
DrivingDojo Dataset: Advancing Interactive and Knowledge-Enriched Driving World Model

Yuqi Wang^{1,2†*} Ke Cheng^{3†} Jiawei He^{1,2†} Qitai Wang^{1,2†}
Hengchen Dai³ Yuntao Chen^{4✉} Fei Xia³ Zhaoxiang Zhang^{1,2,4}

¹ New Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences

² School of Artificial Intelligence, University of Chinese Academy of Sciences

³ Meituan Inc. ⁴ Centre for Artificial Intelligence and Robotics, HKISI, CAS

Project page: <https://drivingdojo.github.io>

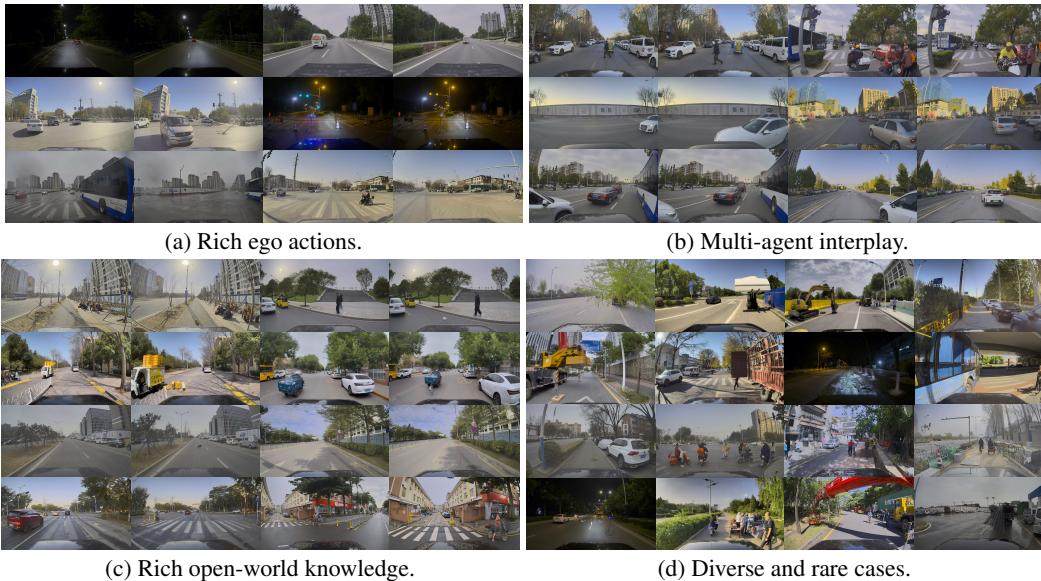


Figure 1: **Examples on DrivingDojo.** (a) showcases various driving actions, such as lane changes, abrupt braking at traffic control, and turning at intersections. (b) illustrates the ego-car’s interactions with other dynamic agents, including cutting-in and cutting-off maneuvers. (c) displays encounters with rolling or falling objects, moving or floating unknown objects, and interactions with traffic lights and boom barriers. (d) presents diverse cases encountered in real-world driving scenarios.

Abstract

1 Driving world models have gained increasing attention due to their ability to model
2 complex physical dynamics. However, their superb modeling capability is yet to
3 be fully unleashed due to the limited video diversity in current driving datasets.
4 We introduce DrivingDojo, the first dataset tailor-made for training interactive
5 world models with complex driving dynamics. Our dataset features video clips
6 with a complete set of driving maneuvers, diverse multi-agent interplay, and rich
7 open-world driving knowledge, laying a stepping stone for future world model
8 development. We further define an action instruction following (AIF) benchmark
9 for world models and demonstrate the superiority of the proposed dataset for
10 generating action-controlled future predictions.

*Work done during an internship at Meituan. † equal contributions. ✉ Corresponding author

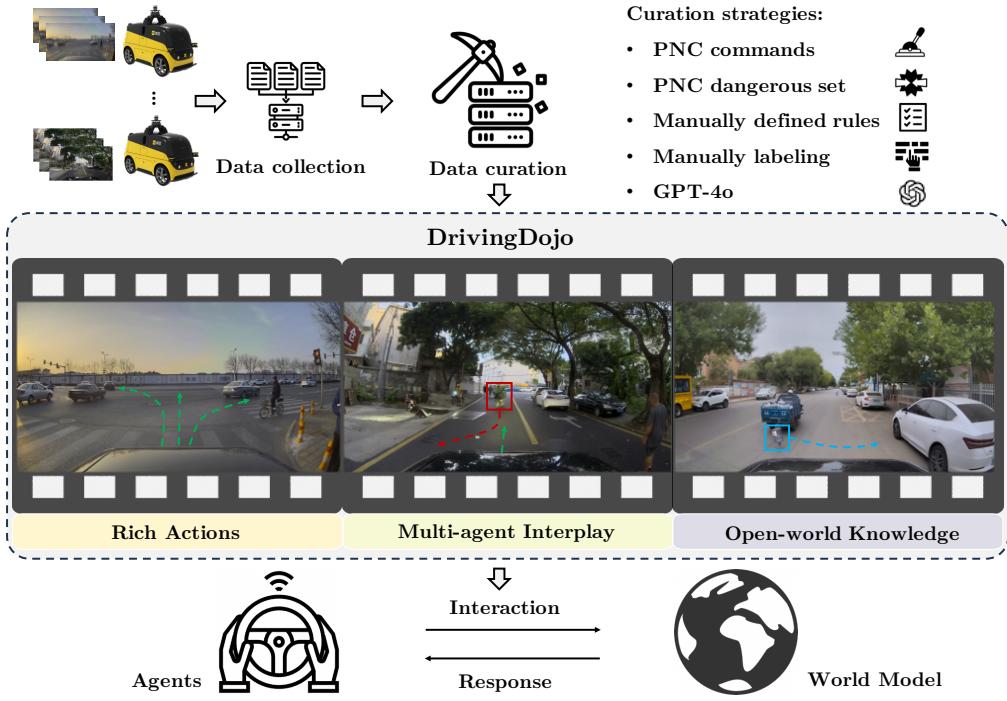


Figure 2: **Enhancing interactive and knowledge-enriched learning of world models.** Data plays a crucial role in modeling the world. DrivingDojo is a large-scale video dataset curated from millions of daily collected videos, designed to investigate real-world visual interactions. DrivingDojo features comprehensive actions, multi-agent interplay, and rich open-world driving knowledge, serving as a superb platform for studying world models.

1 Introduction

World models [16, 19, 32, 20] has received increasing attention due to its ability to model complex real-world physical dynamics and its potential to serve as a general-purpose real-world simulator which is capable of predicting future states in response to diverse action instructions. Facilitated by advancements in video generation techniques [52, 23, 3, 2], models like Sora have achieved remarkable success in producing high-quality videos, thereby opening up a new avenue that treats video generation as real-world dynamics modeling problem [46, 18]. Generative world models, in particular, hold significant promise as real-world simulators and have garnered extensive research in the field of autonomous driving [27, 47, 29, 48, 53, 58].

However, existing driving world models fall short of meeting the requirements of model-based planning in autonomous driving, which aims to improve driving safety in scenarios with diverse ego maneuvers and intricate interaction between the ego vehicle and other road users. These models perform well for non-interactive in-lane maneuvers but have shown limited capability in following more challenging action instructions like lane change. One significant roadblock to building next-generation driving world models lies in the datasets. Autonomous driving datasets commonly used in current world model literature like nuScenes [6], Waymo [44], and ONCE [36], are primarily designed and curated in a perception-oriented manner. As a result, it contains limited driving patterns and multi-agent interactions, which may not fully capture the complexities of real-world driving scenarios. The scarcity of interaction data limits the ability of models to accurately simulate and predict the complex dynamics of real-world driving environments.

In this paper, we propose **DrivingDojo**, a large-scale driving video dataset designed to simulate real-world visual interaction. As illustrated in Figure 1, DrivingDojo features action completeness, multi-agent interplay, and open-world driving knowledge. Our dataset aims to unleash the full action instruction following the ability of world model by including rich lateral maneuvers like acceleration, emergency braking and stop-and-go as well as longitudinal ones like U-turn, overtaking, and lane

Table 1: A comparison of driving datasets for world model. * denotes that the videos are curated from our data pool of around 7500 hours.

Dataset	Videos	Duration (hours)	Ego Trajectory	Complete Actions	Multi-agent Interplay	Open-world Knowledge
nuScenes [6]	1k	5.5	✓			
Waymo [44]	1k	11	✓			
OpenDV-2k [53]	2k	2059		✓		
nuPlan [7]	-	1500	✓	✓	✓	
DrivingDojo (Ours)	18k	150*	✓	✓	✓	✓

36 change. Besides, we explicitly curate the dataset to include a large volume of trajectories containing
 37 multi-agent interplays like cut-in, cut-off, and head-to-head merging. Finally, DrivingDojo taps
 38 into the open-world driving knowledge by including videos containing rare events sampled from
 39 tens of millions of driving video clips, including crossing animals, falling bottles and debris. As
 40 shown in Figure 2, we hope that DrivingDojo could serve as a solid stepping stone for developing
 41 next-generation driving world models.

42 To measure the progress of driving scene modeling, we propose a new action instruction following
 43 (AIF) benchmark to assess the ability of world models to perform plausible future rollouts. The AIF
 44 benchmark measures the visual and structural fidelity of videos generated by world models in an
 45 action-conditioned manner. We propose the AIF errors calculated on the withheld validation data to
 46 evaluate the long-term motion controllability for generated videos. The error is defined as the mean
 47 error between the actions estimated from the generated video and the given action instructions. Then
 48 the baseline world model is evaluated on our DrivingDojo AIF benchmark, for in-domain data and
 49 out-of-domain images or action conditions.

50 Our major contributions are as follows. (1) We design a large-scale driving video dataset to facilitate
 51 research in world model for autonomous driving. Compared to previous datasets in Table 1, our
 52 dataset features complete driving actions, diverse multi-agent interplay, and rich open-world driving
 53 knowledge. (2) We design an action instruction following task for driving world model and provide
 54 corresponding video world model baseline methods. (3) Benchmark results on both driving video
 55 generation and action instruction following show that there are plenty of new opportunities for future
 56 driving world model development on our new dataset.

