographic expressiveness remains constrained. Datasets such as MVImgNet [42], RealEstate10K [47], and DL3DV-10K [26] provide calibrated trajectories through structured capture methods, but predominantly focus on basic paths around static objects or scenes. These datasets lack the sophisticated cinematographic language necessary for narrative-driven sequencing and intentional viewpoint control. CCD [21] and E.T. [7] emphasize human-centric tracking but are confined to reactive tracking mechanisms. In contrast, DataDoP’s camera movement is driven by the compositional logic of the scene and the narrative demands. We underscore DataDoP’s unique contribution to the field of artistic camera trajectory generation.

**Camera trajectory generation.** Early efforts in trajectory generation generally consist optimization-based motion planning [4, 10, 27, 28] and learning-based camera control [5, 8, 17, 20]. Recent progress focuses on integrating camera motion with scene and character dynamics. CCD [21] introduced a camera diffusion model using text and keyframe controls, generating motion in character-centric coordinates. E.T. [7] improves to incorporate both character trajectories and camera-character text descriptions as control and generates trajectories in the global coordinates. On the other hand, Director3D [24]

Dataset	Traj Type	Domain	Caption				Statistics		
			Traj	Scene	Intent	#Vocab	#Sample	#Frame	#Avg (s)
MVImgNet [42]	Object/Scene-Centric	Captured	×	×	×	-	22K	6.5M	10
RealEstate10k [47]	Object/Scene-Centric	Youtube	×	×	×	-	<b>79K</b>	11M	5.5
DL3DV-10K [26]	Object/Scene-Centric	Captured	×	×	×	-	10K	<b>51M</b>	<b>85</b>
CCD [21]	Tracking	Synthetic	✓	×	×	48	25K	4.5M	7.2
E.T. [7]	Tracking	<b>Film</b>	✓	×	×	1790	<b>115K</b>	<b>11M</b>	3.8
DataDoP (Ours)	<b>Free-Moving</b>	<b>Film</b>	✓	✓	✓	<b>8698</b>	29K	<b>11M</b>	<b>14.4</b>

Table 1. **DataDoP Dataset.** We compare the DataDoP dataset to other datasets containing camera trajectories. DataDoP is a large dataset focusing on artistic, free-moving trajectories, each accompanied by high-quality caption annotations. The provided captions detail the camera movements, their interactions with scene content, and the underlying directorial intent. To capture more intricate camera movements, each video clip spans 10-20 seconds, averaging 14.4 seconds.

trains DiT-based framework on object/scene-level multi-view datasets to generate object/scene-centric camera trajectories. NWM [2] employs conditional DiT to plan camera trajectory via agents’ egocentric views. Concurrent work [15] employs an auto-regressive transformer to predict the next frame’s camera movement based on past camera paths and images in aerial videography. Our approach goes further by incorporating both text instructions and RGBD spatial information, enabling precise control in generating camera trajectories for cinematic storytelling.

**Auto-regressive models.** Auto-regressive (AR) modeling employs tokenizers to transform inputs into discrete tokens and formulates generation as a next-token prediction task with transformers. In recent years, great advancements are witnessed in auto-regressive modeling in image [11, 31, 38, 41], video [22, 39, 40], and 3D generation [6, 35, 37]. Early approaches [31, 41] serialize images into patch tokens and train a transformer to auto-regressively model the text and image tokens in a sequential data stream. VAR [38] reformulates auto-regressive image generation as coarse-to-fine next-scale prediction. VideoPoet[22] leverages bidirectional attention for multi-modal input conditioning in auto-regressive video generation. Our work extends auto-regressive modeling to camera trajectory generation controlled by text and geometry cues, leveraging the discrete nature of camera tokens. Compared to diffusion-based methods, our model generates precise, coherent, and intricately detailed artistic trajectories for long camera pose sequences.

### 3. DataDoP Dataset

We introduce **DataDoP**, a camera trajectory dataset extracted from long shots in artistic films, including both movies and documentaries, designed to capture free-moving, intricate, and expressive camera movements. As shown in Fig. 1, each sample in DataDoP consists of a shot-level camera trajectory, accompanied by the corresponding RGBD images and two types of trajectory captions: **Motion** captions, which accurately describe the camera motion alone, and **Directorial** captions, which detail the camera movements, their interaction with the scene, and the directorial intent. We describe the data construction pipeline in

Sec. 3.1 and the dataset statistics in Sec. 3.2.

#### 3.1. Dataset Construction

**Pre-processing.** We curate and filter artistic videos from the internet, which are then segmented into shots using PySceneDetect<sup>1</sup>. Captions are removed using VSR<sup>2</sup>, after which the shots are merged with a publicly available subset from MovieShots [32]. A filtering process is applied to retain shots between 10 and 20 seconds in length, while removing those that are excessively dark or nearly static. Since our dataset focuses on free-moving camera trajectories, which enable unrestricted 3D camera motion within scenes and events, rather than tracking moving people or objects, we specifically filter for this category of data. GPT-4o [18] was used to categorize the shots, removing those with static cameras or object-tracking motion. For details, please refer to Appendix Sec. A.2.

**Trajectory extraction.** We then utilize MonST3R [43] to estimate the geometry of dynamic scenes. Camera trajectories are extracted along with the corresponding depth maps. The trajectories are subsequently cleaned, smoothed, and interpolated into fixed-length sequences.

**Motion tagging.** We then partition the camera trajectories into segments of motion tags. Compared to existing datasets [7, 21], our captions explicitly incorporate descriptions of camera rotation, enabling more fine-grained characterization of camera movements. As a result, our motion tags include both translation and rotation components (see Fig. 2a). For camera translation, excluding the static state, we consider six fundamental motions across three degrees of freedom: lateral (left/static/right), vertical (up/static/down), and depth (forward/static/backward). Each translation motion can be categorized into one, two, or three motions, resulting in a total of 27 possible combinations. For camera rotation, aside from the static state, we consider six fundamental motions across three degrees of freedom: pitch (up/down), yaw (left/right), roll(left/right), resulting in 7 base actions. We do not consider the combination of these rotations, as in practical scenarios, rotation typically involves only one of these basic motions at a time.

