000 CTRL-V: HIGH FIDELITY VIDEO GENERATION WITH 001 **BOUNDING-BOX CONTROLLED OBJECT MOTION** 002 003

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023

024

025

026 027 028 Paper under double-blind review

ABSTRACT

Controllable video generation has attracted significant attention, largely due to advances in video diffusion models. In domains like autonomous driving in particular it can be critical to develop highly accurate predictions for object motions. This paper tackles a crucial challenge of how to exert precise control over object motion for realistic video synthesis in a safety critical setting. To achieve this, we 1) use a separate, specialized model to predict object bounding box trajectories given the past and optionally future locations of bounding boxes, and 2) generate video conditioned on these high quality trajectory predictions. This formulation allows us to test the quality of different model components separately and together. To address the challenges of conditioning video generation on object trajectories in settings where objects may disappear and appear within a scene, we propose an approach based on rendering 2D or 3D boxes as videos. Our method, Ctrl-V, leverages modified and fine-tuned Stable Video Diffusion (SVD) models to solve both trajectory and video generation. Extensive experiments conducted on the KITTI, Virtual-KITTI 2, BDD 100k, and nuScenes datasets validate the effectiveness of our approach in producing realistic and controllable video generation.



047 Figure 1: Overview of Ctrl-V's generation pipeline: Inputs (left): Our inputs include an initial frame, 048 its corresponding bounding-box image and the final frame's bounding-box image. Bounding-box generation samples (middle): We illustrate three different sequences generated from our diffusion based bounding-box motion generation model. Videos sampled from our Box2Videodiffusion model (right): Our Box2Video model conditions on the generated bounding-box videos to produce the final video clips. 052

054 1 INTRODUCTION

055

Recent advances in controllable *image* generation have enabled the creation of highly realistic images from various conditioning inputs, including points, bounding boxes, scribbles, segmentation maps, and skeleton poses. Yet, translating this control to *video* generation is markedly more challenging due to the added temporal dimension. Incorporating time dynamics into diffusion models significantly complicates controllable video generation, as it requires accounting for object interactions, physical consistency, and coherent motion across frames.

062 Numerous recent studies have examined different forms of controllability for video generation. Re-063 searchers have used an array of methods for control, including conditioning on information such as 064 canny edge and depth maps (Zhang et al. (2023b)), similar visual information (Chen et al. (2023)), 065 optical flow (Hu & Xu (2023)), and pose sequences (Karras et al. (2023)). These control inputs are 066 often expensive to produce, especially when sequences of them are required in order to condition a video. Models that use accessible conditioning such as bounding boxes require additional input 067 such as text to help with the generation process (Wang et al. (2024)). A controllable video generation 068 model with an accessible and simple mode of control is greatly desired. 069

In this work, we focus on creating such a model. Specifically, we aim to generate higher fidelity 071 videos controlled by, at the minimum, the beginning and ending positions of 2D and 3D bounding 072 boxes without the help of other modes of control. Our two-part method includes a diffusion-based model that generates the motions and dynamics of objects in the form of bounding box videos (2D 073 images of the bounding boxes evolving over time), and a generative model of videos according to 074 those bounding box videos. To this end, we choose to train and test our model on driving datasets 075 as they contain challenging scenes rich with different types of bounding boxes as well as complex 076 movement and irregular appearing and disappearing objects. In our experiments, we show that our 077 model generates videos that adhere tightly to the desired bounding box motion conditioning, accurately depicting desired object movements. Additionally, through our novel pixel-level bounding 079 box generator and conditioning, our method robustly handles the appearance and disappearance of different objects in a scene, including cars, pedestrians, bikers, and others. 081

In this paper, we present Ctrl-V, a diffusion-based bounding box conditional video generation method that addresses multiple challenges and makes the following contributions to generate higherfidelity videos using diffusion techniques:

- 1. **Bounding box Motion Generations with Diffusion:** We devise a novel diffusion based approach for generating 2D/3D bounding box *trajectories* at the pixel-level (as 2D videos) based on their initial and final states, and the first frame.
- 2. **2D-bounding box and 3D-bounding box Conditioning:** We condition on 2D or 3D Bounding boxes in order to provide a fine-grained control over the generated videos.
- 3. Uninitialized Object Generation: Tracking boxes coordinates outside the current window (boxes that will eventually appear or that are leaving the view) is extremely difficult. With only the first frame, we cannot easily predict these outside-view coordinate movements. This is why, most coordinate-based bounding box generations methods do not account for non-persisting or new bounding boxes (Wang et al., 2024). In this work, we propose a simple solution to this difficult problem: by utilizing on bounding boxes rendered at the pixel-level, we train our model to be sensitive to all bounding boxes, whether present from the first frame or appearing in the middle of the video.
- 4. A New Benchmark for a New Problem Formulation: Given the novelty of our problem formulation, there is no existing standard way to evaluate models that seek to predict vehicle video with high fidelity. We therefore present a new benchmark consisting of a particular way of evaluating video generation models using the KITTI (Geiger et al., 2013), Virtual KITTI 2 (vKITTI) (Cabon et al., 2020), the Berkeley Driving Dataset (BDD 100k) (Yu et al., 2020) and nuScenes (Caesar et al., 2019).
- 103

085

087

2 RELATED WORK

104 105

Video latent diffusion models (VLDMs) extend latent image diffusion techniques (Rombach et al.,
 2022) to video generation. Early VLDMs (Blattmann et al., 2023ba; He et al., 2023; Zeng et al.,
 2023; Wu et al., 2023) shows temporally consistent frame generation and are tailored for text-

prompted or image-prompted video generation. However, these models often struggle with complex scenes and lack the capability for precise local control.

