# Generalized Predictive Model for Autonomous Driving

Jiazhi Yang[1*]    Shenyuan Gao[2,1*]    Yihang Qiu[1*]    Li Chen[3,1†]    Tianyu Li[1]    Bo Dai[1]

Kashyap Chitta[4,5]    Penghao Wu[1]    Jia Zeng[1]    Ping Luo[3]    Jun Zhang[2♮]

Andreas Geiger[4,5♮]    Yu Qiao[1♮]    Hongyang Li[1†]

[1] OpenDriveLab and Shanghai AI Lab    [2] Hong Kong University of Science and Technology
[3] University of Hong Kong    [4] University of Tübingen    [5] Tübingen AI Center

Figure 1. **Overview of the GenAD paradigm**. We aim to establish a generalized video prediction paradigm for autonomous driving by presenting the largest multimodal driving video dataset to date, **OpenDV-2K**, and a generative model that predicts the future given past visual and textual input, **GenAD**. The strong generalization and controllability of GenAD is validated spanning a diverse spectrum of tasks, including zero-shot domain transfer, language-conditioned prediction, action-conditioned prediction, and motion planning.

## Abstract

*In this paper, we introduce the first large-scale video prediction model in the autonomous driving discipline. To eliminate the restriction of high-cost data collection and empower the generalization ability of our model, we acquire massive data from the web and pair it with diverse and high-quality text descriptions. The resultant dataset accumulates over 2000 hours of driving videos, spanning areas all over the world with diverse weather conditions and traffic scenarios. Inheriting the merits from recent latent diffusion models, our model, dubbed GenAD, handles the challenging dynamics in driving scenes with novel temporal reasoning blocks. We showcase that it can generalize to various unseen driving datasets in a zero-shot manner, surpassing general or driving-specific video prediction counterparts. Furthermore, GenAD can be adapted into an action-conditioned prediction model or a motion planner, holding great potential for real-world driving applications.*

## 1. Introduction

Autonomous driving agents, as a promising application of high-level artificial intelligence, perceive the surrounding environment, build internal world model representations,

---

*Equal contribution, ordered by coin toss. ♮Equal co-advising.
†Project lead. Primary contact: yangjiazhi@opendrivelab.com

make decisions, and take actions in response [18, 113]. However, despite dedicated efforts in academia and industry for decades, their deployment is still restricted to certain areas or scenarios, and they cannot be applied over the world seamlessly. One critical reason is the limited generalization ability of learned models in structured autonomous driving systems. Typically, perception models face challenges of generalizing to diverse environments with changes in geographical locations, sensor configurations, weather conditions, open-set objects, *etc.*; prediction and planning models fail to generalize to nondeterministic futures with rare scenarios and different driving intentions [3, 35, 135].

Motivated by how humans learn to perceive and cognize the world [57, 58, 111], we advocate employing driving videos as the universal interface that generalizes to diverse environments with dynamic futures. Based on this, a driving video predictive model is preferred to fully capture the world knowledge about driving scenarios (Fig. 1). By predicting the future, the video predictor essentially learns two vital aspects of autonomous driving: how the world operates, and how to maneuver safely in the wild.

Recently, the community has begun to adopt video as the interface to represent observation behavior and action for various robot tasks [22]. For domains such as classical video prediction and robotics, the video backgrounds are mostly static, the movement of robots is slow, and the resolution of videos is low. In contrast, for the driving scenarios, it struggles with outdoor environments being highly dynamic, agents encompassing much larger motions, and the sensory resolution covering a large range of view. These distinctions lead to substantial challenges for autonomous driving applications. Fortunately, there are some preliminary attempts on developing a video predictive model in the driving domain [8, 34, 44, 48, 52, 71, 80, 104, 106]. Despite promising progress in terms of prediction quality, these attempts have not achieved desirable capability of generalization as in classical robot tasks (*e.g.*, manipulation), being confined to either limited scenarios such as highways with low traffic density [8] and small-scale datasets [34, 48, 71, 104, 106], or restricted conditions that raises difficulties to generate diverse environments [80]. How to unveil the potential of video prediction models for driving remains seldom explored.