<sup>1</sup><https://github.com/Breakthrough/PySceneDetect>

<sup>2</sup><https://github.com/YaoFANGUK/video-subtitle-remover>

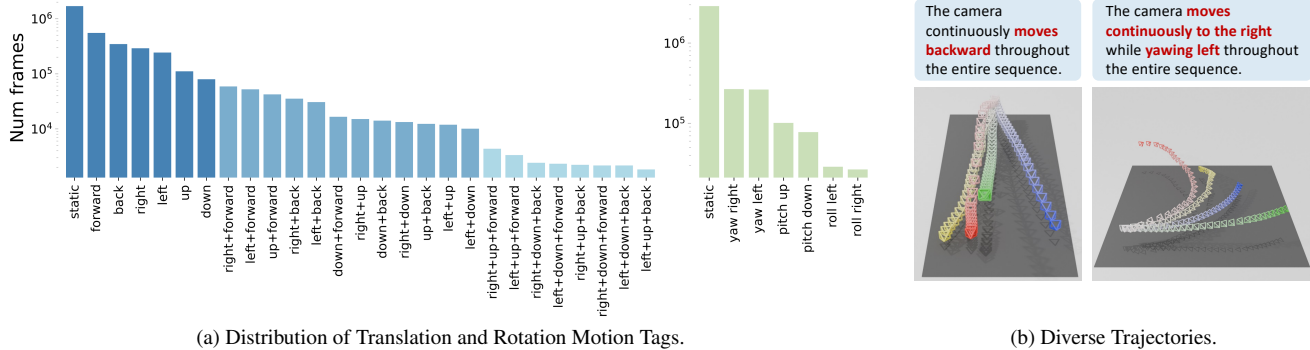


Figure 2. **Dataset Statistics.** (a) The figure illustrates the composition and distribution of 27 translation motions (left) and 7 rotation motions (right), emphasizing the complexity and diversity of trajectories in our DataDoP dataset. (b) Based on the same caption, our dataset includes diverse trajectories that still conform to the given caption. As shown in the figure, the trajectories exhibit variations in terms of length, direction, and speed, effectively showcasing the diversity within our dataset.

Score	Alignment			Quality
	Video-Traj	Traj-Motion	Traj-Directorial	
Acc	0.863	0.913	0.858	0.945
Kappa	0.642	0.530	0.502	0.551

Table 2. **Dataset User Study.** Our user study demonstrates that our dataset exhibits excellent quality and human-alignment, with proven reliability of the results.

We simplify by assuming that camera translation and rotation are completely independent, which results in a total of  $27 \times 7$  possible combinations for camera motion tags.

We adopt the motion tagging method from E.T. [7] to process the camera trajectories. For translation, we use an initial velocity threshold and velocity difference thresholds in different directions to determine the dominant velocity direction combinations. For rotation, we use an initial rotational velocity threshold to identify the unique dominant rotational direction. Finally, we combine the translation and rotation information to generate the complete tags, and apply smoothing to remove noise and sparse tags. These methods provide a coarse temporal description of the camera trajectories, as shown in Fig. 1.

**Caption generation.** Finally, we generate two types of trajectory captions based on the motion tags obtained in the previous stage, as shown in Fig. 1. First, we structure the motion tags by incorporating context, instructions, constraints, and examples, and then leverage GPT-4o to generate **Motion** captions that describe the camera motion alone. Next, we extract 16 evenly spaced frames from the shots to create a  $4 \times 4$  grid and prompt GPT-4o to consider both the previous caption and the image sequence. This enables GPT-4o to generate **Directorial** captions that describe the camera movement, the interaction between the camera and the scene, as well as the directorial intent. Further details can be found in Appendix Sec. A.2.

### 3.2. Dataset Statistics

**Trajectory types.** We classify camera trajectories into four types: *Static*, *Object/Scene-Centric*, *Tracking*, and *Free-*

*Moving*. Static shots keep the camera fixed. Object/Scene-Centric shots capture multi-view data focusing on specific objects or scenes. Tracking shots track a moving subject. Free-Moving shots allow unrestricted 3D camera motion, enabling complex scene exploration and dynamic framing, crucial for cinematic storytelling and creative expression. As shown in Tab. 1, DataDoP stands out by uniquely focusing on artistic, free-moving trajectories, capturing the director’s creative vision and offering significant cinematic and artistic value. Unlike tracking shots, where the camera follows a specific object, free-moving shots fluidly navigate the scene, enhancing visual storytelling without constraints.

**Data scale.** DataDoP is built on long shots from the Internet. As shown in Tab. 1, it consists of 29K samples, spanning 12M frames and totaling 113 hours of footage, all with high-quality trajectory annotations. The dataset focuses on long shots averaging 14.4 seconds, capturing more complex camera movements compared to other datasets. While DL3DV-10K [26] has a longer average duration, its camera trajectories lack directorial intent, emphasizing scene-level consistency rather than creative camera work.

**Statistics.** We present the dataset statistics across four dimensions: *Alignment*, *Quality*, *Complexity*, and *Diversity*. To evaluate Alignment and Quality, we conducted a user study with 8 experts. We selected 100 samples, including original videos, camera trajectories, and two captions: Motion and Directorial. The samples were split into two sets, each labeled by four users. For *Alignment*, we assess the consistency between the trajectory and video (Video-Traj), the motion caption and trajectory (Traj-Motion), and the directorial caption with both the trajectory and video scene (Traj-Directorial). For *Quality*, we assess whether the camera trajectory is free of breaks, roughness, or jitter. We use Fleiss’ Kappa [9] to measure inter-rater agreement among multiple users. As shown in Tab. 2, our dataset achieves high accuracy in both Alignment and Quality, with all Kappa values exceeding 0.4, confirming the reliability of the results. For *Complexity*, as illustrated in Fig. 2a,