Conditional Video Diffusion techniques providing a certain degree of control. Methods like Video-111 Compose (Wang et al., 2023a), Dreamix (Molad et al., 2023), Pix2Video (Ceylan et al., 2023), and 112 DreamPose (Karras et al.) 2023) propose various designs of novel adapters on top of VLDMs in 113 order to incorporate different conditioning to achieve frame-level control. ControlNet Adapted 114 Video Diffusion, on the other hand, achieve precise regional or pixel-level control in video genera-115 tion by utilizing ControlNet (Zhang et al., 2023a) adapters within VLDM frameworks. Models such 116 as Control-A-Video (Chen et al., 2023), Video ControlNet (Hu & Xu, 2023; Chu et al., 2023), Con-117 trolVideo (Zhang et al., 2023b), and ReVideo (Mou et al.) 2024) show that these adapters are highly 118 adaptable to various types of conditioning, easy to train, and allow for more precise manipulation and enhanced accuracy in editing and creating video content. 119

120 Motion Control with bounding box Conditioning There are many strategies of control that have 121 been explored in controllable video generation research. Notably, ControlVideo (Zhang et al.) 122 2023b) utilizes a training-free strategy that employs pre-trained image LDMs and ControlNets to 123 generate videos based on canny edge and depth maps. Control-A-Video (Chen et al., 2023) leverages a controllable video LDM that combines a pre-trained text-to-video model with ControlNet to 124 manipulate videos using *similar visual information*. Video ControlNets (Hu & Xu, 2023) Chu et al., 125 (2023) uses optical flow information to enhance video generation, while ReVideo (Mou et al., 2024) 126 depends on extracted video trajectories. DreamPose (Karras et al., 2023) injects pose sequence in-127 formation into the initial noise. VideoComposer (Wang et al., 2023a) uses an array of sketch, depth, 128 mask, and motion vectors as conditioning. 129

Many of these conditions, such as edge, depth, and optical flow maps, are costly to produce and lack the flexibility needed for customization. Bounding boxes emerge as a conditioning that are easily customizable and can be edited into different shape, size, locations and classes efficiently. To the best of our knowledge, six other research projects are currently exploring the use of bounding boxes for motion control in video generation. However, it is important to note that our work is distinct from these in several critical respects.

136 Direct-A-Video, TrailBlazer (Ma et al., 2024) and Peekaboo (Jain et al., 2024) are different training-free approaches that employ attention map adjustments to direct the model in generating 137 a particular object within a defined region. Direct-A-Video, in particular, is a text-to-video model 138 that learns to control camera motion during training and then adopts a training-free approach to ma-139 nipulate object movements using bounding boxes. FACTOR (Huang et al., 2023) augmented the 140 transformer-based generation model, Phenaki (Villegas et al., 2022), by integrating a box control 141 module. TrailBlazer, Peekaboo and FACTOR necessitate textual descriptions for individual boxes, 142 thus lacking direct visual grounding. 143

Our task setup shares mild similarities with **Boximator** (Wang et al., 2024) and **TrackDiffu**-144 sion(Fischer et al., 2023) because we also utilize bounding box conditioning during training without 145 relying on text descriptions for individual boxes. However, our approach diverges from these text-146 to-video models, as our primary focus is on generating realistic videos conditioned only on a couple 147 frames of bounding boxes, whereas Boximator and TrackDiffusion are designed to be conditioned 148 on text information as they both are text-to-video models. Boximator and TrackDiffusion enhance 149 their models by introducing new self-attention layers to 3D U-Net blocks. These layers incorpo-150 rate additional conditional information, such as box coordinates and object IDs, into the pretrained 151 VLDM model. Their bounding box information is processed using a Fourier embedder (Mildenhall 152 et al., 2020), which is then passed through multi-layer perceptron layers to encode. In contrast, 153 our approach uses ControlNet and does not involve training additional encoding layers or utilizing Fourier embedder to handle the bounding box information. Moreover, Boximator introduces a 154 self-tracking technique to ensure adherence to the bounding boxes in generated outputs, a technique 155 also adopted by TrackDiffusion. This enables the network to learn the object tracking task along-156 side video generation, but requires a two-stage training process: one with target bounding boxes 157 in frames, and another with the boxes removed. They demonstrate that without this technique, the 158 model's performance markedly declines. Conversely, our model achieves alignment with the bound-159 ing box conditions without additional training. 160

162 Vehicle Oriented Generative Models DriveDreamer (Wang et al., 2023b) presents noteworthy con-163 tribution from autonomous driving domain. It takes an action-based approach to video simulation. 164 It also makes use of bounding boxes and generate actions along with a video rendering. Within the DriveDreamer framework, Fourier embeddings (Mildenhall et al., 2020) are also employed to 165 166 encode bounding box information, and CLIP embeddings (Radford et al., 2021) are used for box categorization. They focus on generating multiple camera views and do not condition on bounding 167 box sequences, so cannot be directly compared with our problem setting. In contrast, the DriveGAN 168 work of Kim et al. (2021) aims to learn a GAN based driving environment in pixel-space, complete with actions and an implicit model of dynamics encoded using the latent space of a VAE. While 170 driving oriented, the approach does not focus on controlling the generation of vehicle video that 171 respects well-defined object trajectories with high fidelity. 172

173 3 OUR METHOD: CTRL-V

3.1 PRELIMINARIES

175

186

191

192

193

194 195

196 197

176 In this section, we provide an overview of the Stable Video Diffusion (SVD) (Blattmann et al., 177 2023ab) model, due to its importance in our approach. SVD is an image-to-video (I2V) model 178 that employs video diffusion. Using an image $f^{(0)}$ as initial condition, SVD is able to extend that single frame into a video $f = [f^{(0)}, \ldots, f^{(N)}]$ where N is the length of the sequence. Notably, 179 180 SVD operates in latent space, where the diffusion and denoising process act upon the latents z of 181 the video f. Here, SVD employs an image encoder (\mathcal{E}) and an image decoder (\mathcal{D}) to translate each 182 frame into and out of latent space: $\mathcal{D}(\mathcal{E}(f^{(i)})) = \mathcal{D}(z^{(i)}) \approx f^{(i)}$. At each diffusion step, SVD 183 progressively introduces noise into the latent representations. In this work, the amount of noise is 184 dictated by Euler discrete noise scheduling method (EDM) introduced in Karras et al. (2022). 185

A UNet based denoiser network within the SVD is used to predict this noise in order to recover the original latent representations. The UNet, \mathbb{U}_{θ} , is parameterized as:

$$\mathbb{U}_{\theta}(\hat{\boldsymbol{z}}_t, \boldsymbol{z}_{\text{pad}}^{(0)}, \boldsymbol{c}^{(0)}, t), \tag{1}$$

• $\hat{z}_t \in \mathbb{R}^{N \times C' \times H' \times W'}$: latent representation of frames corrupted by noise at noise level t.

• $\boldsymbol{z}^{(0)} \in \mathbb{R}^{1 \times C' \times H' \times W'}$: latent representation of the initial frame.