57 2 Related Works

58 2.1 Autonomous Driving Datasets

59 **Datasets for perception.** The driving dataset has played a crucial role in advancing computer vision
 60 in recent years, aiming to achieve comprehensive perception and understanding surrounding the ego
 61 vehicle. Initially, perception in autonomous driving relied on 2D image-based perception. Datasets
 62 like Cityscapes [10], Mapillary Vistas [38], and BDD100k [56] provided instance-level masks for
 63 learning tasks. With the integration of LiDAR sensors and advancements in 3D perception, datasets
 64 like KITTI [13], nuScenes [6], and Waymo [44] have emerged as standard benchmarks for various
 65 3D perception tasks. Additionally, datasets like ONCE [36], Argoverse [8, 49], and others [28, 14, 1]
 66 are also utilized for studying various perception tasks.

67 **Datasets for prediction and planning.** In recent years, there's been increasing attention on
 68 prediction and planning in autonomous driving. Prediction involves anticipating the behavior of other
 69 agents, while planning relates to the behavior of the ego vehicle. Prediction methods typically rely on
 70 semantic maps and dynamic traffic light statuses to anticipate future vehicle motions. Notable datasets
 71 in this area include Argoverse Motion Forecasting [8], Waymo Open Motion Dataset [12], Lyft Level
 72 5 Prediction Dataset [25], and nuScenes Prediction [6] challenge. Additionally, the Interaction
 73 dataset [57] provides interactive driving scenarios with semantic maps derived from drones and traffic
 74 cameras, enriching the understanding of complex driving interactions. Transitioning to planning,
 75 CARLA [11] stands out as an open-source simulator designed to simulate real-world traffic scenarios,

76 providing a platform for testing and validating planning algorithms. Complementing this, nuPlan [7]
 77 introduces the first closed-loop planning benchmark for autonomous vehicles, closely mirroring
 78 real-world scenarios.

79 2.2 World Model

80 **Learning world models.** World models [16, 32] enable next-frame prediction based on action
 81 inputs, aiming to build general simulators of the physical world. However, learning dynamic modeling
 82 in pixel space is challenging, leading previous image-based world models to focus on simplistic
 83 gaming environments or simulations [17, 19, 9, 51, 43, 42, 20]. With advances in video generation,
 84 models like Sora can now produce high-definition videos up to one minute long with natural, coherent
 85 dynamics. This progress has encouraged researchers to explore world models in real-world scenarios.
 86 DayDreamer [50] applies the Dreamer algorithm to four robots, allowing them to learn online and
 87 directly in the real world without simulators, demonstrating that world models can facilitate faster
 88 learning on physical robots. Genie [5] demonstrates interactive generation capabilities using vast
 89 internet gaming videos and shows potential for robotics applications. UniSim [54] aims to create a
 90 universal simulator for real-world interactions using generative modeling, with applications extending
 91 to real-robot executions.

92 **World model for autonomous driving.** World models serving as real-world simulators have
 93 garnered widespread attention [15] and can be categorized into two main branches. The first branch
 94 explores agent policies in virtual simulators. MILE [26] employed imitation learning to jointly learn
 95 the dynamics model and driving behavior in CARLA [11]. Think2Drive [33] proposed a model-based
 96 RL method in CARLA v2, using a world model to learn environment transitions and acting as
 97 a neural simulator to train the planner. The second branch focuses on simulating and generating
 98 real-world driving scenarios. GAIA-1 [27] introduced a generative world model for autonomous
 99 driving, capable of simulating realistic driving videos from inputs like images, texts, and actions.
 100 DriveDreamer [47] emphasized scenario generation, leveraging HD maps and 3D boxes to enhance
 101 video quality. Drive-WM [48] was the first to propose a multiview world model for generating high-
 102 quality, controllable multiview videos, exploring applications in end-to-end planning. ADriver-I [29]
 103 constructed a general world model based on MLLM and diffusion models, using vision-action pairs
 104 to auto-regressively predict current frame control signals. DriveDreamer2 [58] leveraged LLMs and
 105 text prompts to generate diverse driving videos in a user-friendly manner. Unlike previous methods
 106 that focused on model design, OpenDV-2K [53] addressed the issue of training data by collecting
 107 over 2000 hours of driving videos from the internet. Previous research has predominantly addressed
 108 static scene generation, with limited emphasis on multi-agent interplays. Our dataset enables the
 109 exploration of world model predictions within dynamic, interactive driving scenarios.

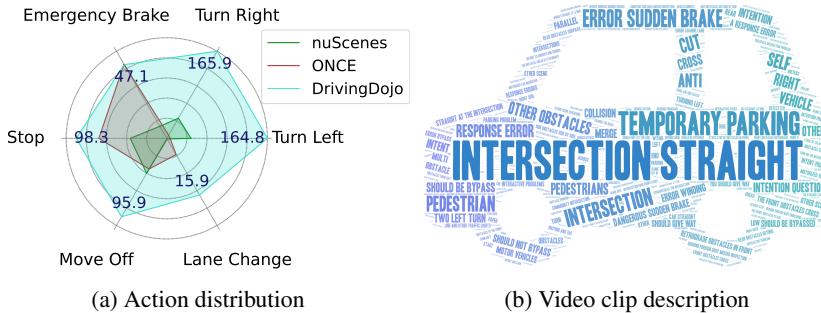


Figure 3: **The strengths of the DrivingDojo dataset.** (a) illustrates a comparison of action distributions among nuScenes, ONCE, and our DrivingDojo. We compare the average hourly event counts of driving actions. (b) presents the distribution of text descriptions for the video clips in DrivingDojo.

110 3 The DrivingDojo Dataset

111 Our goal is to provide a large and diverse action-instructed driving video dataset DrivingDojo to
 112 support the development of driving world models. To accomplish this, we extract highly informative

113 clips from a video pool collected through fleet data, spanning several years and comprising more than
114 500 operating vehicles across multiple major Chinese cities. As a result, our DrivingDojo features
115 diverse ego actions, rich interactions with road users, and rare driving knowledge which are crucial
116 for high-quality future forecasting as shown in Table 2.

117 We begin with the design principles of DrivingDojo and its uniqueness compared with existing
118 datasets in Section 3.1- 3.3. We then describe the data curation procedure and statistics in Section 3.4.
119 Here, we only describe the design principles. More detailed information refer to the Appendix.

Table 2: **DrivingDojo constitution.** We also organize the dataset into three subsets: DrivingDojo-Action, DrivingDojo-Interplay, and DrivingDojo-Open, to facilitate research on specific tasks.

Dataset	Videos	Type	Camera	Ego Trajectory	Text Description
DrivingDojo	18.2k	total	✓	✓	✓
DrivingDojo-Action	7.9k	rich ego-actions	✓	✓	
DrivingDojo-Interplay	6.4k	multi-agent interplay	✓	✓	
DrivingDojo-Open	3.9k	open-world knowledge	✓	✓	✓

120 3.1 Action Completeness

121 Using the driving world model as a real-world simulator requires it to follow action prompts accurately.
122 Existing autonomous driving datasets, such as ONCE [36] and nuScenes [6], are generally curated
123 for developing perception algorithms and thus lack diverse driving maneuvering.

124 To enable the world model to generate an infinite number of high-fidelity, action-controllable virtual
125 driving environments, we create a subset called DrivingDojo-Action that features a balanced
126 distribution of driving maneuvers. This subset includes a diverse range of both longitudinal
127 maneuvers, such as acceleration, deceleration, emergency braking, and stop-and-go driving, as well as
128 lateral maneuvers, including lane-changing and lane-keeping. As demonstrated in Figure 3a, our
129 DrivingDojo-Action subset offers a significantly more balanced and complete set of ego actions
130 compared to existing autonomous driving datasets.

131 3.2 Multi-agent Interplay

132 Besides navigating in a static road network environment, modeling the dynamics of multi-agent
133 interplay like merge and yield is also a crucial task for world models. However, current datasets
134 are either built without considering multi-agent interplays, such as nuScenes [6] and Waymo [44],
135 or are constructed from large-scale internet videos that lack proper curation and balancing, like
136 OpenDV-2K [53].

137 To address this issue, we design the DrivingDojo-Interplay subset focusing on interactions with
138 dynamic agents as a core component of the dataset. As shown in Figure 1b, we curate this subset to
139 include at least one of the following driving scenarios: cutting in/off, meeting, blocked, overtaking,
140 and being overtaken. These scenarios encompass a variety of realistic situations, such as vehicles
141 cutting into lanes, encounters with oncoming traffic, and the necessity for emergency braking. By
142 incorporating these diverse scenarios, our dataset enables world models to better understand and
143 anticipate complex interactions with dynamic agents, thereby improving their performance in real-
144 world driving conditions.

145 3.3 Rich Open-world Knowledge

146 In contrast to perception and prediction models, which compress high-dimensional sensor input
147 into low-dimensional vector representations, world models exhibit a superior modeling capacity by
148 operating in the pixel space. This increased capacity enables world models to effectively capture the
149 intricate dynamics of open-world driving scenarios, such as animals unexpectedly crossing the road
150 or parcels falling off the trunks of vehicles.

151 However, existing datasets, either perception-oriented ONCE [36] or planning-oriented ones like
152 nuPlan [7], do not have adequate data for developing and assessing the long-tail knowledge modeling

ability of world models. Therefore, we place a unique emphasis on including rich open-world knowledge video clips and construct the DrivingDojo-Open subset. As shown in Figure 1c, describing open-world driving knowledge like this is challenging due to its complexity and variability, but these scenarios are crucial for ensuring safe driving.

The DrivingDojo-Open subset consists of 3.9k video clips about the open-world knowledge in driving scenarios. This subset is curated from fleet data that includes unusual weather, foreign objects on the road surface, floating obstacles, falling objects, taking over cases, and interactions with traffic lights and boom barriers. A word cloud of video descriptions for DrivingDojo-Open are shown in Figure 3b. DrivingDojo-Open serves as an invaluable supplementary for driving world modeling by including driving knowledge beyond simply interacting with structured road networks and other regular road users.

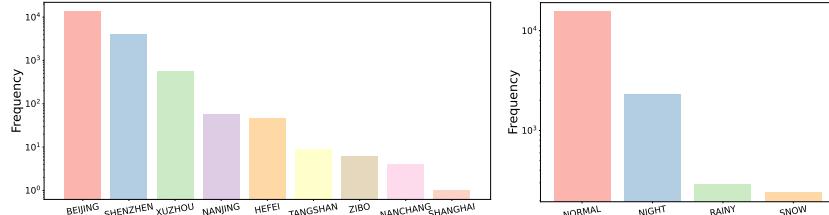


Figure 4: **Descriptive statistics of the DrivingDojo dataset.** The dataset was collected from various regions across China, including nighttime and rainy/snowy conditions.