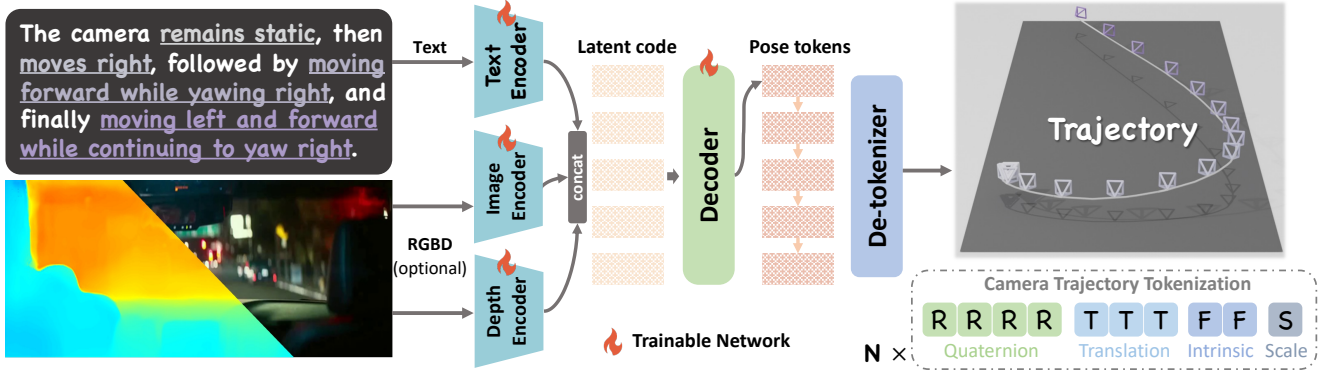


Figure 3. **Our Auto-regressive Generation Model.** Our model supports multi-modal inputs and generates trajectories based on these inputs. By treating the task as an auto-regressive next-token prediction problem, the model sequentially generates trajectories, with each new pose prediction influenced by previous camera states and input conditions.

we present the composition and distribution of motion tags within the dataset. For *Diversity*, as shown in Fig. 2b, the trajectories, while remaining consistent with the caption, exhibit significant variations in length, direction, and speed, effectively showcasing the diversity within our dataset.

## 4. Method

### 4.1. Overview

We introduce **GenDoP** here, an auto-regressive method for camera trajectory generation. Previous trajectory generation methods [2, 16, 21, 24, 25, 37] largely relied on diffusion models [33], which often result in discontinuous and unstable trajectories (See Fig. 4). In contrast, we pioneer the application of auto-regressive models to trajectory generation. Auto-regressive models are well-suited for this task due to their ability to capture sequential dependencies. In trajectory, each pose’s position and orientation depend on the previous one, making the framework ideal for modeling the temporal and spatial continuity of trajectories. By conditioning each pose on its predecessor, the model effectively generates realistic and coherent 3D camera trajectories.

GenDoP automatically constructs the camera’s 3D motion path based on an input caption or appearance and geometry from the initial frame, capturing changes in both position and orientation. As illustrated in Fig. 3, GenDoP takes a text description  $T$ , optionally combined with the initial frame’s RGBD image  $(I_0, D_0)$ , as input and generates the corresponding camera trajectory  $\mathcal{C}$ . A camera trajectory  $\mathcal{C} = \{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{N-1}\}$  is defined as a sequence of  $N$  consecutive camera poses, where each pose  $\mathbf{x}_i = [\mathbf{R}_i | \mathbf{t}_i | \mathbf{K}_i]$  comprises a rotation matrix  $\mathbf{R}_i$  (orientation), a translation vector  $\mathbf{t}_i$  (position), and an intrinsic matrix  $\mathbf{K}_i$  (projection parameters). The intrinsic matrix  $\mathbf{K}_i$  can be simplified to  $(f_x, f_y)$ , assuming a fixed principal point  $(c_x, c_y)$  and image dimensions  $(H, W)$ . Our goal is to derive an auto-regressive generation function  $f$  such that  $\mathcal{C} = f(T, (I_0, D_0))$ . The trajectory tokenization process

is detailed in Sec. 4.2, while the generation method is comprehensively described in Sec. 4.3.

### 4.2. Camera Trajectory Tokenization

Auto-regressive models commonly process information as discrete token sequences, making compact tokenization essential for efficient representation without sacrificing accuracy. Videos are naturally serialized into discrete frames and in this sense, camera trajectories from videos can be easily tokenized into discrete camera pose  $\mathbf{x}_i = [\mathbf{R}_i | \mathbf{t}_i | \mathbf{K}_i]$  at each frame. This simplicity facilitates efficient tokenization, enabling a compact encoding of the trajectory.

**Canonical normalization.** We first establish a scale-invariant trajectory representation via canonical normalization. The camera frame is aligned as the world reference, setting  $\mathbf{R}_0^{\text{norm}} = \mathbf{I}$  and  $\mathbf{t}_0^{\text{norm}} = \mathbf{0}$ . Subsequent poses are relativized through rigid transformation:  $\mathbf{R}_i^{\text{norm}} = \mathbf{R}_0^T \mathbf{R}_i$ ,  $\mathbf{t}_i = \mathbf{R}_0^T (\mathbf{t}_i - \mathbf{t}_0)$  for  $i \in [1, N]$ . Scale normalization then computes  $s = \max_{1 \leq i < N} \|\mathbf{t}_i\|_2$  and projects translations to unit space via  $\mathbf{t}_i^{\text{norm}} = \mathbf{t}_i / (s + \epsilon)$  with  $\epsilon = 10^{-5}$ , maintaining geometric consistency and numerical stability.

**Trajectory tokenization.** For the resulting normalized parameters  $\mathbf{R}_i^{\text{norm}}$  and  $\mathbf{t}_i^{\text{norm}}$ , we compute the corresponding quaternion representation for  $\mathbf{R}_i^{\text{norm}}$  and normalize all parameters to the range  $[0, 1]$ , resulting in the vector  $(r_1, r_2, r_3, r_4, t_1, t_2, t_3)$ . Subsequently, the focal lengths  $f_x, f_y$  and the scale size  $s$  are also normalized, yielding  $(f_1, f_2, s)$ , which are then concatenated with the previously computed values. Finally, these parameters are multiplied by the discrete bin size  $B$  and converted into integer values. Thus, for each  $\mathbf{x}_i$ , we can tokenize it into an integer vector of length 10, where the values are within the range  $[0, B]$ . As a result, each camera trajectory can be tokenized into an integer vector of length  $10N$ .

**Auxiliary tokens.** Similar to prior auto-regressive approaches [6, 35, 37], we prepend a **BOS** token at the beginning of a trajectory sequence, append an **EOS** token at the end, and use **PAD** tokens to fill the necessary positions.