• $\boldsymbol{z}_{pad}^{(0)} \in \mathbb{R}^{N \times C' \times H' \times W'}$: Padded $\boldsymbol{z}^{(0)}$ by repeating itself along the first dimension N times.

• $c^{(0)}$: CLIP encoding (Radford et al., 2021) of the initial frame.

The full denoiser network, \mathbb{D}_{θ} , with an EDM noise scheduler, is formulated as

$$\mathbb{D}_{\theta}(\boldsymbol{z};\boldsymbol{c}^{(0)},\sigma_t) = \lambda_{\text{skip}}(\sigma_t)\boldsymbol{z} + \lambda_{\text{out}}(\sigma_t)\mathbb{U}_{\theta}\left(\lambda_{\text{in}}(\sigma_t)\boldsymbol{z}, \boldsymbol{z}_{\text{pad}}^{(0)}, \boldsymbol{c}^{(0)}; \lambda_{\text{noise}}(\sigma_t)\right)$$
(2)

Here λ_{skip} , λ_{out} , λ_{in} and λ_{noise} denote scaling functions, while σ_t represents the computed noise at level t. The precise mathematical definitions of these terms are detailed in Appendix F

Note that 3D UNet \mathbb{U}_{θ} in Equation [] is a re-parameterized version of the one in Equation 2 (Ronneberger et al., 2015). The scaling terms are absorbed and the inputs are simplified for clarity. In the following sections, we follow the re-parameterized version in Equation [] when refering to the UNets in our model.

205 206

3.2 OVERVIEW OF OUR METHOD: CTRL-V

207
 208 Our controllable video generation method is illustrated in Figure 2 It consists of two sequential steps:

- 210
 1. First, we generate bounding box frames using our diffusion based bounding box predictor net211 work, the **BBox Generator**, which is shown on the left side of Figure 2 These frames contain
 212 only bounding boxes. They make up a video of moving (or stationary) bounding boxes and it
 213 serve as the "skeleton" for the generated video.
- 214 2. Then, we generate a video using our video generator network, Box2Video, shown on the right side of Figure 2 where the bounding boxes frames act as the control signal. The bounding boxes in each frame determine the objects generated in the corresponding frames of the video.

BBox Generator and Box2Video each utilizes a modified SVD backbone – illustrated by the SVD backbone in Figure [2]. These backbones are adapted to their respective generation tasks. Details of each model are presented in their individual sections: BBox Generator – Section [3.3] and Box2Video – Section [3.4].



Figure 2: The diagram illustrates two components of **Ctrl-V**: (left) the **BBox Generator** and (right) **Box2Video**. For both models, we use a **frozen**, **off-the-shelf VAE** to encode images into latent space (\mathcal{E}) and decode them back into pixel space (\mathcal{D}). During training, (1) the **BBox Generator** (Sec. 3.3) learns to denoise the noisy bounding box frame latents \hat{b}_t , conditioned on the first ($b^{(0)}$) and last ($b^{(N-1)}$) bounding box frame latents and the padded initial frame latent $z_{pad}^{(0)}$ and (2) the **Box2Video** (Sec. 3.4) denoises the target frame latents \hat{z}_t by conditioning on the initial frame's latent $z_{pad}^{(0)}$ (input to the SVD UNet) and the bounding box frame latents b (input to the ControlNet).

243 244 245

246

3.3 CTRL-V: BBOX GENERATOR

The BBox Generator shown on the left in Figure 2 aims to predict object bounding boxes across 247 all video frames using an SVD backbone. The four inputs to the model are \hat{b}_t , $b^{(0)}$, $b^{(N-1)}$, $z^{(0)}$, 248 where: \hat{b}_t is the encoded "video" of bounding boxes with t levels of noise added; $b^{(0)}$ is the encoded 249 initial bounding box frame(s); $b^{(N-1)}$ is the encoded final bounding box frame; $z^{(0)}$ is the encoded 250 initial video frame. During training, the model learns to predict the noise added in \hat{b}_t according to 251 the EDM noise scheduler. This allows the recovery of b after subtracting the predicted noise from 252 b_t and passing it through scaling functions. We opt to abstract this detail in the model diagram for 253 readability. 254

In practice, the four inputs are transformed and concatenated into a vector format accepted by the UNet adapter within the SVD backbone. Specifically, as shown in Figure 2 $z^{(0)} \in \mathbb{R}^{1 \times C' \times H' \times W'}$ is replicated along the first dimension, and its front and end (in the first dimension) are replaced by $b^{(0)}, b^{(N-1)}$ respectively. This forms $z_{pad}^{(0)} = \operatorname{concat}(b^{(0)}, z^{(0)}, ..., z^{(0)}, b^{(N-1)}) \in \mathbb{R}^{N \times C' \times H' \times W'}$.

The noise-added encoding of bounding box video \hat{b}_t is then concatenated with $z_{pad}^{(0)}$ to form the final input to the UNet adapter.

The network incorporates additional conditioning inputs, including a CLIP-encoded embedding of the initial frame $c^{(0)}$ and a noise-level embedding t. These embeddings are individually integrated into every sub-block of the U-Net through a self-attention mechanism.

266 3.4 CTRL-V: BOX2VIDEO

267

265

Box2Video is shown on the right in Figure 2 and it aims to generate high-fidelity videos controlled
by bounding box frames, such as those generated by the BBox Generator network. Box2Video consists of an SVD backbone for video generation, and an adapted ControlNet module to process the

bounding box control signal. ControlNet is a widely used network for controlling image generation.
In this work, we modify ControlNet and adapt it to the video diffusion framework (as shown on the right in Figure 2). This architecture allows us to train Box2Video in a single stage without the need for additional optimization criteria, in contrast to previous work such as Boximator and TrackDiffusion (Wang et al.) 2024; [Li et al.] 2024), which require multi-stage learning with extra criteria to train their models.

The SVD component takes two inputs: $z^{(0)}$ and \hat{z}_t . Here, $z^{(0)}$ is the encoded initial video frame and \hat{z}_t is the encoded full video with t levels of noise added to it. As shown in Figure 2 we process these inputs by padding $z^{(0)}$ by repeating it along the first dimension before concatenating it with \hat{z}_t to create the final input to the UNet adapter of the SVD.