3.4 Data Curation and Statistics

Dataset statistics. The DrivingDojo dataset contains 18k videos with resolution of 1920×1080 and frame rate at 5 fps. Our video clips are collected from major Chinese cities including Beijing, Shenzhen, Xuzhou, etc., as shown in Figure 4. Furthermore, these videos are recorded in diverse weather conditions at different daylight conditions. All videos are paired with synced camera poses derived from the HD-Map powered high precision localization stack onboard. Videos in the DrivingDojo-Open subset are paired with text descriptions about the rare event happening in each video. More details are in the Appendix.

Data collection. We collected multi-modal fleet data using the platform of Meituan’s autonomous delivery vehicles. Our dataset consists of video clips recorded by the front-view camera with a horizontal field of view of 120° to capture comprehensive visual information. The raw data is collected from multiple Chinese cities between May 2022 and May 2024, amassing a total of 900,000 videos and approximately 7,500 hours of driving footage pre-filtered before recording.

Data curation. In order to ensure both the data diversity as well as balanced ego action and multi-agent interplay distribution, we include fleet data with different criteria. The data sources of DrivingDojo include 1) intervention data from safety inspectors during vehicle operation, 2) emergency brake data from automatic emergency braking, 3) randomly sampled 30-second general videos from collected videos, 4) selected distinct scenarios such as traffic light changes, barrier opening, left and right turns, straight crossings, vehicle encounters, lane changes, and pedestrian interactions, 5) manually sorted rare data containing moving and static foreign objects on the road, floating obstacles, falling and rolling objects. The curation details are in the Appendix.

Personal Identification Information (PII) removal. To avoid privacy infringement and obey the regulation laws, we employ a high precision license plate and face detectors [30] to detect and blur these PII for each frame of all videos. An in-house annotation team and the authors have manually double-checked that the PII removal procedure is correctly carried out for all the videos.

4 DrivingDojo for World Model

To facilitate the study of world models in autonomous driving, we define a novel action instruction following (AIF) task. We provide baseline methods (Section 4.2) and evaluation metrics (Section 4.3), enabling further investigations. More details are described in the Appendix.

193 **4.1 Action Instruction Following**

194 Action-controllable video forecasting is the core ability of world models [5]. Instead of solely focusing
 195 on predicting high-quality video frames, action instruction following requires world models to take
 196 both the initial video frame and ego action prompts into consideration for predicting corresponding
 197 world responses. Given the initial image I_t and a sequence of actions $\{A_t, \dots, A_{t+k}\}$, the model f_θ
 198 predicts future states $\{I_{t+1}, \dots, I_{t+k}\}$ as:

$$\{I_t, \dots, I_{t+k}\} = f_\theta(I_t, \{A_t, \dots, A_{t+k}\}). \quad (1)$$

199 Here, $\{A_t, \dots, A_{t+k}\}$ refers to the action prompts for each frame, with trajectories $A_t = (\Delta x_t, \Delta y_t)$
 200 in our experiment. f_θ represents the world model, and $\{I_{t+1}, \dots, I_{t+k}\}$ signifies the visual prediction
 201 for subsequent k frames.

202 **4.2 Model Architecture**

203 We propose DrivingDojo baseline, a video generation model based on Stable Video Diffusion
 204 (SVD) [2]. While SVD is a latent diffusion model for image-to-video generation, we extend its
 205 capability to generate videos conditioned on action. For the AIF task, we encode the value of each
 206 action sequence into a 1024-dimensional vector using a Multilayer Perceptron (MLP). Subsequently,
 207 the action feature is concatenated with the first-frame image feature and passed into the U-Net [39].

208 **4.3 Evaluation Metrics**

209 **Visual quality.** To evaluate the quality of the generated video, we utilize FID (Frechet Inception
 210 Distance) [22] and FVD (Frechet Video Distance) [45] as the main metrics.

211 **Action instruction following.** We propose the action instruction following (AIF) errors E_x^{AIF}
 212 and E_y^{AIF} to measure the consistency between the generated video and the input action conditions.
 213 Given the generated video sequences $\{I_t, \dots, I_{t+k}\}$, we estimate vehicle trajectories in the generated
 214 videos with the offline visual structure-from-motion (SfM) implementation like COLMAP [40, 41]:
 215 $\{\tilde{A}_t, \dots, \tilde{A}_{t+k}\} = \text{SfM}(\{I_t, \dots, I_{t+k}\})$, where $\{\tilde{A}_t, \dots, \tilde{A}_{t+k}\}$ are estimated trajectories of unknown
 216 scale. We estimated the scale factor \hat{S} for the predicted trajectory by minimizing the error between
 217 estimated and input ego-motion in the first N frames. We compare the estimated actions with the
 218 ground-truth action instructions $\{A_t, \dots, A_{t+k}\}$ and report the mean absolute error for both lateral
 219 (E_y^{AIF}) and longitudinal (E_x^{AIF}) actions:

$$(E_x^{\text{AIF}}, E_y^{\text{AIF}}) = \frac{\sum_{i=0}^k |A_{t+i} - \tilde{A}_{t+i} * \hat{S}|}{k+1}, \quad (2)$$

220 where the scale factor $\hat{S} = \arg \min_S \sum_{i=0}^N |A_{t+i} - \tilde{A}_{t+i} * S|$.

Table 3: **Comparison of visual prediction fine-tuning across different datasets.**, † indicates using camera sweeps data. The performance is zero-shot evaluated on the OpenDV-2K dataset.

Method	Fine-tuning	Evaluation	FID	FVD
SVD [2]	-	OpenDV-2K	24.17	580.94
SVD [2]	nuScenes†	OpenDV-2K	21.05	395.04
SVD [2]	DrivingDojo	OpenDV-2K	19.20	343.91

221 **5 Experiments**

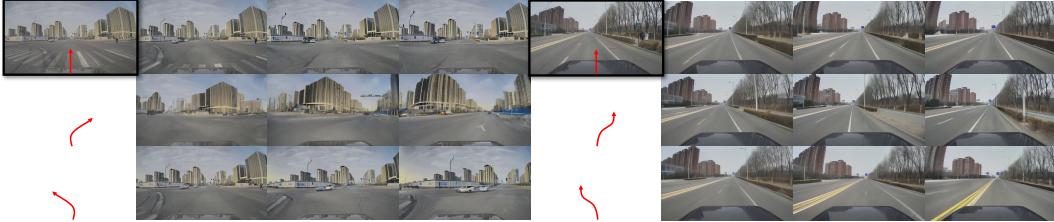
222 **5.1 Results of Visual Prediction**

223 To illustrate the richness of behaviors and dynamics within our dataset, we compare video fine-tuning
 224 quality across various datasets. In Table 3, we random selected 256 video segments from the OpenDV-

225 2K dataset [53] as our test set and evaluated fine-tuning performance of SVD [2] model across various
 226 datasets. The results indicate that models trained on our dataset exhibit better visual quality.

227 5.2 Results of Action Instruction Following

228 **Diverse driving behaviors.** Based on different sequences of actions, our model is able to generate
 229 multiple possible futures. As shown in Figure 5, we showcase the model’s capability to execute
 230 forward, left turn, and right turn maneuvers at intersections, as well as lane-changing to the left or
 231 right, and maintaining on straight roads.



232 **Figure 5: Predicting multiple futures based on different actions.** Left: going straight, turning
 233 left, and turning right at a crossing; Right: changing to the left lane, staying in the current lane, and
 234 changing to the right lane.

235 **Action instruction following.** Although qualitative evaluations demonstrate the powerful generative
 236 ability of our model, we also endeavor to measure the accuracy of action instruction following
 237 quantitatively. We seek to evaluate whether the video trajectories generated by the model closely
 238 adhere to our expected route paths. This serves as a fundamental assurance for the future application
 239 of world model. As shown in Table 4, with the in-domain actions (original action sequences of the
 240 test video) as conditions, videos generated by the baseline world model trained on DrivingDojo
 241 exhibit strong loyalty towards the action instructions. The mean action error in each video frame
 242 is limited to only 6-7 cm in the lateral or longitudinal directions. In row 2, feeding the model with
 243 the same initial images and randomly sampled action instructions slightly increases the mean action
 244 errors. When the model is applied zero-shot to initial images from OpenDV-2K [53] and fed with
 245 randomly sampled action instructions, its generated videos still demonstrate considerable consistency
 246 to the action instructions. Note that the proposed action instruction following errors can sensitively
 247 reflect the impact of out-of-domain inputs on the performance of the model.

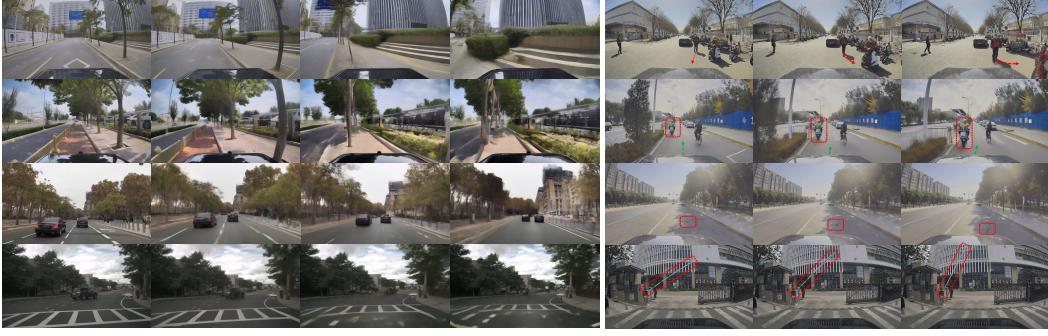
Table 4: **Action instruction following on the DrivingDojo dataset.** * denotes the model is applied
 zero-shot to this dataset without fine-tuning.