Motivated by the discussions above, we target at building a video predictive model for autonomous driving, capable of generalizing to new conditions and environments. To this end, we have to answer the following questions: *(1) What data can be obtained in a feasible and scalable manner? (2) How can we formulate a predictive model to capture the complex evolution of dynamic scenarios? (3) How can we apply the (foundation) model for downstream tasks?*

**Scaled Data.** To achieve powerful generalization ability, a substantial and diverse corpus of data is necessary. Inspired by the success of learning from Internet-scale data in foundation models [2, 55, 82], we construct our driving dataset from both the web and publicly licensed datasets. Compared to existing options, which are limited in scale and diversity due to their regulated collection processes, online data owns great diversity in several aspects: geographic locations, terrains, weather conditions, safety-critical scenarios, sensor settings, traffic elements, *etc*. To guarantee the data is of high-quality and desirable for large-scale training, we exhaustively collect driving recordings on YouTube and remove unintended corruption frames via rigorous human verification. Furthermore, videos are paired with diverse text-level conditions, including descriptions generated and refined with the aid of existing foundation models [62, 77], and high-level instructions inferred by a video classifier. Through these steps, we construct **OpenDV-2K**, the *largest* public driving dataset to date, containing more than 2000 hours of driving videos and being 374 times larger than the widely used nuScenes counterpart. Our dataset is publicly available at https://github.com/OpenDriveLab/DriveAGI.

**Generalized Predictive Model.** Learning a generalized driving video predictor bears several key challenges: generation quality, training efficiency, causal reasoning, and drastic view shift. We address these aspects by presenting a novel temporal generative model with two-stage learning. To capture the environment details, enhance generation quality, and maintain training efficiency simultaneously, we build upon the recent success of *latent diffusion models* (LDMs) [79, 87]. In the first stage, we transfer the generation distribution of LDM from its pre-trained general vision domain to the driving domain by fine-tuning it on OpenDV-2K images. In the second stage, we interleave the proposed temporal reasoning blocks into the original model and learn to predict the future given past frames and conditions. Contrary to conventional temporal modules [8, 42] that suffer from causal confusion and large motion, our solution consists of causal temporal attention and decoupled spatial attention to efficiently model the drastic spatiotemporal shift in highly dynamic driving scenes. After sufficient training, our **Gen**erative model for **A**utonomous **D**riving (**GenAD**)[1] can generalize to various scenarios in a zero-shot fashion.

**Extensions for Simulation and Planning.** After large-scale pre-training of video prediction, GenAD essentially understands how the world evolves and how to drive. We show how to adapt its learned knowledge for real-world driving problems, *i.e.*, simulation and planning. For simulation, we fine-tune the pre-trained model with future ego trajectories as additional conditions, to associate future imaginations with different ego actions. We also empower GenAD to perform planning on challenging benchmarks by

---

[1]Note that GenAD is abbreviated from both **Gen**erative models and **Gen**eralized capabilities.

| | Dataset | Duration (hours) | Front-view Frames | Geographic Diversity Countries | Cities | Sensor Setup |
|---|---|---|---|---|---|---|
| ✗ | KITTI [30] | 1.4 | 15k | 1 | 1 | fixed |
| ✗ | Cityscapes [21] | 0.5 | 25k | 3 | 50 | fixed |
| ✗ | Waymo Open⋆ [98] | 11 | 390k | 1 | 3 | fixed |
| ✗ | Argoverse 2⋆ [110] | 4.2 | 300k | 1 | 6 | fixed |
| ✓ | nuScenes [12] | 5.5 | 241k | 2 | 2 | fixed |
| ✓ | nuPlan⋆ [13] | 120 | 4.0M | 2 | 4 | fixed |
| ✓ | Talk2Car [24] | 4.7 | - | 2 | 2 | fixed |
| ✓ | ONCE [72] | 144 | 7M | 1 | - | fixed |
| ✓ | Honda-HAD [51] | 32 | 1.2M | 1 | - | fixed |
| ✓ | Honda-HDD-Action [84] | 104 | 1.1M | 1 | - | fixed |
| ✓ | Honda-HDD-Cause [84] | 32 | - | 1 | - | fixed |
| ✓ | OpenDV-YouTube (Ours) | 1747 | 60.2M | $\geq 40^\dagger$ | $\geq 244^\dagger$ | uncalibrated |
| - | **OpenDV-2K (Ours)** | **2059** | **65.1M** | $\mathbf{\geq 40^\dagger}$ | $\mathbf{\geq 244^\dagger}$ | **uncalibrated** |

Table 1. **OpenDV-2K comparison at a glance to existing counterparts in terms of scale and diversity**. Note that datasets with ✓ are included in OpenDV-2K (last row). ⋆Perception subset in Waymo Open, Argoverse 2, and nuPlan. †Estimated by GPT [78] from video titles.



Figure 2. **Geographic distribution of OpenDV-2K**. Our dataset covers ample driving scenarios around the world.

using a lightweight planner to translate latent features into the future trajectory of the ego vehicle. On account of its pre-trained ability to predict accurate future frames, our algorithm exhibits promising results in both simulation consistency and planning reliability.