The same input is also sent to the ControlNet module through its own UNet adapter layers. Additionally, ControlNet also receives the encoded bounding box frames, *b*, as input, through ControlNet adapter layers. Both of these transformed input is then added together before processed by the ControlNet module. The output signal of the ControlNet module then goes through a zero-convolution before being sent to the SVD UNet decoder layers through residual paths as control signal.

During training, the weights of the SVD model (θ) are frozen, while only the weights in the ControlNet (ξ) are updated.

288 289

3.5 BOUNDING BOX REPRESENTATION

290 The choice of taking bounding boxes information and rendering them out in pixel space is an impor-291 tant detail in Ctrl-V. The method of incorporating bounding box information as a control signal is not 292 trivial. For example, prior work such as Boximator Wang et al. (2024) represents bounding boxes 293 as a Fourier transformed concatenated vector of their raw coordinates, ID and other information. In 294 contrast, in our work we choose to render bounding boxes into frames while maintaining minimal 295 loss of meta information. We encode information such as track ID, object type, and orientation for 296 each bounding box using a combination of visual attributes, including border color, fill color, and markings. Specifically, the *track ID* represents a unique identifier for each tracked object across 297 frames, the object type specifies the category of the object (e.g., car, pedestrian), and the orienta-298 tion indicates the direction the object is facing. Details about how these bounding box frames are 299 rendered can be found in Appendix B.1. Our approach allows us to leverage the highly effective 300 ControlNet approach to provide pixel-level guidance to influence diffusion generated imagery. 301

302 303

304 305

306 307

312 313 314

315 316 317

4 EXPERIMENTAL ANALYSIS AND ABLATION STUDIES



Figure 3: Visualizing video samples generated using the Ctrl-V pipeline: bounding box frame predictions (GB) alongside motion-controlled video generation (GF). *GT:* The ground truth frame sequence from BDD dataset. *GF:* Frames generated based on the predicted bounding box frames. *GB:* Predicted bounding box frames with ground-truth bounding box overlaid+track IDs included.



324 1. The overall visual quality of the generated results (Section 4.2) 325

2. The alignment of the predicted bounding box trajectories with the ground truth (Section 4.3)

3. The fidelity of the generated objects in the video to the bounding box control signal (Section 4.4)

For visual assessment, Figure 3 and Appendix D showcase sample demonstrations generated by our 328 model. To assess video quality, we randomly select 200 initial frames from each dataset's testing 329 set and generate videos. The results in this section are based on analyses of these 200 generated 330 videos per dataset. Furthermore, we explored different bounding box conditioning options: one or three initial bounding box frames, followed by a single final bounding box. Additional variations 332 are discussed in Appendix D.8 333

4.1 DATASETS

326

327

331

334

335 336

337

339

346

348

349

352

353

We evaluate the performance of our models across four autonomous-vehicle datasets: KITTI (Geiger et al., 2013), Virtual KITTI 2 (vKITTI) (Cabon et al., 2020), Berkeley Driving Dataset (BDD) (Yu 338 et al., 2020) with Multi-object Tracking labels (MOT2020), and the nuScenes Dataset (Caesar et al.) 2019).

340 KITTI comprises 22 real-world driving clips with 3D object labelling. vKITTI consists of 5 virtual 341 simulated driving scenes, each offering 6 weather variants, all including 3D object labelling. BDD is 342 a large-scale real-world driving dataset, featuring 1603 2D-labeled sequences of driving videos. The 343 nuScenes dataset is a large-scale driving dataset that includes 1000 scenes 20-second scenes anno-344 tated with 3D bounding boxes, multiple sensor data (lidar, radar and cameras) and map information. Further details on dataset configurations are provided in Appendix B.3. 345

347 4.2 GENERATION QUALITY

To assess the quality of video generation, we compare videos generated through 4 distinct pipelines:

- 350 1. Pre-trained Stable Video Diffusion (SVD) baselines¹ without fine-tuning (initial frame \rightarrow 351 video)
 - 2. Fine-tuned Stable Video Diffusion (SVD) baselines on the provided dataset (initial frame \rightarrow video)
- 354 3. Teacher-forced Box2Video generation (initial frame and all bounding box frames \rightarrow video)
- 355 4. bounding box generation with BBox Generator and Box2Video (initial frame, one or three 356 initial and one last bounding box frames \rightarrow in-between bounding box frames and video).

357 We evaluate our generation across four metrics: Fréchet Video Distance (FVD) (Unterthiner et al., 358 2019), Learned Perceptual Image Patch Similarity (LPIPS) (Zhang et al., 2018), Structural Similarity 359 Index Measure (SSIM) (Wang et al., 2004b) and Peak Signal-to-Noise Ratio (PSNR). These metrics 360 either measure the consistency of frame pixels with the ground truth or the consistency of the frame 361 latents extracted by another network. FVD^2 is an exception; it evaluates the generation distribution against the ground truth's distribution. It is important to note that while many papers report their 362 best-out-of-K results on these metrics, due to computational constraints, we evaluate our model on 363 a single sample for each input. 364

365 The evaluated results are reported in Table 1 and visualizations are available in Appendix D.1. These 366 results indicate that the generation quality improves as we condition on more ground-truth bounding 367 box frames. Details regarding the metrics and their limitations are discussed in Appendix C.1.

368 369

370

4.3 BBOX GENERATOR: QUANTITATIVE EVALUATION

To evaluate the quality of our bounding box generations, we create mask images for both the ground-371 truth and generated bounding box sequences. The mask images are generated by converting the 372

373 ¹Stable Video Diffusion (SVD) baseline is an image-to-video (I2V) model that generates a video sequence 374 conditioned on a single video frame.

375 ²FVD is highly sensitive to video configuration parameters—such as frame rate, clip duration, and spatial 376 resolution-making direct comparisons of FVD values across studies challenging. Additionally, the metric's 377 sensitivity to sample sizes raises concerns, as some datasets may lack sufficient samples for convergence, leading to unreliable estimates.