Action Type	Test Dataset	FID	FVD	$E_x^{\text{AIF}}(\downarrow)$	$E_y^{\text{AIF}}(\downarrow)$
In-Domain	DrivingDojo	37.07	658.72	0.068m	0.060m
Out-of-Domain	DrivingDojo	38.30	716.44	0.133m	0.099m
Out-of-Domain	OpenDV-2K*	24.27	442.67	0.175m	0.121m

245 5.3 Real-world Simulation

246 **Action generalization.** Our model demonstrates robust generalization capabilities in two key aspects.
 247 As illustrated in Figure 6a, firstly, it effectively generalizes to out-of-domain (OOD) actions, such
 248 as forcefully driving on pedestrian walkways, showcasing its adaptability to some unreasonable
 249 actions. Secondly, it successfully extends its capabilities to other datasets, executing tasks such as
 250 lane changes on the OpenDV-2K [53] dataset and backing-the-car maneuvers on the nuScenes [6]
 251 dataset without requiring further fine-tuning. This underscores the model’s potential as a real-world
 252 simulator, capable of adapting to diverse driving scenarios.

253 **Dynamic agents.** We showcase our model’s ability to simulate interactions with dynamic agents in
 254 Figure 6b. The results indicate that the model can provide reasonable responses based on our actions.



(a) Action generalization

(b) Interaction simulation

Figure 6: Qualitative examples of our model’s capability.

255 The first scenario depicts a pedestrian opting to yield as our vehicle continues forward, resulting in a
 256 change in trajectory. In the second scenario, a delivery person opts to stop and wait at a narrow road.

257 **Open-world dynamics.** In Figure 6b, our model showcases the simulations of rare scenarios
 258 encountered on the road, including interactions with moving birds and parking lot barriers.

259 5.4 Limitations and Future Work

260 This paper primarily delves into the value of the dataset, with the model aspect serving as a base-
 261 line without any special design. Although the DrivingDojo dataset significantly improves model
 262 capabilities, there are still several limitations that require further investigation in future studies.

263 **Hallucination.** As shown in Figure 7, we observed that the model exhibits some hallucinations, such
 264 as the sudden disappearance of objects, and when an action is unrealistic given the scene, such as
 forcefully turning right, the model sometimes imagines a new road.



265 Figure 7: **Examples of hallucination.** Top: object suddenly disappears. bottom: a non-existed road.

266 **Long-horizon visual prediction.** Our baseline model is only capable of generating short videos,
 267 which can be used to simulate short-term interaction events. Longer predictions [4, 55, 21] and faster
 268 generation [37, 35] are left for future research.

269 6 Conclusion

270 In this work, we present DrivingDojo, a large-scale video dataset aimed at advancing the study of
 271 driving world models. DrivingDojo offers a testbed for studying diverse real-world interactions. Our
 272 findings indicate that simulating interactions and rare dynamics observed in open-world environments
 273 remains an unsolved challenge, highlighting significant opportunities for future research.

274 **Societal impacts.** By providing a comprehensive dataset covering diverse driving scenarios and
 275 behaviors, researchers can develop and refine algorithms that increase the safety, reliability, and
 276 efficiency of autonomous vehicles. However, the development of driving world model requires large
 277 and diverse driving videos, introducing privacy issues.

278 **References**

- 279 [1] Mina Alibeigi, William Ljungbergh, Adam Tonderski, Georg Hess, Adam Lilja, Carl Lindström,
280 Daria Motorniuk, Junsheng Fu, Jenny Widahl, and Christoffer Petersson. Zenseact open dataset:
281 A large-scale and diverse multimodal dataset for autonomous driving. In *ICCV*, 2023. 3
- 282 [2] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Do-
283 minik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion:
284 Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
285 2, 7, 8, 25
- 286 [3] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja
287 Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent
288 diffusion models. In *CVPR*, 2023. 2
- 289 [4] Tim Brooks, Janne Hellsten, Miika Aittala, Ting-Chun Wang, Timo Aila, Jaakko Lehtinen,
290 Ming-Yu Liu, Alexei Efros, and Tero Karras. Generating long videos of dynamic scenes.
291 *NeurIPS*, 35, 2022. 9
- 292 [5] Jake Bruce, Michael Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes,
293 Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative
294 interactive environments. *arXiv preprint arXiv:2402.15391*, 2024. 4, 7
- 295 [6] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu,
296 Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal
297 dataset for autonomous driving. In *CVPR*, 2020. 2, 3, 5, 8, 25
- 298 [7] Holger Caesar, Juraj Kabzan, Kok Seang Tan, Whye Kit Fong, Eric Wolff, Alex Lang, Luke
299 Fletcher, Oscar Beijbom, and Sammy Omari. nuplan: A closed-loop ml-based planning
300 benchmark for autonomous vehicles. *arXiv preprint arXiv:2106.11810*, 2021. 3, 4, 5
- 301 [8] Ming-Fang Chang, John Lambert, Patsorn Sangkloy, Jagjeet Singh, Slawomir Bak, Andrew
302 Hartnett, De Wang, Peter Carr, Simon Lucey, Deva Ramanan, et al. Argoverse: 3d tracking and
303 forecasting with rich maps. In *CVPR*, 2019. 3
- 304 [9] Chang Chen, Jaesik Yoon, Yi-Fu Wu, and Sungjin Ahn. Transdreamer: Reinforcement learning
305 with transformer world models. In *Deep RL Workshop NeurIPS 2021*, 2021. 4
- 306 [10] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo
307 Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic
308 urban scene understanding. In *CVPR*, 2016. 3
- 309 [11] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla:
310 An open urban driving simulator. In *CoRL*, 2017. 3, 4
- 311 [12] Scott Ettinger, Shuyang Cheng, Benjamin Caine, Chenxi Liu, Hang Zhao, Sabeek Pradhan,
312 Yuning Chai, Ben Sapp, Charles R Qi, Yin Zhou, et al. Large scale interactive motion forecasting
313 for autonomous driving: The waymo open motion dataset. In *ICCV*, 2021. 3
- 314 [13] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the
315 kitti vision benchmark suite. In *CVPR*, 2012. 3
- 316 [14] Jakob Geyer, Yohannes Kassahun, Mentar Mahmudi, Xavier Ricou, Rupesh Durgesh, Andrew S
317 Chung, Lorenz Hauswald, Viet Hoang Pham, Maximilian Mühlegg, Sebastian Dorn, et al. A2d2:
318 Audi autonomous driving dataset. *arXiv preprint arXiv:2004.06320*, 2020. 3
- 319 [15] Yanchen Guan, Haicheng Liao, Zhenning Li, Jia Hu, Runze Yuan, Yunjian Li, Guohui Zhang,
320 and Chengzhong Xu. World models for autonomous driving: An initial survey. *IEEE Transac-*
321 *tions on Intelligent Vehicles*, 2024. 4

- 322 [16] David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018. 2,
323 4
- 324 [17] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control:
325 Learning behaviors by latent imagination. In *ICLR*, 2019. 4
- 326 [18] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and
327 James Davidson. Learning latent dynamics for planning from pixels. In *ICML*, 2019. 2
- 328 [19] Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with
329 discrete world models. *arXiv preprint arXiv:2010.02193*, 2020. 2, 4
- 330 [20] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains
331 through world models. *arXiv preprint arXiv:2301.04104*, 2023. 2, 4
- 332 [21] Roberto Henschel, Levon Khachatryan, Daniil Hayrapetyan, Hayk Poghosyan, Vahram Tade-
333 vosyan, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Streamingt2v: Consistent,
334 dynamic, and extendable long video generation from text. *arXiv preprint arXiv:2403.14773*,
335 2024. 9
- 336 [22] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.
337 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *NeurIPS*, 30,
338 2017. 7, 22
- 339 [23] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko,
340 Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High
341 definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 2
- 342 [24] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop
343 on Deep Generative Models and Downstream Applications*, 2021. 22
- 344 [25] John Houston, Guido Zuidhof, Luca Bergamini, Yawei Ye, Long Chen, Ashesh Jain, Sammy
345 Omari, Vladimir Iglovikov, and Peter Ondruska. One thousand and one hours: Self-driving
346 motion prediction dataset. In *CoRL*, 2021. 3
- 347 [26] Anthony Hu, Gianluca Corrado, Nicolas Griffiths, Zachary Murez, Corina Gurau, Hudson Yeo,
348 Alex Kendall, Roberto Cipolla, and Jamie Shotton. Model-based imitation learning for urban
349 driving. *NeurIPS*, 2022. 4
- 350 [27] Anthony Hu, Lloyd Russell, Hudson Yeo, Zak Murez, George Fedoseev, Alex Kendall, Jamie
351 Shotton, and Gianluca Corrado. Gaia-1: A generative world model for autonomous driving.
352 *arXiv preprint arXiv:2309.17080*, 2023. 2, 4
- 353 [28] Xinyu Huang, Peng Wang, Xinjing Cheng, Dingfu Zhou, Qichuan Geng, and Ruigang Yang. The
354 apolloscape open dataset for autonomous driving and its application. *TPAMI*, 42(10):2702–2719,
355 2019. 3
- 356 [29] Fan Jia, Weixin Mao, Yingfei Liu, Yucheng Zhao, Yuqing Wen, Chi Zhang, Xiangyu Zhang,
357 and Tiancai Wang. Adriver-i: A general world model for autonomous driving. *arXiv preprint
358 arXiv:2311.13549*, 2023. 2, 4
- 359 [30] Glenn Jocher, Ayush Chaurasia, and Jing Qiu. Ultralytics YOLO, January 2023. 6
- 360 [31] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of
361 diffusion-based generative models. *NeurIPS*, 35, 2022. 22
- 362 [32] Yann LeCun. A path towards autonomous machine intelligence. 2022. 2, 4
- 363 [33] Qifeng Li, Xiaosong Jia, Shaobo Wang, and Junchi Yan. Think2drive: Efficient reinforcement
364 learning by thinking in latent world model for quasi-realistic autonomous driving (in carla-v2).
365 *arXiv preprint arXiv:2402.16720*, 2024. 4