## 2. OpenDV-2K Dataset

We introduce OpenDV-2K, a large-scale multimodal dataset for autonomous driving, to support the training of a generalized video prediction model. The main component is a vast corpus of high-quality YouTube driving videos, which are collected from all over the world, and are gathered into our dataset after a careful curation process. We automatically create language annotations for these videos using vision-language models. To further improve its diversity in sensor configurations and language expressions, we merge 7 publicly licensed datasets into our OpenDV-2K, as shown in Tab. 1. As a result, OpenDV-2K occupies a total of 2059 hours of videos paired with texts, including 1747 hours from YouTube and 312 hours from public datasets. We use OpenDV-YouTube and OpenDV-2K to specify the YouTube split and the overall dataset, respectively.

### 2.1. Diversity over Prior Datasets

A brief comparison with other public datasets is provided in Tab. 1. Beyond its significant scale, the proposed OpenDV-2K represents *diversity* across various aspects as follows.

**Globe-wise Geographic Distribution.** Due to the global nature of online videos, OpenDV-2K covers more than 40 countries and 244 cities worldwide. This is a tremendous improvement over previous public datasets, which are typically gathered in a small number of restricted areas. We plot the specific distribution of OpenDV-YouTube in Fig. 2.

**Open-world Driving Scenarios.** Our dataset provides a huge amount of realistic driving experience in the open world, covering rare environments like forests, extreme weather conditions like heavy snow, and appropriate driving behaviors in response to interactive traffic situations. These data are crucial for diversity and generalization yet are seldom collected in existing public datasets.

**Unrestricted Sensor Configurations.** Current driving datasets are confined to specific sensor configurations, including intrinsic and extrinsic camera parameters, image, sensor type, optics, *etc*., which poses great challenges for deploying the learned models with different sensors [65]. In contrast, YouTube driving videos are recorded in various types of vehicles with flexible camera setups, which aids in the robustness of the trained model when deployed using a novel camera setting.

### 2.2. Towards High-quality Multimodal Dataset

**Driving Video Collection and Curation.** Finding clean driving videos from the vast pool of the web is a tedious and costly task. To simplify the process, we start by selecting certain video uploaders, *i.e*., YouTubers. Judging from the average length and overall quality, we collect 43 YouTubers with 2139 high-quality front-view driving videos. To make sure there is no overlap between training and validation sets, we take all videos from 3 YouTubers for validation, with the remaining videos as the training set. To rule out non-driving frames like video introductions and subscription reminders, we discard a certain length of segments at the beginning and end of each video. Each frame is then described with language contexts using a VLM model, BLIP-2 [62]. We further remove the black frames and transition frames, which are not ideal for training, by manually checking if there are certain keywords in these contexts. An illustration of the dataset construction pipeline is in Fig. 3, and we introduce how to generate the contexts below.

Figure 3. **Dataset construction of OpenDV-YouTube with quality check in the loop.** We collect videos from YouTubers with qualified driving videos, and dispose of those with inappropriate viewpoints or involving scene transitions. Then each frame is described with language contexts using VLM followed by keyword checks on texts, such as "words", "watermark", "dark", "blurry", *etc*. Through this process, distorted or entirely black images are wiped out. A classifier tags videos with high-level intentions as commands, incubating the final data corpus of high-quality video-text pairs being 1747 hours long.

**Language Annotation for YouTube Videos.** To create a predictive model that can be controlled by natural language to simulate different futures accordingly, To make the predictive model controllable and improve the sample quality [6], it is crucial to pair the driving videos with meaningful and varied language annotations. We construct two types of texts for OpenDV-YouTube, *i.e.*, driving commands for ego-vehicle and frame descriptions, namely "command" and "context", to help the model comprehend ego actions and open-world concepts, respectively. For commands, we train a video classifier on Honda-HDD-Action [84] for 14 types of actions to label ego behaviors in a 4s sequence. These categorical commands will be further mapped to multiple free-form expressions from a predefined dictionary. For contexts, we leverage an established vision-language model, BLIP-2 [62], to describe the main objects and scenarios for each frame. For more details on annotations, please refer to Appendix C.1.2.

**Enlarging Language Spectrum with Public Datasets.** Considering that BLIP-2 annotations are generated for static frames without comprehension of dynamic driving scenarios such as the traffic light transitions, we exploit several public datasets that provide linguistic descriptions for driving scenarios [12, 13, 24, 51, 72, 84]. However, their metadata is relatively sparse with only a few words such as "sunny road". We further enhance their text quality using GPT [78] to form a descriptive "context" and generate a "command" by categorizing the logged trajectory for each video clip. Ultimately, we integrate these datasets with OpenDV-YouTube to establish OpenDV-2K dataset, as shown in the last row of Tab. 1.