8			Pipeline	# Cond. BBox	FVD↓	LPIPS↓	SSIM ↑	PSNR ↑
9			Stable Video Diffusion Baseline (Blattmann et al. 2023a)	None	1118.4	0.4575	0.2919	10.63
0			Stable Video Diffusion Fine-tuned (Blattmann et al. 2023a)	None	552.7	0.3504	0.4030	13.01
		È	Ctrl-V: BBox Generator + Box2Video(Ours)	1-to-1	467.7	0.3416	0.3241	13.21
		\mathbf{X}	Ctrl-V: BBox Generator + Box2Video(Ours)	3-to-1	422.2	0.2981	0.4277	13.85
			Ctrl-V: Teacher-forced Box2Video(Ours)	All	435.6	0.2963	0.4394	14.10
			Stable Video Diffusion Baseline (Blattmann et al. 2023a)	None	922.7	0.3636	0.4740	14.61
		E	Stable Video Diffusion Fine-tuned (Blattmann et al., 2023a)	None	331.0	0.2852	0.5540	16.60
		E	Ctrl-V: BBox Generator + Box2Video(Ours)	1-to-1	400.2	0.3179	0.4714	15.78
		٧k	Ctrl-V: BBox Generator + Box2Video(Ours)	3-to-1	341.4	0.2645	0.5841	17.60
			Ctrl-V: Teacher-forced Box2Video(Ours)	All	313.3	0.2372	0.6203	18.41
			Stable Video Diffusion Baseline (Blattmann et al. 2023a)	None	933.6	0.4880	0.3349	12.70
		\cap	Stable Video Diffusion Fine-tuned (Blattmann et al. 2023a)	None	409.0	0.3454	0.5379	16.99
		ā	Ctrl-V: BBox Generator + Box2Video(Ours)	1-to-1	412.8	0.2967	0.5470	17.52
		щ	Ctrl-V: BBox Generator + Box2Video(Ours)	3-to-1	373.1	0.3071	0.5407	17.37
			Ctrl-V: Teacher-forced Box2Video(Ours)	All	348.9	0.2926	0.5836	18.39
			Stable Video Diffusion Baseline (Blattmann et al., 2023a)	None	1179.4	0.5004	0.2877	13.31
		8	Stable Video Diffusion Fine-tuned (Blattmann et al., 2023a)	None	316.6	0.2730	0.4787	18.58
		<u>Vie</u>	Ctrl-V: BBox Generator + Box2Video(Ours)	1-to-1	285.3	0.2647	0.5050	18.93
		<u>-</u>	Ctrl-V: BBox Generator + Box2Video(Ours)	3-to-1	235.0	0.2235	0.5500	20.33
		ng	Ctrl-V: Teacher-forced Box2Video(Ours)	All	235.5	0.2104	0.5705	23.36
	cs	S	DriveGAN (Kim et al., 2021)	None	390.8	-	-	-
	cen		DriveDreamer (Wang et al., 2023b)	All	340.8	-	-	-
	Sut		WoVoGen (Lu et al., 2023)	All	417.7	-	-	-
	ч	3	Drivingdiffusion (Li et al., 2023)	All	332.0	-	-	-
		vie	Drive-WM (Lu et al., 2023)	None	212.5	-	-	-
		Ē	BEV World (Zhang et al., 2024)	None	154.0	-	-	-
		[n]	Panacea (Wen et al., 2024)	All	139.0	-	-	-
		4	Drive-WM (Lu et al., 2023)	All	122.7	-	-	-
			DriveDreamer-2 (Zhao et al., 2024)	None	105.1	-	-	-

401 Table 1: Comparing the quality and diversity of the generated video models. The generated videos 402 consist of 25 frames (except for our nuScenes models which consist of 11 frames videos at 4 Hz) at 403 a resolution of 312×520 , while the reported metrics from this table are evaluated at a resolution of 404 256×410 . The "# Cond. BBox" column reports the number of ground-truth input bounding box 405 frames used by the generation pipelines. "None" indicates that no ground-truth frames are used, 406 while "All" indicates that all ground-truth bounding box frames are utilized. If "# Cond. BBox" is 407 *n*-to-*m*, then it represents the number of initial bounding box frames used by the pipeline is n and the number of final bounding box frames used by the pipeline is m. 408

409 410

bounding box frames into binary masks (details can be found in Appendix C.2). We then calculate 411 the generated averaged mask Intersection over Union (maskIoU) scores, averaged mask Precision 412 (maskP) scores, and averaged mask Recall (maskR) scores against the ground-truth bounding box 413 masks. To assess our bounding box trajectories, we applied the "best-out-of-K" method, selecting 414 the model with the highest maskIoU score for evaluation. In this instance, K equals 5. We compare 415 our results with a baseline referred to as the "Trajeglish-Style" model, an autoregressive GPT-like 416 encoder-decoder that models the bounding box trajectories as a sequence of discrete motion tokens. 417 This baseline is inspired by the work of Philion et al. (2023) with implementation details provided 418 in Appendix E. We present our findings in Table 2, and demonstrate examples of our bounding box generations on each dataset in Appendix D. 419

In the bounding box generation figures, our generator model achieves the closest alignment with the ground-truth in the first and last frames. This near-perfect alignment is primarily attributed to conditioning the model on the bounding boxes of these key frames. When considering all generated frames, the alignment scores decrease, as shown by the plotted demonstrations and metric results in Table 2. This is because objects in frames do not move deterministically. *The role of the bounding box generator is to generate a plausible trajectory for moving objects from the initial bounding box frame to the last.*

Despite the disparity between the ground-truth trajectory and the generated trajectory, our
 Box2Video consistently generates high-fidelity videos based on either trajectory provided. Further
 analysis of this aspect is provided in the subsequent sections.