- 366 [34] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2018. 22
- 367 [35] Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models:
368 Synthesizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*,
369 2023. 9
- 370 [36] Jiageng Mao, Minzhe Niu, Chenhan Jiang, Hanxue Liang, Jingheng Chen, Xiaodan Liang,
371 Yamin Li, Chaoqiang Ye, Wei Zhang, Zhenguo Li, et al. One million scenes for autonomous
372 driving: Once dataset. *arXiv preprint arXiv:2106.11037*, 2021. 2, 3, 5, 25
- 373 [37] Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho,
374 and Tim Salimans. On distillation of guided diffusion models. In *CVPR*, 2023. 9
- 375 [38] Gerhard Neuhold, Tobias Ollmann, Samuel Rota Bulo, and Peter Kontschieder. The mapillary
376 vistas dataset for semantic understanding of street scenes. In *ICCV*, 2017. 3
- 377 [39] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks
378 for biomedical image segmentation. In *Medical image computing and computer-assisted
379 intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5–9,
380 2015, proceedings, part III 18*, pages 234–241, 2015. 7
- 381 [40] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In
382 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 7, 22
- 383 [41] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise
384 view selection for unstructured multi-view stereo. In *European Conference on Computer Vision
385 (ECCV)*, 2016. 7, 22
- 386 [42] Younggyo Seo, Danijar Hafner, Hao Liu, Fangchen Liu, Stephen James, Kimin Lee, and Pieter
387 Abbeel. Masked world models for visual control. In *CoRL*, 2023. 4
- 388 [43] Younggyo Seo, Kimin Lee, Stephen L James, and Pieter Abbeel. Reinforcement learning with
389 action-free pre-training from videos. In *ICML*, 2022. 4
- 390 [44] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul
391 Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for
392 autonomous driving: Waymo open dataset. In *CVPR*, 2020. 2, 3, 5
- 393 [45] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski,
394 and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges.
395 *arXiv preprint arXiv:1812.01717*, 2018. 7, 22
- 396 [46] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics.
397 *NeurIPS*, 2016. 2
- 398 [47] Xiaofeng Wang, Zheng Zhu, Guan Huang, Xinze Chen, and Jiwen Lu. Drivedreamer: Towards
399 real-world-driven world models for autonomous driving. *arXiv preprint arXiv:2309.09777*,
400 2023. 2, 4
- 401 [48] Yuqi Wang, Jiawei He, Lue Fan, Hongxin Li, Yuntao Chen, and Zhaoxiang Zhang. Driving
402 into the future: Multiview visual forecasting and planning with world model for autonomous
403 driving. *CVPR*, 2024. 2, 4
- 404 [49] Benjamin Wilson, William Qi, Tanmay Agarwal, John Lambert, Jagjeet Singh, Siddhesh
405 Khandelwal, Bowen Pan, Ratnesh Kumar, Andrew Hartnett, Jhony Kaesemel Pontes, et al.
406 Argoverse 2: Next generation datasets for self-driving perception and forecasting. *arXiv preprint
407 arXiv:2301.00493*, 2023. 3
- 408 [50] Philipp Wu, Alejandro Escontrela, Danijar Hafner, Pieter Abbeel, and Ken Goldberg. Day-
409 dreamer: World models for physical robot learning. In *CoRL*, 2023. 4

- 410 [51] Ziyi Wu, Nikita Dvornik, Klaus Greff, Thomas Kipf, and Animesh Garg. Slotformer: Unsuper-
411 vised visual dynamics simulation with object-centric models. In *ICLR*, 2023. 4
- 412 [52] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. Videogpt: Video generation
413 using vq-vae and transformers. *arXiv preprint arXiv:2104.10157*, 2021. 2
- 414 [53] Jiazhi Yang, Shenyuan Gao, Yihang Qiu, Li Chen, Tianyu Li, Bo Dai, Kashyap Chitta, Penghao
415 Wu, Jia Zeng, Ping Luo, et al. Generalized predictive model for autonomous driving. *arXiv*
416 *preprint arXiv:2403.09630*, 2024. 2, 3, 4, 5, 8, 25
- 417 [54] Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Dale Schuurmans, and
418 Pieter Abbeel. Learning interactive real-world simulators. *arXiv preprint arXiv:2310.06114*,
419 2023. 4
- 420 [55] Shengming Yin, Chenfei Wu, Huan Yang, Jianfeng Wang, Xiaodong Wang, Minheng Ni,
421 Zhengyuan Yang, Linjie Li, Shuguang Liu, Fan Yang, et al. Nuwa-xl: Diffusion over diffusion
422 for extremely long video generation. *arXiv preprint arXiv:2303.12346*, 2023. 9
- 423 [56] Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht
424 Madhavan, and Trevor Darrell. Bdd100k: A diverse driving dataset for heterogeneous multitask
425 learning. In *CVPR*, 2020. 3
- 426 [57] Wei Zhan, Liting Sun, Di Wang, Haojie Shi, Aubrey Clausse, Maximilian Naumann, Julius
427 Kummerle, Hendrik Konigshof, Christoph Stiller, Arnaud de La Fortelle, et al. Interaction
428 dataset: An international, adversarial and cooperative motion dataset in interactive driving
429 scenarios with semantic maps. *arXiv preprint arXiv:1910.03088*, 2019. 3
- 430 [58] Guosheng Zhao, Xiaofeng Wang, Zheng Zhu, Xinze Chen, Guan Huang, Xiaoyi Bao, and Xing-
431 gang Wang. Drivedreamer-2: Llm-enhanced world models for diverse driving video generation.
432 *arXiv preprint arXiv:2403.06845*, 2024. 2, 4

433 **Checklist**

- 434 1. For all authors...
 - 435 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
436 contributions and scope? [Yes] See section 1
 - 437 (b) Did you describe the limitations of your work? [Yes] See section 5.4
 - 438 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See
439 section 6
 - 440 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
441 them? [Yes]
- 442 2. If you are including theoretical results...
 - 443 (a) Did you state the full set of assumptions of all theoretical results? [N/A] No theoretical
444 results
 - 445 (b) Did you include complete proofs of all theoretical results? [N/A] No theoretical results
- 446 3. If you ran experiments (e.g. for benchmarks)...
 - 447 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
448 perimental results (either in the supplemental material or as a URL)? [Yes] In the
449 supplemental material
 - 450 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
451 were chosen)? [Yes] In the supplemental material
 - 452 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
453 ments multiple times)? [Yes] We repeat evaluation multiple times and report the mean
454 performance.
 - 455 (d) Did you include the total amount of compute and the type of resources used (e.g., type
456 of GPUs, internal cluster, or cloud provider)? [Yes]
- 457 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - 458 (a) If your work uses existing assets, did you cite the creators? [Yes]
 - 459 (b) Did you mention the license of the assets? [Yes] See supplemental material
 - 460 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
461 See <https://drivingdojo.github.io>
 - 462 (d) Did you discuss whether and how consent was obtained from people whose data you're
463 using/curating? [Yes] The public release of the data has been approved and authorized
464 by Meituan Inc.
 - 465 (e) Did you discuss whether the data you are using/curating contains personally identifi-
466 able information or offensive content? [Yes] See Personal Identification Information
467 Removal in Section 3.4
- 468 5. If you used crowdsourcing or conducted research with human subjects...
 - 469 (a) Did you include the full text of instructions given to participants and screenshots, if
470 applicable? [N/A]
 - 471 (b) Did you describe any potential participant risks, with links to Institutional Review
472 Board (IRB) approvals, if applicable? [N/A]
 - 473 (c) Did you include the estimated hourly wage paid to participants and the total amount
474 spent on participant compensation? [N/A]

475 Appendix

476	A Dataset	16
477	A.1 Overview	16
478	A.2 Curation	19
479	B Implementation Details	21
480	B.1 Experiment Setup	21
481	B.2 Training	22
482	B.3 Evaluation	22
483	C Additional Experiments	22
484	C.1 Action Instruction Following	22
485	D Visualizations	23
486	D.1 Diverse Actions	23
487	D.2 Dynamic Interaction	23
488	D.3 Open-world Knowledge	23
489	D.4 AIF Visualization	23
490	E License of Assets	25
491	F Datasheet	26
492	F.1 Motivation	26
493	F.2 Composition	26
494	F.3 Collection Process	27
495	F.4 Preprocessing/cleaning/labeling	28
496	F.5 Uses	28
497	F.6 Distribution	29
498	F.7 Maintenance	29

499 **A Dataset**

500 **A.1 Overview**

501 We will publish the DrivingDojo dataset, data format and annotation instructions, AIF benchmark,
502 and code for the baseline method on our project page: <https://drivingdojo.github.io>.

503 **Terms of use and License.** Our dataset is released under the **CC BY-NC 4.0** license, allowing
504 everyone to use it for non-commercial research purposes.

505 **Data maintenance.** The data is stored on Google Drive for global accessibility, and we will supply
506 various links (e.g., Baidu Cloud Drive and Hugging Face) for researchers' convenience. We will
507 maintain the data long-term and periodically verify its accessibility.

508 In the following, we showcase more video examples in our DrivingDojo dataset, the corresponding
509 videos are better illustrated on our project page.



Figure 8: Examples of rich ego-actions on the DrivingDojo dataset.

510 **Action completeness.** We include more dataset visualizations depicting various ego-actions in
 511 Figure 8. From top to bottom, the images show the ego vehicle performing left turns, right turns,
 512 going straight, lane-changing, and making emergency brakes during the driving.



Figure 9: Examples of multi-agent interplay on the DrivingDojo dataset.

513 **Multi-agent interplay.** Interaction plays a crucial role in driving scenarios. It usually means that
 514 the ego vehicle has engaged with other road users, leading to changes in the behavior of either the
 515 ego vehicle or the other road users. As shown in Figure 9, we present a series of interaction examples
 516 in our dataset. In the first scenario, the car suddenly encounters another vehicle crossing its path
 517 while moving forward, prompting an abrupt braking maneuver. The second scenario portrays the car
 518 encountering an electric scooter unexpectedly crossing its path. Illustrating the third scenario, the
 519 car comes across a vehicle in front opening its door, forcing an abrupt brake. In the fourth scenario,
 520 the challenge involves encountering a bicycle approaching from the opposite direction, while the
 521 fifth scenario involves navigating around a stroller. The sixth scenario showcases encountering
 522 road construction ahead, followed by encountering a street sweeper in the seventh scenario. The

523 eighth scenario presents a situation where a car suddenly makes a U-turn from the opposite direction,
 524 prompting an urgent braking response from our vehicle. Subsequent scenarios involve interactions
 525 with pedestrians. These diverse interaction scenarios provide a crucial foundation for studying the
 526 interaction of real-world simulators.



Figure 10: Examples of diverse open-world objects on the DrivingDojo dataset.