Condition	Method	Dataset	Text-Trajectory Alignment		Trajectory Quality		User Study (AUR)		
			F1-Score $\uparrow$	CLaTr-CLIP $\uparrow$	Coverage $\uparrow$	CLaTr-FID $\downarrow$	Alignment $\uparrow$	Quality $\uparrow$	Complexity $\uparrow$
Motion	CCD [21]	Pre-trained	0.297	5.288	0.332	357.822	2.153	2.033	2.373
	E.T. [7]	Pre-trained	0.315	7.604	0.606	103.799	2.960	3.060	2.953
	Director3D [24]	Pre-trained	0.058	0.000	0.171	542.385	1.407	1.873	1.400
	Director3D [24]	DataDoP	0.391	31.689	0.839	31.979	3.800	3.313	3.520
	<b>GenDoP(Ours)</b>	DataDoP	<b>0.400</b>	<b>36.179</b>	<b>0.872</b>	<b>22.714</b>	<b>4.680</b>	<b>4.720</b>	<b>4.753</b>
Directorial	CCD [21]	Pre-trained	0.315	4.247	0.416	240.216	2.120	2.193	2.287
	E.T. [7]	Pre-trained	0.303	6.127	0.613	94.818	2.813	2.713	2.887
	Director3D [24]	Pre-trained	0.126	0.000	0.348	348.312	1.513	1.747	1.440
	Director3D [24]	DataDoP	0.361	23.505	0.802	35.538	3.913	3.713	3.747
	<b>GenDoP(Ours)</b>	DataDoP	<b>0.399</b>	<b>32.408</b>	<b>0.854</b>	<b>34.275</b>	<b>4.640</b>	<b>4.633</b>	<b>4.640</b>
RGBD & Text	<b>GenDoP(Ours)</b>	DataDoP	<b>0.388</b>	<b>30.231</b>	<b>0.855</b>	<b>33.653</b>	-	-	-

Table 3. **Quantitative Results.** We present the quantitative results of our GenDoP across two text-conditional generation tasks and an RGBD & Text-conditioned task, comparing it with human-tracking methods CCD [21] and E.T. [7], as well as the object/scene-centric method Director3D [24]. Our model consistently outperforms all baselines across all metrics and caption subsets, confirming the effectiveness of both our dataset and auto-regressive framework, positioning GenDoP as a state-of-the-art trajectory generation model.

During the tokenization process, we obtain  $B + 1$  integer values. Consequently, this tokenized representation can be discretized through a learnable codebook  $\mathcal{V} \in \mathbb{R}^{(B+4) \times L}$ , where  $L$  is the latent dimension.

### 4.3. Auto-regressive Generation

We employ a transformer-based auto-regressive architecture to establish a bidirectional mapping between fixed-length camera trajectories and their compact latent representations. Although raw camera trajectories may vary in length, their spatial paths remain consistent after interpolation, allowing us to target fixed-length camera trajectories.

**Text-conditioned encoder.** Our architecture comprises a text encoder  $\mathcal{E}_T$  and an auto-regressive decoder  $\mathcal{D}$  for the base text-conditioned model, as shown in Fig. 3. The text encoder  $\mathcal{E}_T$  utilizes the pretrained and learnable text encoder from Stable Diffusion 2.1 (SD2.1) [33] to extract semantic features, which are then processed through an MLP to generate a textual latent code  $\mathbf{Z}_T \in \mathbb{R}^{M_T \times L}$ , where  $M_T$  denotes the textual latent size and  $L$  is the latent dimension.

**RGBD-conditioned encoder.** For the RGBD-conditioned model, we introduce two separate encoders:  $\mathcal{E}_I$  for RGB image and  $\mathcal{E}_D$  for depth. Specifically, we expand the depth to  $\mathbb{R}^{3 \times H \times W}$  to ensure it can be processed by the encoder. Both encoders use the pretrained and learnable CLIP Vision Model [19, 30, 34] to extract features, which are then passed through MLPs to generate latent codes  $\mathbf{Z}_I \in \mathbb{R}^{M_I \times L}$  and  $\mathbf{Z}_D \in \mathbb{R}^{M_D \times L}$ . The final latent representation is the concatenation of the textual, RGB, and depth codes:

$$\mathbf{Z} = [\mathbf{Z}_T; \mathbf{Z}_I; \mathbf{Z}_D] \in \mathbb{R}^{M \times L}, M = M_T + M_I + M_D. \quad (1)$$

This combined representation integrates both visual and geometric modalities, conditioning the trajectory generation on the accompanying textual information.

**Auto-regressive decoder.** The decoder  $\mathcal{D}$  is an auto-regressive transformer designed to generate a trajectory token sequence from the latent code  $\mathbf{Z}$  and previously token

IDs. We adopt the OPT architecture [45] as the decoder, as utilized in prior works [6, 37]. The latent code  $\mathbf{Z}$  is prepended to the input sequence, positioned before the **BOS** token. For each token prediction, the decoder  $\mathcal{D}$  queries the learnable codebook  $\mathcal{V} \in \mathbb{R}^{(B+4) \times L}$  using the previous token IDs  $\mathbf{y}_{0:P-1}$ , producing the corresponding continuous token embeddings  $V[\mathbf{y}_{0:P-1}] \in \mathbb{R}^{P \times L}$ , where  $P$  denotes the length of the previous token sequence. The input embeddings for the decoder are then computed as:

$$\mathbf{X}_P = \text{PosEmbed}([\mathbf{Z}; V[\mathbf{y}_{0:P-1}]]) \in \mathbb{R}^{(M+P) \times L}. \quad (2)$$

Stacked causal self-attention layers are then employed to predict the next feature based on  $\mathbf{X}_P$ . A linear projection is applied to map the predicted feature to classification logits, which are subsequently used to retrieve the corresponding token ID  $\mathbf{y}_P$ . This process ultimately generates a fixed-length trajectory token sequence.

**Loss function.** The model is optimized with a weighted sum of cross-entropy loss and a regularization term:

$$L = \text{CrossEntropy}(S[1:], \hat{S}[:, -1]) + \lambda \|\mathbf{Z}\|_2^2, \quad (3)$$

where  $S$  is the one-hot ground truth token sequence,  $\hat{S}$  is the predicted logits, and  $\mathbf{Z}$  is the latent code.