430 431

4.4 BOX2VIDEO: MOTION CONTROL EVALUATION

	Method	# Cond. BBox	maskIoU↑	maskP↑	maskR↑	maskIoU ↑ (first+last)	maskP ↑ (first+last)	maskR ↑ (first+last)
ITI	BBox Generator (ours) Trajeglish-Style	1-to-1	$\begin{array}{c} \textbf{.629} \pm .212 \\ .447 \pm .154 \end{array}$.758 ± .176 .568 ± .172	$\begin{array}{c} \textbf{.763} \pm .188 \\ .679 \pm .177 \end{array}$	$\begin{array}{c} \textbf{.986} \pm .012 \\ .561 \pm .151 \end{array}$.994 ± .008 .663 ± .150	.992 ± .009 .789 ± .165
Kľ	BBox Generator (ours) Trajeglish-Style	3-to-1	$\begin{array}{c} \textbf{.795} \pm .112 \\ .491 \pm .164 \end{array}$	$\begin{array}{c} \textbf{.881} \pm .082 \\ .622 \pm .173 \end{array}$	$\begin{array}{c} \textbf{.884} \pm .078 \\ .691 \pm .175 \end{array}$	$\begin{array}{c} \textbf{.986} \pm .010 \\ .576 \pm .154 \end{array}$	$.992 \pm .007$ $.684 \pm .149$	$\begin{array}{c} \textbf{.994} \pm .005 \\ .784 \pm .163 \end{array}$
vKITTI	BBox Generator (ours) Trajeglish-Style	1-to-1	$\begin{array}{c} \textbf{.710} \pm .205 \\ .471 \pm .171 \end{array}$	$.828 \pm .178 \\ .578 \pm .200 $.809 ± .171 .700 ± .187	$\begin{array}{c} \textbf{.943} \pm .048 \\ .557 \pm .171 \end{array}$.946 ± .046 .628 ± .194	.997 ± .006 .835 ± .135
	BBox Generator (ours) Trajeglish-Style	3-to-1	.767 ± .131 .520 ± .162	.881 ± .126 .630 ± .186	.853 ± .078 .741 ± .176	$\begin{array}{c} \textbf{.944} \pm .039 \\ .575 \pm .154 \end{array}$	$.948 \pm .036$ $.657 \pm .182$.996 ± .006 .836 ± .143
BDD	BBox Generator (ours) Trajeglish-Style	1-to-1	$\begin{array}{c} \textbf{.587} \pm .214 \\ .305 \pm .183 \end{array}$	$.747 \pm .187 \\ .372 \pm .213$	$.712 \pm .194 \\ .658 \pm .207$	$.954 \pm .047$ $.432 \pm .171$	$.955 \pm .047$ $.483 \pm .192$	$.999 \pm .002$ $.840 \pm .166$
	BBox Generator (ours) Trajeglish-Style	3-to-1	$\begin{array}{c} \textbf{.647} \pm .176 \\ .373 \pm .185 \end{array}$.784 ± .150 .454 ± .206	.783 ± .156 .686 ± .193	.955 ± .043 .492 ± .190	$.955 \pm .042$ $.553 \pm .208$.997 ± .001 .842 ± .154
nuScenes	BBox Generator (ours) Trajeglish-Style	1-to-1	$\begin{array}{c}.364\pm.242\\\textbf{.405}\pm.202\end{array}$	$\begin{array}{c} .433 \pm .278 \\ \textbf{.506} \pm .220 \end{array}$	$\begin{array}{c} \textbf{.740} \pm .186 \\ .661 \pm .216 \end{array}$	$\begin{array}{c} \textbf{.983} \pm .013 \\ .511 \pm .168 \end{array}$	$\begin{array}{c} \textbf{.985} \pm .0112 \\ .603 \pm .172 \end{array}$.997 ± .003 .789 ± .195
	BBox Generator (ours) Trajeglish-Style	3-to-1	.827 ± .150 .448 ± .194	$\begin{array}{c} \textbf{.892} \pm .120 \\ .554 \pm .213 \end{array}$.906 ± .099 .695 ± .196	.983 ± .013 .529 ± .172	$\begin{array}{c} \textbf{.985} \pm .012 \\ .623 \pm .177 \end{array}$.998 ± .003 .791 ± .192

Table 2: Comparing real and generated bounding boxes. We condition on 1 or 3 initial bounding box frame(s) and 1 final bounding box or trajectory frame. The first three columns show evaluations on the entire generated bounding box sequence, measuring the alignment scores between our generated bounding box generations and ground-truth labels. The last three columns focus on testing the auto-encoding capability of the network, evaluating only the first and last frames of the generated sequence. "BBox Generator" is our method and "Trajeglish-Style" is a baseline inspired from Philion et al. (2023) (see Appendix E for implementation details on this baseline).



Figure 4: Illustrations of the generations conditioned on ground truth 3D bounding boxes (2D for BDD) across various datasets. The 2D outlines of bounding boxes are overlayed on top.

Our Box2Video is trained to control object motions through bounding boxes using a teacher-forcing approach, where only ground-truth bounding box frames are provided during the training phase. In this section, we analyze the fidelity of our Box2Video generations to the ground-truth bounding box conditions. To access the consistency of objects' locations between our generated content and ground-truth, we compute the average precision of the bounding boxes in the generated frames and the ground-truth frames.

Average precision (AP) scores gauge the alignment of predicted/generated bounding boxes with the ground-truth labeling. In all related prior studies, average precision (AP) scores have been consis-tently reported. However, it is important to acknowledge that AP scores can vary across studies, depending on the specifics of the task setup. Boximator (Wang et al.) [2024)'s motion control model predicts object locations in the scene, focusing solely on objects with consistent appearances across all frames. Their AP implementation disregards the object locations in the intermediate frames, comparing the objects' locations only in the final frame. In contrast, TrackDiffusion (Li et al., 2024) uses TrackAP for evaluation, employing a QDTrack model (Fischer et al., 2023) to track instances in generated videos and comparing them to ground-truth labels. However, these evaluated datasets

486 have limited instances, and TrackAP requires consistent tracking across frames, making it unsuit-487 able for our project without modifications. Therefore, our AP score differs slightly from those in 488 previous works.

489 Autonomous driving datasets often contain numerous object instances within a scene, with objects 490 continuously entering, exiting, and interacting with each other. In line with this complexity, we 491 have introduced our own version of the AP metric in this work. Our AP metric is designed to 492 comprehensively compare all objects across every scene: encompassing those that newly enter, 493 those that exist during the intermediate frames, and those that overlap with others. 494

First, we utilize the state-of-the-art object detection tool, YOLOv8 (Reis et al., 2024), to obtain the 495 objects' trackings from the generated and ground-truth scenes. Detailed information about the tool 496 and our configurations is reported in Appendix C.3 Next, we match objects in each generated-vs-497 ground-truth frame pair based on *spatial similarity* – calculating the intersection over union (IoU) 498 score to determine the similarity in location between objects' bounding boxes. Our metric disregards 499 object type and tracking IDs equivalence – assuming that objects close in location should naturally 500 have the same type and IDs. Finally, we compute the average precision score following MS COCO 501 protocol (Lin et al.) 2015). Details are provided in Appendix C.4 and results are listed in Table 3. These results indicate that our Box2Video model is particularly adept at adhering to the specified 502 conditions, especially when evaluated with a more lenient metric (i.e., a lower IoU threshold for the AP computation). 504