527 **Open-world knowledge.** In complex driving environments, we often encounter a wide variety of
 528 open-world situations. These scenarios can include sudden appearances of unexpected obstacles such
 529 as fallen trees, construction barriers, or abandoned vehicles. Typically belonging to the tail end of a
 530 long-tail distribution, these scenarios are rare yet crucial for ensuring safe driving. In Figure 10, we
 531 showcase a series of examples from the dataset, which fully demonstrate the richness of our dataset
 532 in capturing long-tail scenarios. From top to bottom, the examples illustrate encounters with a crane,
 533 a towing rope, construction barriers, a fallen roadblock, a vehicle transporting iron pipes, a vehicle
 534 transporting tree branches, a herd of sheep, an excavator, a bonfire, and power lines.

535 **A.2 Curation**

536 In this section, we provide the details of the curation procedure of each subset of DrivingDojo dataset
537 and the descriptions of curated actions and interactions. This section supplements the details for
538 Section 3.4 in the main paper.

539 **Action Completeness** Ego maneuvers for a car, particularly in the context of autonomous driving,
540 refer to the actions and decisions the vehicle makes to navigate its environment safely and efficiently.
541 Here is an exhaustive list of common ego maneuvers, and some examples in our datasets are shown
542 in Figure 1a in the main paper:

- 543 • **Acceleration:** Increasing speed to match traffic flow.
- 544 • **Deceleration:** Gradual slowing down for stop signs, traffic lights, or traffic congestion.
- 545 • **Lane Keeping:** Maintaining the current lane.
- 546 • **Lane Changing:** Changing lanes to overtake slower vehicles or merge into traffic.
- 547 • **Turning:** Left/right or U-turns at intersections or roundabouts.
- 548 • **Stop and Move on:** Stopping/proceeding at traffic signals or stop signs.
- 549 • **Emergency brake:** Abrupt and sudden braking maneuver to avoid a collision or mitigate the
550 impact of a potential hazard.

551 So, in the DrivingDojo-Action set, the videos follow different action commands, and the actions are
552 mainly from the planning and control (PNC) signals, such as left and right turns, straight crossings,
553 and lane changes. Each curated video clip begins with the PNC issuing a specific command and ends
554 when the command is completed.

555 **Multi-agent Interplay** The examples of multi-agent interplay are shown in Figure 1b in the main
556 paper. Then we describe the detailed cases of the interactions with dynamic agents.

- 557 • **Cutting in/off:** Another vehicle abruptly changes lanes and enters the path of the autonomous
558 vehicle. Ego vehicle changes lanes and enters the path of the other vehicles.
- 559 • **Meeting:** Ego vehicle encounters other vehicles traveling in the opposite direction.
- 560 • **Blocked:** Ego vehicle is stopped by other agents, such as vehicles, motorcycles, and pedestrians.
- 561 • **Overtaking and being overtaken:** Ego vehicle attempting to pass another vehicle and being
562 passed by another vehicle.

563 In the DrivingDojo-Interplay set, the core data curation strategy is to find the interaction with other
564 agents. The interaction is determined using PNC signals and manually defined rules. The main
565 interaction videos are from PNC dangerous interaction data. PNC conducts a deduction between the
566 ego vehicle and obstacles. When the ego vehicle cannot avoid collision by turning the steering wheel
567 or slowing down slightly, it is a PNC interaction case.

568 **Open-world Knowledge** Here, we select some representative and interesting examples from these
569 rare cases and show them in Figure 1c in the main paper. Based on the provided image and the given
570 descriptions, here are the detailed descriptions of each rare case in autonomous driving:

- 571 (a) A worker's helmet rolls on the sidewalk next to the road. (b) A soccer ball is seen flying across
572 the road. (c) A water bucket is depicted falling onto the road. (d) Parcel boxes have fallen onto the
573 road. (e) A dog is crossing the road. (f) A rope is floating over the road. (g) The traffic light turns red.
574 (h) A boom barrier blocks the vehicle from moving forward.

575 As mentioned above, we curated DrivingDojo-Open set in which the videos are more carefully
576 categorized and labeled with text descriptions. The sources are unusual weather, foreign objects on
577 the road surface, floating obstacles, falling objects, taking over cases, and interactions with traffic
578 lights and boom barriers. For curating the foreign objects/obstacles, we manually check and label
579 them by a large number of data annotators.

580 **Dataset Format** DrivingDojo dataset provides a file named ‘dataset_info.json’ that stores information corresponding to each video segment, including the information shown in Table 5. The ‘type’
 581 represents the major category, ‘tag’ represents the minor category, and ‘remark’ provides detailed
 582 descriptions of the reasons for hard braking and intervention.
 583

Table 5: The explanation of the information in dataset_info.json.

Information	Detailed explanation
meta_info	weather, location, time, frame number
description	type, tag, remark
videos	the image path for each frame
camera_info	the camera intrinsic parameters and extrinsic matrix for each frame
action_info	the coordinates of the next frame’s camera position in the current camera coordinate system

584 The following is an example directory structure for a dataset:

```

585 .
586 dataset_info.json
587 action_info
588   062959_s20-370_1712024694.0_1712024714.0
589     0023_next_frame_position_at_current_camera.txt
590     0025_next_frame_position_at_current_camera.txt
591     0027_next_frame_position_at_current_camera.txt
592     ...
593   145325_s20-190_1683790938.0_1683790958.0
594     0024_next_frame_position_at_current_camera.txt
595     0026_next_frame_position_at_current_camera.txt
596     0028_next_frame_position_at_current_camera.txt
597     ...
598   ...
599 camera_info
600   062959_s20-370_1712024694.0_1712024714.0
601     0023_camera_parameters.txt
602     0025_camera_parameters.txt
603     0027_camera_parameters.txt
604     ...
605   145325_s20-190_1683790938.0_1683790958.0
606     0024_camera_parameters.txt
607     0026_camera_parameters.txt
608     0028_camera_parameters.txt
609     ...
610   ...
611 videos
612   062959_s20-370_1712024694.0_1712024714.0
613     0023_CameraFpgaP0H120.jpg
614     0025_CameraFpgaP0H120.jpg
615     0027_CameraFpgaP0H120.jpg
616     ...
617   145325_s20-190_1683790938.0_1683790958.0
618     0024_CameraFpgaP0H120.jpg
619     0026_CameraFpgaP0H120.jpg
620     0028_CameraFpgaP0H120.jpg
621     ...
622   ...

```

623 **Camera info.** The ‘camera info’ refers to the extrinsic and intrinsic matrices of each frame of a
 624 fisheye camera. The world coordinate system is chosen as the East-North-Up (ENU) coordinate
 625 system. In the camera coordinate system, the x, y, and z axes respectively point to the right, down,
 626 and forward. We normalize the world coordinate system of the first frame to the origin, which means
 627 that the translation variables in the extrinsic matrices of each frame are subtracted by the translation
 628 variables of the first frame.

629 **Action info.** The ‘action info’ represents the coordinates of the next frame’s camera position in
 630 the current camera coordinate system. Let the transformation matrix from the camera to the world
 631 coordinate system be $\begin{pmatrix} R & T \\ 0^3 & 1 \end{pmatrix}$. The calculation method for the action info A_n of the n -th frame
 632 is shown in formula 3. The orientation of xyz axes in matrix A_n is consistent with the camera
 633 coordinate system, where the x, y, and z axes respectively point to the right, down, and forward.

$$\begin{pmatrix} A_n \\ 1 \end{pmatrix} = \begin{pmatrix} R_n & T_n \\ 0^3 & 1 \end{pmatrix}^{-1} \begin{pmatrix} T_{n+1} \\ 1 \end{pmatrix} \quad (3)$$

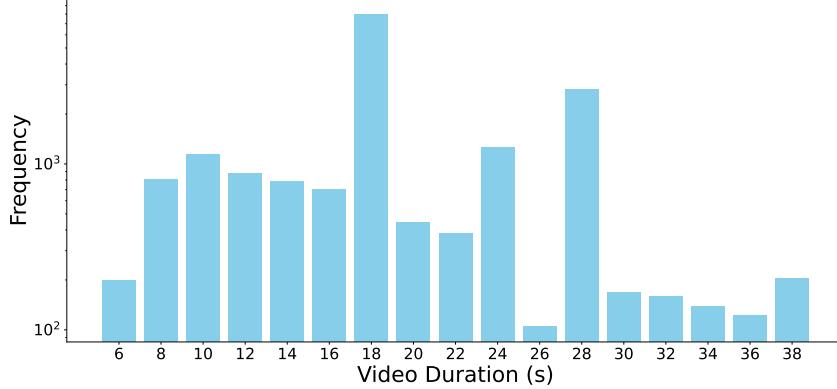


Figure 11: **Distribution of DrivingDojo video duration.**

634 **Video info.** The video is stored as a sequence of individual image frames. The distribution of video
 635 duration is shown in Figure 11, with the majority of videos lasting around 20 seconds.

636 **B Implementation Details**

637 **B.1 Experiment Setup**

638 During the experiment, we employed two settings for training the model. In the first setting, the focus
 639 is on visual prediction: the model predicts subsequent video content based solely on the initial frame
 640 image. In the second setting, we employ action-controlled video generation. Here, alongside the
 641 initial frame image, action information for the subsequent frames is provided to the model, enabling
 642 it to predict the ensuing video content.

643 **Visual prediction.** In this setup, we trained a high-resolution version of the model, 1024×576
 644 resolution for 11 frames, aimed at better capturing the generation of long-tail objects. Additionally,
 645 we developed a low-resolution version of the model, 576×320 resolution for 30 frames, to simulate
 646 various vehicle behaviors and interaction events. We fine-tune all parameters of the U-Net model.

647 **Action instruction following.** In this setup, we trained model using 576×320 resolution for 30
 648 frames. We fine-tune all U-Net parameters together with a new action encoder.

649 **B.2 Training**

650 We initialize the model using the SVD-XT checkpoint. Following SVD, our model is trained with
651 the EDM framework [31]. During training, we set the fps to 5 and the motion_bucket_id to 127. We
652 utilize the AdamW optimizer [34] with a learning rate of 1×10^{-5} . The training process is conducted
653 on 16 NVIDIA A100 (80G) GPUs with 32 batch size for 50K iterations. To allow classifier-free
654 guidance [24], we drop out action feature with a ratio of 20%.

655 **B.3 Evaluation**

656 During inference, we generate videos using the DDIM sampler for 25 steps.

657 **Visual Quality.** To evaluate the quality of the generated video, we utilize FID (Frechet Inception
658 Distance) [22] and FVD (Frechet Video Distance) [45] as the main metrics. For FID calculation
659 on videos, we randomly select 5,000 frames for evaluation. Additionally, for FVD calculation, we
660 generate 256 videos for evaluation. The results are the average of 10 calculations. We use the official
661 UCF FVD evaluation code².