## 5. Experiments

### 5.1. Experimental Setting

Our GenDoP framework implements three conditional generation paradigms: (1) **Motion** captions for isolated camera motions, (2) **Directorial** captions for scene-synchronized trajectories, and (3) **RGBD & Text**, a novel approach that integrates images and depth maps with *Directorial* captions through hierarchical feature fusion.

All experiments for both training and inference are carried out with an Intel(R) Xeon(R) Gold 6248R CPU @ 3.00GHz and a single NVIDIA A100-SXM4-80GB GPU. We maintain consistency in parameters and strategies

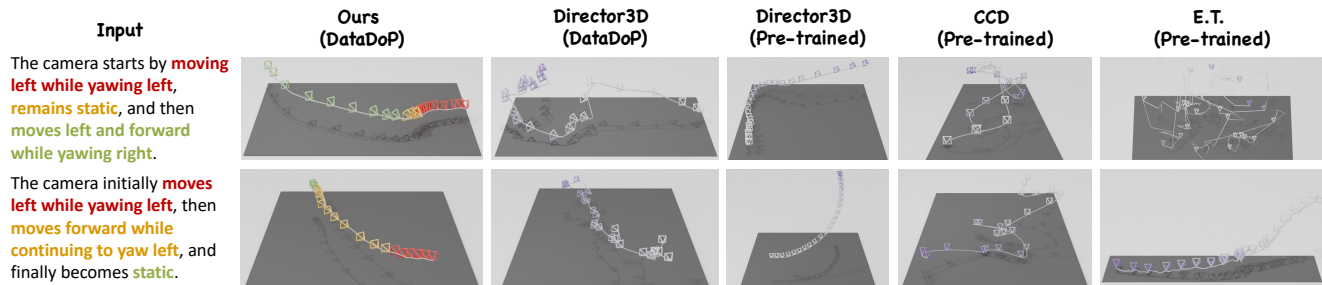


Figure 4. **Qualitative Results of Text-conditioned Trajectory Generation.** We offer a comparative analysis of text-conditioned trajectory generation in the figure. Our model’s trajectories (color-coded to highlight text alignment) remain stable and closely follow the instructions, while other models exhibit significant jitter or fail to match the instructions well.

throughout training to ensure uniformity across the experimental setup. The image resolution is set to  $W = 512$ , with a trajectory length of  $N = 60$ , discrete bin size  $B = 256$ , and latent dimension  $L = 1024$ . The textual latent size is  $M_T = 77$ , the image latent size is  $M_V = 257$ , and the depth latent size is  $M_D = 257$ . The backbone OPT Transformer consists of 12 layers with 12 attention heads each. Training converges after 8 hours on a single A100 GPU, yielding an inference throughput of approximately 3 seconds per trajectory. For evaluation, 3k samples are randomly selected from the DataDoP dataset as the test set, with the remaining data forming the training corpus. Implementation specifics are detailed in Appendix Sec. C.1.

## 5.2. Quantitative Results

**Metrics.** We obtain the Contrastive Language-Trajectory embedding (CLaTr) [7] by leveraging the DataDoP dataset with a CLIP-like approach [30]. A random subset of 3k samples is selected as a test set, from which CLaTr embeddings for both ground truth (GT) and generated data are extracted. Using these embeddings, we evaluate the model with two main metrics. (1) **Text-Trajectory Alignment:** We measure the similarity between text and trajectory embeddings using the CLaTr-CLIP (analogous to CLIP-Score [19]). We also replicate the motion tagging step from Sec. 3.1 to obtain the GT motion tags, which are then compared with the generated tags. Classifier F1-Score is computed by verifying the generated motion tags against the GT labels. (2) **Trajectory Quality:** The alignment between GT and generated trajectories is evaluated using CLaTr-FID (analogous to FID [13]). Additionally, Coverage evaluates how well the generated data spans the range of real data, with higher values indicating a broader representation of the data distribution.

**Results.** We report the quantitative results of our GenDoP across two text-conditional generation paradigms (Motion / Directorial) in Tab. 3. We compare it with previous trajectory generation methods. For the human-tracking methods, CCD [21] and E.T. [7], we assume the character remains static to simplify the camera trajectory inference process.

For the object/scene-centric, text-only conditioned method, Director3D [24], in addition to the pretrained model, we also train a version using DataDoP to emphasize the significance and effectiveness of our dataset for camera trajectory generation tasks. For the RGBD & Text-conditioned task, a novel paradigm introduced by us, we present only the metric results for GenDoP.

Our model demonstrates consistent superiority across all metrics and caption subsets, primarily due to the enhanced trajectory complexity and trajectory-aware captions in our dataset. This innovation enables more precise motion representation, significantly enhancing text-trajectory alignment. This is demonstrated by Director3D models trained on our dataset, which show a dramatic leap in CLaTr-CLIP scores from 0 to over 30, transitioning from object-centric to trajectory-enriched training. Despite sharing the same training data, GenDoP outperforms DataDoP-trained Director3D by 4.5 (Motion) and 9.1 (Directorial) for CLaTr-CLIP, while reducing CLaTr-FID by 9.3 (Motion) and 1.3 (Directorial), as confirmed by user studies. These results validate the effectiveness of our auto-regressive framework. Additionally, GenDoP demonstrates exceptional versatility in handling RGBD & Text-conditional tasks, showing strong multi-modal integration for high-quality trajectory generation under complex constraints. Collectively, the experiments confirm the effectiveness of both our dataset and auto-regressive framework, establishing GenDoP as a state-of-the-art model for trajectory generation.

**User study.** To establish human-aligned evaluation metrics, we engaged 30 domain experts in a user study centered on three critical dimensions: Alignment (trajectory consistency with input text), Quality (smoothness, logical coherence, and seamless connectivity between sequential actions), and Complexity (kinematic sophistication of motion sequences under input constraints). We employed the Average User Ranking (AUR) metric to evaluate model performance, where domain experts assigned ranking scores (1-5) to the five competing models per task. Higher ranking scores indicate superior performance. We comparatively assessed five models on text-conditioned tasks (with 10 sam-

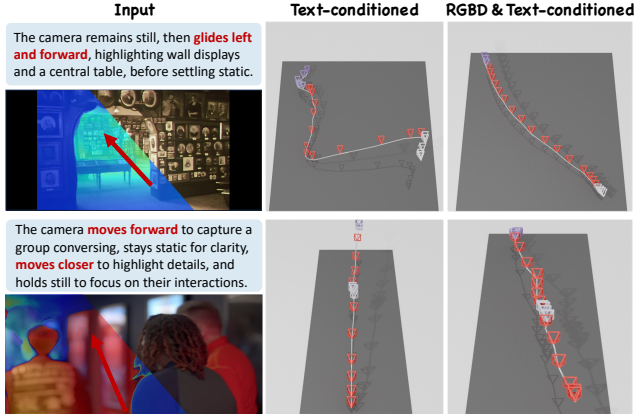


Figure 5. **Qualitative Results of RGBD & Text-conditioned Generation.** This figure compares the impact of incorporating RGBD input on trajectory generation under identical text conditions. While both models generate command-compliant trajectories, the RGBD & Text-conditioned model demonstrates superior scene adaptation by utilizing RGBD data to integrate geometric and contextual constraints.