506	Method	Dataset	Dataset Type	# Frames	mAP↑	$AP_{50}\uparrow$	$AP_{75}\uparrow$	$AP_{90}\uparrow$
507		KITTI	Driving	25	0.547	0.712	0.601	0.327
508	Ctrl-V	vKITTI	Driving-sim	25	0.599	0.776	0.667	0.356
	Curv	BDD	Driving	25	0.685	0.855	0.781	0.401
509		nuScenes	Driving	25	0.661	0.833	0.734	0.381
510	D	MSR-VTT(Xu et al., 2016)	Web videos	16	0.365	0.521	0.384	-
511	Boximator _	ActivityNet (Heilbron et al., 2015)	Human-action	16	0.394	0.607	0.409	-
512	(Wang et al., 2024)	UCF-101 (Soomro et al., 2012)	Human-action	16	0.212	0.343	0.205	-
513	TrackDiffusion	YTVIS (Yang et al., 2019)	YouTube videos	16	0.467	0.656	-	-
514	(Li et al., 2024)	UCF-101 Soomro et al. (2012)	Human-action	16	0.205	0.326	-	-
01-								

Table 3: Average Precision scores obtained by comparing the YOLOv8 bounding box estimations of real and generated samples. Prior works (Wang et al., 2024; Li et al., 2024) do not report results on driving datasets; thus, we draw upon their reported performances on alternative datasets to provide a comparative context. Longer videos are associated with decreased quality and lower detection rates, posing an additional challenge for our model (since it generates 56.25% more frames), yet it obtains higher precision than the other baselines.

5 CONCLUSIONS

524 We present **Ctrl-V**, a novel model capable of generating controllable autonomous vehicle videos via 525 bounding boxes rendering. Our approach demonstrates that the **BBox Generator** model can closely 526 follow generation requirements for the first and last frames and produce a coherent bounding box 527 track for the intermediate frames. Moreover, our Box2Video network generates high-fidelity videos 528 that strictly conform to the provided bounding boxes. Furthermore, our model accommodates both 2D and 3D bounding boxes and handles uninitialized objects appearing in the middle of the videos. 529 Ctrl-V provides future researchers with an efficient way to simulate driving video data with flexible 530 controllability in the form of bounding boxes. In addition, we further define an improved metric to evaluate bounding box conditioned video generation to account for objects that are not present in the 532 first frame, and those that do not remain until the last frame. In Appendix G, we discuss potential 533 future work for this project. With Ctrl-V and an improved metric for more accurate evaluation, we 534 aim to establish a solid foundation for future research in controllable video generation. 535

536

531

505 506 507

515

516

517

518

519

520 521 522

523

REFERENCES 538

Nir Aharon, Roy Orfaig, and Ben-Zion Bobrovsky. Bot-sort: Robust associations multi-pedestrian tracking. arXiv preprint arXiv:2206.14651, 2022.

540 541 542	Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, Varun Jampani, and Robin Rom- bach. Stable video diffusion: Sacling latent video diffusion models to large datasets. 2023a			
543	bach. Stable video diffusion: Scamg fatent video diffusion models to farge datasets, 2025a.			
544	Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidl			
545 546	and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models, 2023b.			
547	Yohann Cabon, Naila Murray, and Martin Humenberger. Virtual kitti 2, 2020.			
548				
549 550 551 552	Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. <i>CoRR</i> , abs/1903.11027, 2019. URL http://arxiv.org/abs/1903.11027.			
553 554 555	Duygu Ceylan, Chun-Hao Paul Huang, and Niloy J. Mitra. Pix2video: Video editing using image diffusion, 2023.			
556 557	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, 2023.			
558 559 560	Ernie Chu, Shuo-Yen Lin, and Jun-Cheng Chen. Video controlnet: Towards temporally consistent synthetic-to-real video translation using conditional image diffusion models, 2023.			
561 562 563	Tobias Fischer, Thomas E. Huang, Jiangmiao Pang, Linlu Qiu, Haofeng Chen, Trevor Darrell, and Fisher Yu. Qdtrack: Quasi-dense similarity learning for appearance-only multiple object tracking, 2023.			
565 566	Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. <i>International Journal of Robotics Research (IJRR)</i> , 2013.			
567 568 569	Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion models for high-fidelity long video generation, 2023.			
570 571 572 573	Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 961–970, 2015. doi: 10.1109/CVPR.2015. 7298698.			
574 575 576	Zhihao Hu and Dong Xu. Videocontrolnet: A motion-guided video-to-video translation framework by using diffusion model with controlnet, 2023.			
577 578	Hsin-Ping Huang, Yu-Chuan Su, Deqing Sun, Lu Jiang, Xuhui Jia, Yukun Zhu, and Ming-Hsuan Yang. Fine-grained controllable video generation via object appearance and context, 2023.			
579 580 581	Yash Jain, Anshul Nasery, Vibhav Vineet, and Harkirat Behl. Peekaboo: Interactive video generation via masked-diffusion, 2024.			
582 583 584	Glenn Jocher, Ayush Chaurasia, and Jing Qiu. Ultralytics YOLO, January 2023. URL https://github.com/ultralytics/ultralytics/			
585 586	Johanna Karras, Aleksander Holynski, Ting-Chun Wang, and Ira Kemelmacher-Shlizerman. Dream- pose: Fashion image-to-video synthesis via stable diffusion, 2023.			
587 588 589	Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion- based generative models, 2022.			
590 591 592	Seung Wook Kim, Jonah Philion, Antonio Torralba, and Sanja Fidler. Drivegan: Towards a control- lable high-quality neural simulation, 2021.			
	Zere Ret 1: 1D. Mite 1: Endine (1) and (1) and (1)			

593 Zoran Kotevski and Pece Mitrevski. Experimental comparison of psnr and ssim metrics for video quality estimation, 01 2010.