662 **Action instruction following (AIF).** For each generated video with action instructions, we estimate
663 the camera poses for each frame in the video, align the scale of the estimated trajectory with the
664 instruction trajectory, and compare the vehicle motion in each frame with the respective action
665 instructions. We estimate the ego trajectories in generated videos with the offline visual structure-
666 from-motion (SfM) implementation COLMAP [40, 41]. For videos generated based on initial images
667 from DrivingDojo, we fix the camera intrinsic parameters as the ground truth values for videos from
668 DrivingDojo. For videos generated from initial images with unknown camera intrinsics (e.g. images
669 from OpenDV-2K), we estimate the camera intrinsics together with the camera extrinsics of images.
670 We perform feature point extraction, feature point matching, and sparse scene reconstruction with
671 the official implementation of COLMAP³ to estimate the poses of cameras. In our experiments,
672 we generate videos in 30 frames and align the scale of estimated trajectories with the instruction
673 trajectories based on the motions in the first $N = 10$ frames. We report the mean value of the absolute
674 error between estimated motions and instruction motions in all video frames.

675 **C Additional Experiments**

Table 6: **Action instruction following under zero-shot evaluation.** * denotes the model is applied
zero-shot to this dataset without fine-tuning.

Training set	Test set	FID	FVD	$E_x^{\text{AIF}}(\downarrow)$	$E_y^{\text{AIF}}(\downarrow)$
DrivingDojo	OpenDV-2K*	24.27	442.67	0.175m	0.121m
ONCE	OpenDV-2K*	28.37	473.59	0.295m	0.308m
nuScenes	OpenDV-2K*	37.90	794.36	0.488m	0.465m

676 **C.1 Action Instruction Following**

677 As shown in Table 6, we compared the performance of models trained on different datasets and their
678 zero-shot generalization performance on new datasets. The results indicate that models trained on our
679 dataset exhibit higher generation quality and significantly improved action-following ability. Espe-
680 cially, we noticed that richer driving actions in the autonomous driving datasets lead to significantly
681 better AIF performance of models trained on them. According to Figure 3a, videos in DrivingDojo
682 averagely contain far richer driving actions compared to ONCE or nuScenes. This leads to the far

²<https://github.com/SongweiGe/TATS/>

³<https://github.com/colmap/colmap>

683 better AIF performance of model trained on DrivingDojo compared to those trained on ONCE or
 684 nuScenes. we observed that the model trained on the ONCE dataset will always generate videos in
 685 which the vehicle moves in a straight line, even with action instructions to turn left/right or change
 686 lanes. This leads to its especially poor AIF performance in the lateral direction (E_y^{AIF}). We speculate
 687 that this is because the driving action of making turns or changing lanes is very rare in the ONCE
 688 dataset, as shown in Figure 3a, which results in the lack of ability of the model trained on the ONCE
 689 dataset to follow the lateral motion instructions. Moreover, the even more lacking driving actions in
 690 the nuScenes dataset lead to a worse AIF performance of the world model.

691 D Visualizations

692 In this section, we show the model generation demos trained on the DrivingDojo dataset. As shown
 693 in Figure 12, our model can generate high-resolution, complex driving scenarios.



Figure 12: **Examples of high-resolution and complex scenarios generation.** For illustration purposes, we represent each video example with a single frame.

694 D.1 Diverse Actions

695 As shown in Figure 13, we demonstrate how actions control the generation of different futures, such
 696 as moving forward, backward, and stopping.

697 D.2 Dynamic Interaction

698 As shown in Figure 14, we observe that choosing different actions can influence the behavior of other
 699 vehicles, resulting in different responses from the world model. For instance, in the first example,
 700 if we choose to proceed slowly, the vehicle on the left decides to stop and yield. Conversely, if our
 701 vehicle stops, the left vehicle perceives an obstruction and slightly reverses to make way. In the
 702 second example, when we choose to brake, the right vehicle quickly cuts in front of us, while if we
 703 choose to proceed straight, the right vehicle waits in place.

704 D.3 Open-world Knowledge

705 As illustrated in Figure 15, we demonstrate the model’s ability to simulate various open-world
 706 objects, such as encountering construction zones, rare objects like ladders or balloons on the road,
 707 and simulating a puddle of water on the ground.

708 D.4 AIF Visualization

709 We showcase examples of estimated trajectories from generated videos in Figure 16. In each frame,
 710 the red dot represents the current estimated camera pose and the black dots represent the camera
 711 poses in past frames.

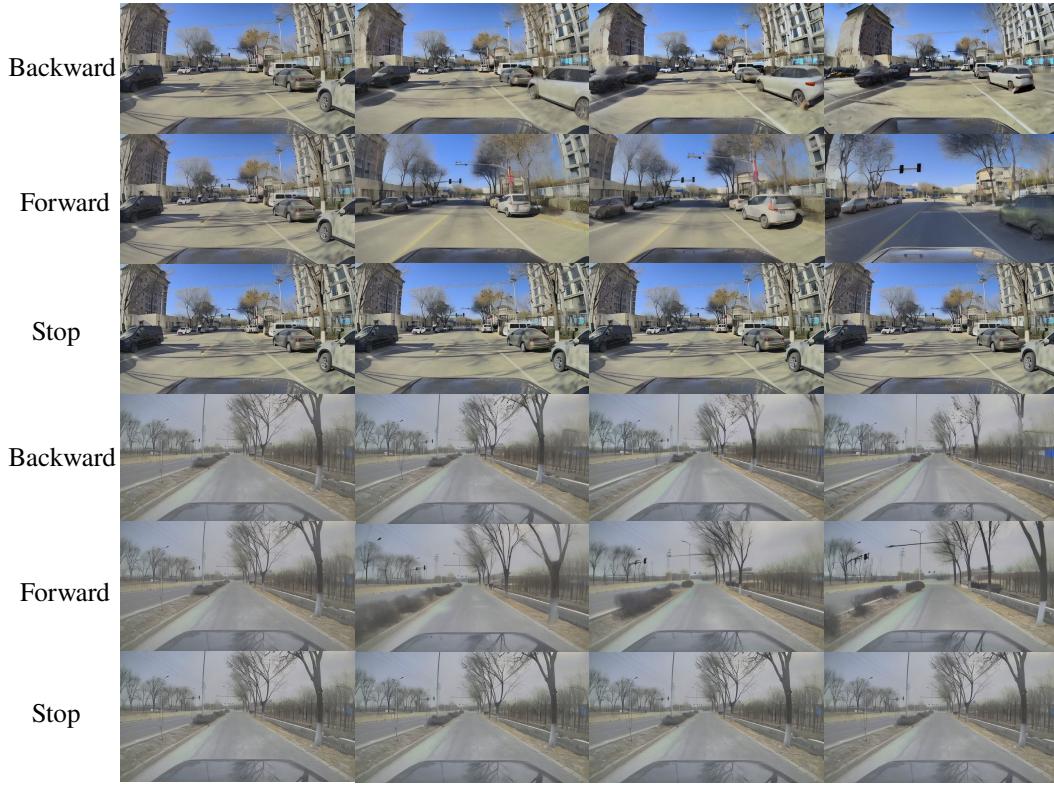


Figure 13: Examples of diverse action-based video generation.



Figure 14: Simulation of interaction with other agents.



Figure 15: Simulation of various open-world objects on the road.



Figure 16: Examples of ego trajectories estimated based on the generated videos.

712 E License of Assets

713 We report licenses of all artifacts used in this work in this section.

714 **Model** We use the pre-trained stable video diffusion [2] checkpoints from the huggingface platform.
 715 These checkpoints are released under the stable video diffusion non-commercial community license
 716 agreement⁴ for research purpose.

717 **Our Dataset** Our dataset is collected and curated by the autonomous driving team of Meituan Inc.
 718 The road test and data collection procedures conform to privacy and security requirements of local
 719 authorities. The authors have obtained the permission for publicly releasing this dataset from both
 720 the management team and the company legal team. All personal identifiable information has been
 721 removed by both algorithm and subsequent manual inspection. We release the dataset under the CC
 722 BY-NC 4.0 license.

723 **Other Datasets** We use other public datasets in this work including nuScenes [6], ONCE [36] and
 724 OpenDV-2k [53]. The nuScenes [6] dataset is released under the CC BY-NC-SA 4.0 license with
 725 Dataset Terms⁵. The ONCE dataset is also released under the CC BY-NC-SA 4.0 license with
 726 Dataset Terms⁶. The OpenDV-2K dataset is constructed from publicly licensed datasets and youtube
 727 videos that the authors claimed to support academic usage licenses.

⁴<https://huggingface.co/stabilityai/stable-video-diffusion-img2vid-xt/blob/main/LICENSE>

⁵<https://www.nuscenes.org/terms-of-use>

⁶https://once-for-auto-driving.github.io/terms_of_use.html

728 **F Datasheet**

729 **F.1 Motivation**

- 730 • **For what purpose was the dataset created?** Was there a specific task in mind? Was there
731 a specific gap that needed to be filled? Please provide a description.

732 We introduce DrivingDojo, the first dataset tailor-made for training interactive world models
733 with complex driving dynamics. Our dataset features video clips with a complete set of
734 driving maneuvers, diverse multi-agent interplay, and rich open-world driving knowledge,
735 laying a stepping stone for future world model development.

- 736 • **Who created the dataset (e.g., which team, research group) and on behalf of which
737 entity (e.g., company, institution, organization)?**

738 Institute of Automation, Chinese Academy of Sciences, University of Chinese Academy of
739 Sciences, Meituan Inc., and Centre for Artificial Intelligence and Robotics, HKISI_CAS.

- 740 • **Who funded the creation of the dataset?** If there is an associated grant, please provide the
741 name of the grantor and the grant name and number.

742 This work was supported in part by the National Key R&D Program of China (No.
743 2022ZD0116500), the National Natural Science Foundation of China (No. U21B2042,
744 No. 62320106010), and in part by the 2035 Innovation Program of CAS, and the InnoHK
745 program.

- 746 • **Any other comments?**

747 No.

748 **F.2 Composition**

- 749 • **What do the instances that comprise the dataset represent (e.g., documents, photos,
750 people, countries)?** Are there multiple types of instances (e.g., movies, users, and ratings;
751 people and interactions between them; nodes and edges)? Please provide a description.