Ablation		Text-Traj Alignment		Trajectory Quality	
Encoder	Norm	F1-Score	CLaTr-CLIP	Coverage	CLaTr-FID
✓	✓	<b>0.400</b>	<b>36.179</b>	<b>0.872</b>	<b>22.714</b>
✓	×	0.322	14.917	0.766	68.590
×	✓	0.389	31.420	0.866	22.841

Table 4. **Ablation Study.** We conduct an ablation study to evaluate the effectiveness of canonical normalization (see Sec. 4.2) and the trainability of the encoder (see Sec. 4.3).

ples per task), excluding RGBD & Text-conditional scenarios with single-model baselines. As evidenced in Tab. 3, our approach outperformed others across all metrics, with results closely matching the earlier quantitative findings, validating its perceptual and technical coherence.

### 5.3. Qualitative Results

**Text-conditioned generation.** We present comparative analysis of Text-conditioned trajectory generation in Fig. 4. Our model not only achieves superior text-trajectory alignment but also maintains high-quality trajectory generation. Furthermore, the intricate input conditions highlight its capacity to produce sophisticated outputs with high-level complexity. In contrast, the DataDoP-trained Director3D captures basic motion patterns but exhibits trajectory jitter and instability. Furthermore, its object-centric variant pre-trained on [42, 47] generates orbit-dominated trajectories that exhibit no text correspondence, despite improved smoothness. Other baselines exhibit notably inferior performance in both text-trajectory alignment and quality.

**RGBD & Text-conditioned generation.** We conduct a comparative analysis of our trajectory generation model un-

der varying input conditions, as shown in Fig. 5. The results demonstrate that both models generate command-compliant trajectories when given identical textual inputs. However, the RGBD & Text-conditioned model shows superior scene adaptation by leveraging RGBD to incorporate geometric and contextual constraints. Specifically, as shown in the first row of Fig. 5, the spatial information from RGBD effectively mitigates ambiguities, i.e., “left and forward” in the text. This multimodal conditioning enables precise alignment with the 3D scene structure.

### 5.4. Ablation Studies

**Canonical normalization.** We experiment with an alternative strategy that skips canonical normalization, directly using trajectories from Monst3r [43] with scale normalization for tokenization feasibility. These trajectories are scene-centered, with the 3D space focused around the scene. In contrast, canonical normalization transforms them into first-person tracking paths. As shown in the table Tab. 4, applying canonical normalization significantly improves both alignment and quality, providing more consistent camera movements align with the instructions.

**Trainable encoder.** Contrary to conventional practice in text/image-conditional generation where pretrained encoders remain frozen to preserve prior knowledge, our experiments demonstrate comprehensive performance gains (see Tab. 4) by employing trainable encoders. This improvement arises from the encoders’ ability to adapt and bridge cross-modal gaps: through joint optimization, the visual encoder creates geometry-aware trajectory embeddings, while the text encoder learns motion-semantic relationships, resulting in more accurate alignment between text and camera movements.

## 6. Conclusion

We propose **DataDoP**, a pioneering dataset of expressive, free-moving camera trajectories from artistic videos, and **GenDoP**, an auto-regressive multimodal model for trajectory generation. Our approach innovatively incorporates RGBD information as input, enabling spatial data to guide trajectory supervision. This sets a new benchmark, achieving state-of-the-art performance with superior controllability and intent alignment compared to existing methods.

**Limitations and future work.** Currently, our multimodal approach combines text and first-frame RGBD to generate trajectories. Meanwhile, our dataset also extracts 4D point cloud during the extraction process but remains underexplored. Looking ahead, we aim to incorporate more modalities to enhance the adaptability and contextual awareness of the generated trajectories. In addition, we plan to unify trajectory and camera-controlled video creation for iterative creation of both trajectories and video content, establishing a seamless pipeline for automated, artistic film production.

## Acknowledgements

This project is supported by the National Key R&D Program of China (No. 2022ZD0160201) and Shanghai Artificial Intelligence Laboratory, the Ministry of Education, Singapore, under its MOE AcRF Tier 2 (MOE-T2EP20221-0012, MOE-T2EP20223-0002), and under the RIE2020 Industry Alignment Fund – Industry Collaboration Projects (IAF-ICP) Funding Initiative, as well as cash and in-kind contribution from the industry partner(s), Hong Kong RGC TRS T41-603/20-R, the Centre for Perceptual and Interactive Intelligence (CPII) Ltd under the Innovation and Technology Commission (ITC)’s InnoHK. Dahua Lin is a PI of CPII under the InnoHK.