594 Pengxiang Li, Kai Chen, Zhili Liu, Ruiyuan Gao, Lanqing Hong, Guo Zhou, Hua Yao, Dit-Yan 595 Yeung, Huchuan Lu, and Xu Jia. Trackdiffusion: Tracklet-conditioned video generation via dif-596 fusion models, 2024. 597 Xiaofan Li, Yifu Zhang, and Xiaoqing Ye. Drivingdiffusion: Layout-guided multi-view driving 598 scene video generation with latent diffusion model. arXiv preprint arXiv:2310.07771, 2023. 600 Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro 601 Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context, 2015. 602 603 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. 604 605 Jiachen Lu, Ze Huang, Jiahui Zhang, Zeyu Yang, and Li Zhang. Wovogen: World volume-aware dif-606 fusion for controllable multi-camera driving scene generation. arXiv preprint arXiv:2312.02934, 2023. 607 608 Wan-Duo Kurt Ma, J. P. Lewis, and W. Bastiaan Kleijn. Trailblazer: Trajectory control for diffusion-609 based video generation, 2024. 610 Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and 611 Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis, 2020. 612 613 Eyal Molad, Eliahu Horwitz, Dani Valevski, Alex Rav Acha, Yossi Matias, Yael Pritch, Yaniv 614 Leviathan, and Yedid Hoshen. Dreamix: Video diffusion models are general video editors, 2023. 615 Chong Mou, Mingdeng Cao, Xintao Wang, Zhaoyang Zhang, Ying Shan, and Jian Zhang. Revideo: 616 Remake a video with motion and content control, 2024. 617 618 Jonah Philion, Xue Bin Peng, and Sanja Fidler. Trajeglish: Learning the language of driving scenar-619 ios. arXiv preprint arXiv.2312.04535, 2023. 620 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-621 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya 622 Sutskever. Learning transferable visual models from natural language supervision. CoRR, 623 abs/2103.00020, 2021. URL https://arxiv.org/abs/2103.00020. 624 625 Dillon Reis, Jordan Kupec, Jacqueline Hong, and Ahmad Daoudi. Real-time flying object detection with yolov8, 2024. 626 627 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-628 resolution image synthesis with latent diffusion models, 2022. 629 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedi-630 cal image segmentation, 2015. 631 632 Mahmood Sharif, Lujo Bauer, and Michael K. Reiter. On the suitability of lp-norms for creating and 633 preventing adversarial examples. CoRR, abs/1802.09653, 2018. URL http://arxiv.org/ 634 abs/1802.09653 635 Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human 636 actions classes from videos in the wild. CoRR, abs/1212.0402, 2012. URL http://arxiv. 637 org/abs/1212.0402 638 639 Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and 640 Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges, 2019. 641 Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang, 642 Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable 643 length video generation from open domain textual description, 2022. 644 Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Ra-645 sul, Mishig Davaadorj, Dhruv Nair, Sayak Paul, William Berman, Yiyi Xu, Steven Liu, and 646 Thomas Wolf. Diffusers: State-of-the-art diffusion models. https://github.com/ 647

huggingface/diffusers, 2022.

648 Jiawei Wang, Yuchen Zhang, Jiaxin Zou, Yan Zeng, Guoqiang Wei, Liping Yuan, and Hang Li. 649 Boximator: Generating rich and controllable motions for video synthesis, 2024. 650 Xiang Wang, Hangjie Yuan, Shiwei Zhang, Dayou Chen, Jiuniu Wang, Yingya Zhang, Yujun Shen, 651 Deli Zhao, and Jingren Zhou. Videocomposer: Compositional video synthesis with motion con-652 trollability, 2023a. 653 654 Xiaofeng Wang, Zheng Zhu, Guan Huang, Xinze Chen, Jiagang Zhu, and Jiwen Lu. Drivedreamer: 655 Towards real-world-driven world models for autonomous driving, 2023b. 656 Zhou Wang and Alan C. Bovik. Mean squared error: Love it or leave it? a new look at signal 657 fidelity measures. IEEE Signal Processing Magazine, 26(1):98–117, 2009. doi: 10.1109/MSP. 658 2008.930649. 659 Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error 661 visibility to structural similarity. IEEE Transactions on Image Processing, 13(4):600-612, 2004a. 662 doi: 10.1109/TIP.2003.819861. 663 Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error 664 visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004b. 665 doi: 10.1109/TIP.2003.819861. 666 667 Yuqing Wen, Yucheng Zhao, Yingfei Liu, Fan Jia, Yanhui Wang, Chong Luo, Chi Zhang, Tiancai Wang, Xiaoyan Sun, and Xiangyu Zhang. Panacea: Panoramic and controllable video generation 668 for autonomous driving. In Proceedings of the IEEE/CVF Conference on Computer Vision and 669 Pattern Recognition, pp. 6902-6912, 2024. 670 671 Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying 672 Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion 673 models for text-to-video generation, 2023. 674 Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging 675 video and language. In 2016 IEEE Conference on Computer Vision and Pattern Recognition 676 (CVPR), pp. 5288–5296, 2016. doi: 10.1109/CVPR.2016.571. 677 678 Linjie Yang, Yuchen Fan, and Ning Xu. Video instance segmentation. CoRR, abs/1905.04804, 2019. 679 URL https://arxiv.org/abs/1905.04804 680 Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht Madha-681 van, and Trevor Darrell. Bdd100k: A diverse driving dataset for heterogeneous multitask learning, 682 2020. 683 684 Yan Zeng, Guoqiang Wei, Jiani Zheng, Jiaxin Zou, Yang Wei, Yuchen Zhang, and Hang Li. Make 685 pixels dance: High-dynamic video generation, 2023. 686 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image 687 diffusion models, 2023a. 688 689 Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable 690 effectiveness of deep features as a perceptual metric, 2018. 691 Yabo Zhang, Yuxiang Wei, Dongsheng Jiang, Xiaopeng Zhang, Wangmeng Zuo, and Qi Tian. Con-692 trolvideo: Training-free controllable text-to-video generation, 2023b. 693 694 Yumeng Zhang, Shi Gong, Kaixin Xiong, Xiaoqing Ye, Xiao Tan, Fan Wang, Jizhou Huang, Hua 695 Wu, and Haifeng Wang. Beyworld: A multimodal world model for autonomous driving via unified bev latent space. arXiv preprint arXiv:2407.05679, 2024. 696 697 Guosheng Zhao, Xiaofeng Wang, Zheng Zhu, Xinze Chen, Guan Huang, Xiaoyi Bao, and Xingang 698 Wang. Drivedreamer-2: Llm-enhanced world models for diverse driving video generation. arXiv 699 preprint arXiv:2403.06845, 2024. 700