752 The instances of our DrivingDojo dataset are videos with ego actions and DrivingDojo-Open
753 subset is also with text descriptions for each scene.

- 754 • **How many instances are there in total (of each type, if appropriate)?**

755 There are 18.2k videos for the whole DrivingDojo dataset, in which the DrivingDojo-Action
756 subset has 7.9k videos, DrivingDojo-Interplay subset has 6.4k videos, and DrivingDojo-
757 Open has 3.9k videos.

- 758 • **Does the dataset contain all possible instances or is it a sample (not necessarily random)
759 of instances from a larger set?** If the dataset is a sample, then what is the larger set? Is the
760 sample representative of the larger set (e.g., geographic coverage)? If so, please describe
761 how this representativeness was validated/verified. If it is not representative of the larger set,
762 please describe why not (e.g., to cover a more diverse range of instances, because instances
763 were withheld or unavailable).

764 The DrivingDojo dataset is sampled from a data pool of around 7500 hours. About represen-
765 tativeness, please refer to the Data Curation section (Sec. 3.4 and Sec. A.2).

- 766 • **What data does each instance consist of?** “Raw” data (e.g., unprocessed text or images)
767 or features? In either case, please provide a description.

768 DrivingDojo-Action and DrivingDojo-Interplay subsets consist of videos and ego actions,
769 and DrivingDojo-Open subset consists of videos, ego actions, and text descriptions.

- 770 • **Is there a label or target associated with each instance?** If so, please provide a description.

771 Yes. There is a text description label for each instance in DrivingDojo-Open subset, which
772 describes the open-world knowledge in the scene.

- 773 • **Is any information missing from individual instances?** If so, please provide a description,
 774 explaining why this information is missing (e.g., because it was unavailable). This does not
 775 include intentionally removed information, but might include, e.g., redacted text.
 776 No.
- 777 • **Are relationships between individual instances made explicit (e.g., users' movie ratings,
 778 social network links)?** If so, please describe how these relationships are made explicit.
 779 No.
- 780 • **Are there recommended data splits (e.g., training, development/validation, testing)?** If
 781 so, please provide a description of these splits, explaining the rationale behind them.
 782 No. There is no need for the validation/testing split. We care about zero-shot generation.
- 783 • **Are there any errors, sources of noise, or redundancies in the dataset?** If so, please
 784 provide a description.
 785 Yes. The sources of noise may be inaccurate poses, camera noises, and human-sourced text
 786 noises.
- 787 • **Is the dataset self-contained, or does it link to or otherwise rely on external resources
 788 (e.g., websites, tweets, other datasets)?** If it links to or relies on external resources, a) are
 789 there guarantees that they will exist, and remain constant, over time; b) are there official
 790 archival versions of the complete dataset (i.e., including the external resources as they
 791 existed at the time the dataset was created); c) are there any restrictions (e.g., licenses, fees)
 792 associated with any of the external resources that might apply to a dataset consumer? Please
 793 provide descriptions of all external resources and any restrictions associated with them, as
 794 well as links or other access points, as appropriate.
 795 Yes. the DrivingDojo dataset is self-contained.
- 796 • **Does the dataset contain data that might be considered confidential (e.g., data that is
 797 protected by legal privilege or by doctor–patient confidentiality, data that includes the
 798 content of individuals' non-public communications)?** If so, please provide a description.
 799 No.
- 800 • **Does the dataset contain data that, if viewed directly, might be offensive, insulting,
 801 threatening, or might otherwise cause anxiety?** If so, please describe why.
 802 No.

803 **E.3 Collection Process**

- 804 • **How was the data associated with each instance acquired?** Was the data directly ob-
 805 servable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or
 806 indirectly inferred/derived from other data (e.g., part-of-speech tags, model-based guesses
 807 for age or language)? If the data was reported by subjects or indirectly inferred/derived from
 808 other data, was the data validated/verified? If so, please describe how.
 809 DrivingDojo dataset is collected using the platform of Meituan's autonomous delivery
 810 vehicles.
- 811 • **What mechanisms or procedures were used to collect the data (e.g., hardware appa-
 812 ratuses or sensors, manual human curation, software programs, software APIs)?** How
 813 were these mechanisms or procedures validated?
 814 DrivingDojo dataset is collected using the platform of Meituan's autonomous delivery
 815 vehicles with fish-eye RGB cameras. The cameras are calibrated. The text labels are
 816 manually validated.
- 817 • **If the dataset is a sample from a larger set, what was the sampling strategy (e.g.,
 818 deterministic, probabilistic with specific sampling probabilities)?**
 819 Please refer to the Data Curation section (Sec. 3.4 and Sec. A.2).

- 820 • Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and how were they compensated (e.g., how much were crowdworkers paid)?

821 The data collectors are employed by Meituan Inc. and are paid by Meituan Inc.

- 823 • Over what timeframe was the data collected? Does this timeframe match the creation
824 timeframe of the data associated with the instances (e.g., recent crawl of old news articles)?
825 If not, please describe the timeframe in which the data associated with the instances was
826 created.

827 The data are collected from May 2022 to May 2024. This timeframe matches the creation
828 timeframe of the data associated with the instances.

- 829 • Were any ethical review processes conducted (e.g., by an institutional review board)?
830 If so, please provide a description of these review processes, including the outcomes, as well
831 as a link or other access point to any supporting documentation.

832 Yes. The ethical review is conducted before the release by Meituan Inc.

833 E.4 Preprocessing/cleaning/labeling

- 834 • Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucket-
835 ing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances,
836 processing of missing values)? If so, please provide a description. If not, you may skip the
837 remaining questions in this section.

838 No.

- 839 • Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to
840 support unanticipated future uses)? If so, please provide a link or other access point to
841 the “raw” data.

842 N/A.

- 843 • Is the software that was used to preprocess/clean/label the data available? If so, please
844 provide a link or other access point.

845 N/A.

- 846 • Any other comments?

847 No.

848 E.5 Uses

- 849 • Has the dataset been used for any tasks already? If so, please provide a description.

850 The DrivingDojo dataset has been used for driving world models. The experiments are in
851 Sec. 5 in the main paper and Sec. C in the appendix.

- 852 • Is there a repository that links to any or all papers or systems that use the dataset? If
853 so, please provide a link or other access point.

854 Yes. Please refer to the webset:<https://drivingdojo.github.io>.

- 855 • What (other) tasks could the dataset be used for?

856 The DrivingDojo dataset could be used for training end-to-end autonomous driving models.

- 857 • Is there anything about the composition of the dataset or the way it was collected
858 and preprocessed/cleaned/labeled that might impact future uses? For example, is there
859 anything that a dataset consumer might need to know to avoid uses that could result in unfair
860 treatment of individuals or groups (e.g., stereotyping, quality of service issues) or other risks
861 or harms (e.g., legal risks, financial harms)? If so, please provide a description. Is there
862 anything a dataset consumer could do to mitigate these risks or harms?

863 No.

- 864 • Are there tasks for which the dataset should not be used? If so, please provide a
865 description.

866 Due to the known biases of the dataset, under no circumstance should any models be put
867 into production using the dataset as is. It is neither safe nor responsible. As it stands, the
868 dataset should be solely used for research purposes in its uncurated state.

869 • **Any other comments?**

870 No.

871 **F.6 Distribution**

872 • **Will the dataset be distributed to third parties outside of the entity (e.g., company,
873 institution, organization) on behalf of which the dataset was created?** If so, please
874 provide a description.

875 Yes, the dataset will be open-source.

876 • **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?** Does
877 the dataset have a digital object identifier (DOI)?

878 On our website: <https://drivingdojo.github.io>.

879 • **When will the dataset be distributed?**

880 We have released some demos on the project page. The whole DrivingDojo dataset will be
881 public in the camera-ready version.

882 • **Will the dataset be distributed under a copyright or other intellectual property (IP)
883 license, and/or under applicable terms of use (ToU)?** If so, please describe this license
884 and/or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant
885 licensing terms or ToU, as well as any fees associated with these restrictions.

886 DrivingDojo dataset will be distributed under the CC BY-NC 4.0 license.

887 • **Have any third parties imposed IP-based or other restrictions on the data associated
888 with the instances?** If so, please describe these restrictions, and provide a link or other
889 access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees
890 associated with these restrictions.

891 No.

892 • **Do any export controls or other regulatory restrictions apply to the dataset or to
893 individual instances?** If so, please describe these restrictions, and provide a link or other
894 access point to, or otherwise reproduce, any supporting documentation.

895 No.

896 • **Any other comments?**

897 No.

898 **F.7 Maintenance**

899 • **Who will be supporting/hosting/maintaining the dataset?**

900 Institute of Automation, Chinese Academy of Sciences and Meituan Inc. will maintain
901 DrivingDojo dataset.

902 • **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

903 The main maintainer Yuqi Wang's e-mail: wangyuqi2020@ia.ac.cn.

904 • **Is there an erratum?** If so, please provide a link or other access point.

905 There is no erratum for our initial release. Errata will be documented as future releases on
906 the dataset website.

907 • **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete in-
908 stances)?** If so, please describe how often, by whom, and how updates will be communicated
909 to dataset consumers (e.g., mailing list, GitHub)?

910 Yes. We will update the DrivingDojo dataset. Especially, we will adapt to end-to-end
911 autonomous driving tasks in the future. The update will be released on the website and
912 GitHub.

- 913 • **If the dataset relates to people, are there applicable limits on the retention of the data**
914 **associated with the instances (e.g., were the individuals in question told that their data**
915 **would be retained for a fixed period of time and then deleted)?** If so, please describe
916 these limits and explain how they will be enforced.

917 N/A.

- 918 • **Will older versions of the dataset continue to be supported/hosted/maintained?** If so,
919 please describe how. If not, please describe how its obsolescence will be communicated to
920 dataset consumers.

921 Yes. We will maintain the older versions of the dataset on the website and GitHub.

- 922 • **If others want to extend/augment/build on/contribute to the dataset, is there a mech-**
923 **anism for them to do so?** If so, please provide a description. Will these contributions
924 be validated/verified? If so, please describe how. If not, why not? Is there a process for
925 communicating/distributing these contributions to dataset consumers? If so, please provide
926 a description.

927 Yes. The dataset is open source under the CC BY-NC 4.0 license. So it is open to other
928 contributors.

- 929 • **Any other comments?**

930 No.