## References

- [1] Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chattopadhyay, Yongxin Chen, Yin Cui, Yifan Ding, Daniel Dworakowski, Jiaojiao Fan, Michele Fenzi, Francesco Ferroni, Sanja Fidler, Dieter Fox, Songwei Ge, Yunhao Ge, Jinwei Gu, Siddharth Gururani, Ethan He, Jiahui Huang, Jacob Samuel Huffman, Pooya Jannaty, Jingyi Jin, Seung Wook Kim, Gergely Klár, Grace Lam, Shiyi Lan, Laura Leal-Taixé, Anqi Li, Zhaoshuo Li, Chen-Hsuan Lin, Tsung-Yi Lin, Huan Ling, Ming-Yu Liu, Xian Liu, Alice Luo, Qianli Ma, Hanzi Mao, Kaichun Mo, Arsalan Mousavian, Seungjun Nah, Sriharsha Niverty, David Page, Despoina Paschalidou, Zeeshan Patel, Lindsey Pavao, Morteza Ramezani, Fitsum Reda, Xiaowei Ren, Vasanth Rao Naik Sabavat, Ed Schmerling, Stella Shi, Bartosz Stefanik, Shitao Tang, Lyne Tchampi, Przemek Tredak, Wei-Cheng Tseng, Jibin Varghese, Hao Wang, Haoxiang Wang, Heng Wang, Ting-Chun Wang, Fangyin Wei, Xinyue Wei, Jay Zhangjie Wu, Jiashu Xu, Wei Yang, Lin Yen-Chen, Xiaohui Zeng, Yu Zeng, Jing Zhang, Qinsheng Zhang, Yuxuan Zhang, Qingqing Zhao, and Artur Zólkowski. Cosmos world foundation model platform for physical AI. *CoRR*, abs/2501.03575, 2025. 2
- [2] Amir Bar, Gaoyue Zhou, Danny Tran, Trevor Darrell, and Yann LeCun. Navigation world models. *CoRR*, abs/2412.03572, 2024. 3, 5
- [3] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, Varun Jampani, and Robin Rombach. Stable video diffusion: Scaling latent video diffusion models to large datasets. *CoRR*, abs/2311.15127, 2023. 2
- [4] James F. Blinn. Where am i? what am I looking at? [cinematography]. *IEEE Computer Graphics and Applications*, 8(4):76–81, 1988. 2
- [5] Rogerio Bonatti, Wenshan Wang, Cherie Ho, Aayush Ahuja, Mirko Gschwindt, Efe Camci, Erdal Kayacan, Sanjiban Choudhury, and Sebastian A. Scherer. Autonomous aerial cinematography in unstructured environments with learned artistic decision-making. *J. Field Robotics*, 37(4):606–641, 2020. 2
- [6] Yiwen Chen, Tong He, Di Huang, Weicai Ye, Sijin Chen, Jiayang Tang, Xin Chen, Zhongang Cai, Lei Yang, Gang Yu, Guosheng Lin, and Chi Zhang. Meshanything: Artist-created mesh generation with autoregressive transformers. *CoRR*, abs/2406.10163, 2024. 3, 5, 6
- [7] Robin Courant, Nicolas Dufour, Xi Wang, Marc Christie, and Vicky Kalogeiton. E.T. the exceptional trajectories: Text-to-camera-trajectory generation with character awareness. In *ECCV (4)*, pages 464–480. Springer, 2024. 2, 3, 4, 6, 7
- [8] Steven Mark Drucker, Tinsley A. Galyean, and David Zeltzer. CINEMA: A system for procedural camera movements. In *SI3D*, pages 67–70. ACM, 1992. 2
- [9] Joseph L. Fleiss. Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76(5):378–382, 1971. 4
- [10] Quentin Galvane, Marc Christie, Christophe Lino, and Rémi Ronfard. Camera-on-rails: automated computation of constrained camera paths. In *MIG*, pages 151–157. ACM, 2015. 2
- [11] Jian Han, Jinlai Liu, Yi Jiang, Bin Yan, Yuqi Zhang, Zehuan Yuan, Bingyue Peng, and Xiaobing Liu. Infinity: Scaling bit-wise autoregressive modeling for high-resolution image synthesis. *CoRR*, abs/2412.04431, 2024. 3
- [12] Hao He, Yinghao Xu, Yuwei Guo, Gordon Wetzstein, Bo Dai, Hongsheng Li, and Ceyuan Yang. Cameractrl: Enabling camera control for text-to-video generation. *arXiv preprint arXiv:2404.02101*, 2024. 1, 2
- [13] 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 *NIPS*, pages 6626–6637, 2017. 7
- [14] Chen Hou, Guoqiang Wei, Yan Zeng, and Zhibo Chen. Training-free camera control for video generation. *CoRR*, abs/2406.10126, 2024. 1, 2
- [15] Yunzhong Hou, Liang Zheng, and Philip Torr. Learning camera movement control from real-world drone videos. *arXiv preprint*, 2024. 3
- [16] Biaozhang Huang, Xinde Li, Chuanfei Hu, and Heqing Li. Stochastic human motion prediction using a quantized conditional diffusion model. *Knowl. Based Syst.*, 309:112823, 2025. 5
- [17] Chong Huang, Chuan-En Lin, Zhenyu Yang, Yan Kong, Peng Chen, Xin Yang, and Kwang-Ting Cheng. Learning to film from professional human motion videos. In *CVPR*, pages 4244–4253. Computer Vision Foundation / IEEE, 2019. 2
- [18] Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Madry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko,

- Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll L. Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, and Dane Sherburn. Gpt-4o system card. *CoRR*, abs/2410.21276, 2024. 3
- [19] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hananeh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, 2021. If you use this software, please cite it as below. 6, 7
- [20] Hongda Jiang, Bin Wang, Xi Wang, Marc Christie, and Baoquan Chen. Example-driven virtual cinematography by learning camera behaviors. *ACM Trans. Graph.*, 39(4):45, 2020. 2
- [21] Hongda Jiang, Xi Wang, Marc Christie, Libin Liu, and Baoquan Chen. Cinematographic camera diffusion model. *Comput. Graph. Forum*, 43(2):i–iii, 2024. 2, 3, 5, 6, 7
- [22] Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Grant Schindler, Rachel Hornung, Vignesh Birodkar, Jimmy Yan, Ming-Chang Chiu, Krishna Somandepalli, Hassan Akbari, Yair Alon, Yong Cheng, Joshua V. Dillon, Agrim Gupta, Meera Hahn, Anja Hauth, David Hendon, Alonso Martinez, David Minnen, Mikhail Sirotenko, Kihyuk Sohn, Xuan Yang, Hartwig Adam, Ming-Hsuan Yang, Irfan Essa, Huisheng Wang, David A. Ross, Bryan Seybold, and Lu Jiang. Videopoet: A large language model for zero-shot video generation. In *ICML*. OpenReview.net, 2024. 3
- [23] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, Kathrina Wu, Qin Lin, Junkun Yuan, Yanxin Long, Aladdin Wang, Andong Wang, Changlin Li, DuoJun Huang, Fang Yang, Hao Tan, Hongmei Wang, Jacob Song, Jiawang Bai, Jianbing Wu, Jinbao Xue, Joey Wang, Kai Wang, Mengyang Liu, Pengyu Li, Shuai Li, Weiyan Wang, Wenqing Yu, Xincheng Deng, Yang Li, Yi Chen, Yutao Cui, Yuanbo Peng, Zhentao Yu, Zhiyu He, Zhiyong Xu, Zixiang Zhou, Zunnan Xu, Yangyu Tao, Qinglin Lu, Songtao Liu, Daquan Zhou, Hongfa Wang, Yong Yang, Di Wang, Yuhong Liu, Jie Jiang, and Caesar Zhong. Hunyuanvideo: A systematic framework for large video generative models. *CoRR*, abs/2412.03603, 2024. 2
- [24] Xinyang Li, Zhangyu Lai, Linning Xu, Yansong Qu, Liujuan Cao, Shengchuan Zhang, Bo Dai, and Rongrong Ji. Director3d: Real-world camera trajectory and 3d scene generation from text. In *NeurIPS*, 2024. 2, 5, 6, 7
- [25] Jing Liang, Amirreza Payandeh, Daeun Song, Xuesu Xiao, and Dinesh Manocha. DTG : Diffusion-based trajectory generation for mapless global navigation. In *IROS*, pages 5340–5347. IEEE, 2024. 5
- [26] Lu Ling, Yichen Sheng, Zhi Tu, Wentian Zhao, Cheng Xin, Kun Wan, Lantao Yu, Qianyu Guo, Zixun Yu, Yawen Lu, Xuanmao Li, Xingpeng Sun, Rohan Ashok, Aniruddha Mukherjee, Hao Kang, Xiangrui Kong, Gang Hua, Tianyi Zhang, Bedrich Benes, and Aniket Bera. DL3DV-10K: A large-scale scene dataset for deep learning-based 3d vision. In *CVPR*, pages 22160–22169. IEEE, 2024. 2, 3, 4
- [27] Christophe Lino and Marc Christie. Intuitive and efficient camera control with the toric space. *ACM Trans. Graph.*, 34(4):82:1–82:12, 2015. 2
- [28] Xinyi Liu, Tianyi Zhang, Matthew Johnson-Roberson, and Weiming Zhi. Splatraj: Camera trajectory generation with semantic gaussian splatting. *CoRR*, abs/2410.06014, 2024. 2
- [29] Koichi Namekata, Sherwin Bahmani, Ziyi Wu, Yash Kant, Igor Gilitschenski, and David B. Lindell. Sg-i2v: Self-guided trajectory control in image-to-video generation. *arXiv preprint arXiv:2411.04989*, 2024. 2
- [30] Alec Radford, Jong Wook Kim, Chris Hallacy, A. 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. In *ICML*, 2021. 6, 7
- [31] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, pages 8821–8831. PMLR, 2021. 3
- [32] Anyi Rao, Jiaze Wang, Linning Xu, Xuekun Jiang, Qingqiu Huang, Bolei Zhou, and Dahua Lin. A unified framework for shot type classification based on subject centric lens. In *ECCV (11)*, pages 17–34. Springer, 2020. 3
- [33] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10674–10685. IEEE, 2022. 5, 6
- [34] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5b: An open large-scale dataset for training next generation image-text models. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022. 6
- [35] Yawar Siddiqui, Antonio Alliegro, Alexey Artemov, Tatiana Tommasi, Daniele Sirigatti, Vladislav Rosov, Angela Dai, and Matthias Nießner. Meshgpt: Generating triangle meshes with decoder-only transformers. In *CVPR*, pages 19615–19625. IEEE, 2024. 3, 5
- [36] Wenqiang Sun, Shuo Chen, Fangfu Liu, Zilong Chen, Yueqi Duan, Jun Zhang, and Yikai Wang. Dimensionx: Create any

- 3d and 4d scenes from a single image with controllable video diffusion. *arXiv preprint arXiv:2411.04928*, 2024. 2
- [37] Jiayang Tang, Zhaoshuo Li, Zekun Hao, Xian Liu, Gang Zeng, Ming-Yu Liu, and Qinsheng Zhang. Edgerunner: Auto-regressive auto-encoder for artistic mesh generation. *CoRR*, abs/2409.18114, 2024. 3, 5, 6
- [38] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling: Scalable image generation via next-scale prediction. In *NeurIPS*, 2024. 3
- [39] Yuqing Wang, Tianwei Xiong, Daquan Zhou, Zhijie Lin, Yang Zhao, Bingyi Kang, Jiashi Feng, and Xihui Liu. Loong: Generating minute-level long videos with autoregressive language models. *CoRR*, abs/2410.02757, 2024. 3
- [40] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. Videogpt: Video generation using VQ-VAE and transformers. *CoRR*, abs/2104.10157, 2021. 3
- [41] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gungjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive models for content-rich text-to-image generation. *Trans. Mach. Learn. Res.*, 2022, 2022. 3
- [42] Xianggang Yu, Mutian Xu, Yidan Zhang, Haolin Liu, Chongjie Ye, Yushuang Wu, Zizheng Yan, Chenming Zhu, Zhangyang Xiong, Tianyou Liang, Guanying Chen, Shuguang Cui, and Xiaoguang Han. Mvimnet: A large-scale dataset of multi-view images. In *CVPR*, pages 915–9161. IEEE, 2023. 2, 3, 8
- [43] Junyi Zhang, Charles Herrmann, Junhwa Hur, Varun Jampani, Trevor Darrell, Forrester Cole, Deqing Sun, and Ming-Hsuan Yang. Monst3r: A simple approach for estimating geometry in the presence of motion. *CoRR*, abs/2410.03825, 2024. 2, 3, 8
- [44] Mengchen Zhang, Tong Wu, Tai Wang, Tengfei Wang, Ziwei Liu, and Dahua Lin. Omni6d: Large-vocabulary 3d object dataset for category-level 6d object pose estimation. In *ECCV (25)*, pages 216–232. Springer, 2024. 2
- [45] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona T. Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. OPT: open pre-trained transformer language models. *CoRR*, abs/2205.01068, 2022. 6
- [46] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all, 2024. 2
- [47] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: learning view synthesis using multiplane images. *ACM Trans. Graph.*, 37(4): 65, 2018. 2, 3, 8