## 3. GenAD Framework

In this section, we introduce the training and design of the GenAD model. As shown in Fig. 4, GenAD is trained in two stages, *i.e.*, image domain transferring and video prediction pre-training. The first stage adapts the general text-to-image model to the driving domain (Sec. 3.1). The second stage lifts the text-to-image model to a video prediction model with our proposed temporal reasoning block and modified training schemes (Sec. 3.2). In Sec. 3.3, we explore how the predictive model can be extended to action-conditioned prediction and planning.

### 3.1. Image Domain Transfer

On-board cameras capture a large field of views with abundant visual contents, including the road, background buildings, surrounding vehicles, *etc*., which require strong and robust generation capability to produce continuous and realistic driving scenarios. To facilitate the learning process, we start with independent image generation in the first stage. Concretely, we initialize our model with SDXL [79], which is a large-scale latent diffusion model (LDM) for text-to-image generation, to leverage its ability to synthesize high-quality images with plenty of visual details. It is implemented as a denoising UNet $\mathbf{f}_\theta$ with several stacked convolution and attention blocks, which learns to synthesize images by denoising the noisy latents [87]. Specifically, given a noisy input latent $\mathbf{x}_t$ corrupted by the forward diffusion process, it is trained to predict the added noise $\boldsymbol{\epsilon}$ of $\mathbf{x}_t$ via the following objective:

$$\mathcal{L}_{\text{img}} := \mathbb{E}_{\mathbf{x},\boldsymbol{\epsilon}\sim\mathcal{N}(0,1),\mathbf{c},t}\left[\|\boldsymbol{\epsilon} - \mathbf{f}_\theta(\mathbf{x}_t;\mathbf{c},t)\|_2^2\right], \quad (1)$$

where $\mathbf{x}$ and $\mathbf{x}_t$ are the clean and noisy latent, respectively, $t$ denotes the timestep for different noise scales, and $\mathbf{c}$ is the text condition that guides the denoising process, which is a concatenation of context and command. For training efficiency, the learning process takes place in a compressed latent space [25, 79, 87] instead of pixel space. During sampling, the model generates images from standard Gaussian noise by denoising the last-step predictions iteratively.

However, the original SDXL is trained on data in the general domain, such as portraits and artistic paintings, which are not concerned with autonomy systems. To adapt

Figure 4. **Framework of GenAD**. (**a**) The two-stage learning for GenAD is composed of transferring the image domain of an image diffusion model to the driving field (a.1 Stage one), and video prediction pre-training for modeling the temporal dependency of videos (a.2 Stage two). (**b**) One transformer block in GenAD for the second stage training has interleaved temporal reasoning blocks before each frozen layer to align spatiotemporal features. (**c**) The proposed Temporal Reasoning Block includes one causal temporal attention (TA) and two decoupled spatial attention (SA) layers to extract features in different axes. A query grid attends to itself as well as blue grids while the dark gray grid is masked out in causal attention. 'Zero init' is appended at the end of each attention block to stabilize training.

the model to synthesize images for driving, we fine-tune it on text-to-image generation using image-text pairs in OpenDV-2K with the same objective as Eq. (1). Following the original training of SDXL, all parameters $\theta$ of the UNet are fine-tuned at this stage, whereas the CLIP text encoders [82] and the autoencoder [25] remain frozen.

### 3.2. Video Prediction Pre-training

In the second stage, with a few frames of a consecutive video as past observations, GenAD is trained to reason about all visual observations and predict several future frames in plausible ways. Similar to the first stage stage, the prediction process can also be guided by text conditions. However, predicting the highly dynamic driving world temporally is challenging due to two fundamental barriers.

1. *Causal Reasoning*: To predict plausible futures following the temporal causality of the driving world, the model needs to comprehend the intentions of all other agents together with the ego vehicle, and understand underlying traffic rules, *e.g.*, how the traffic will change with the transition of traffic lights.
2. *Drastic View Shift*: Contrary to typical video generation benchmarks which mainly have a static background with slow motion of centered objects, the view of driving changes drastically over time. Each pixel in every frame may move to a distant location in the next frame.

We propose temporal reasoning blocks to address these problems. As illustrated in Fig. 4(c), each block is composed of three successive attention layers, *i.e.*, the causal

temporal attention layer and two decoupled spatial attention layers, which are tailored for the causal reasoning and modeling large shifts in the driving scenes, respectively.

**Causal Temporal Attention.** Since the model after the stage-one training can only process each frame independently, we leverage temporal attention to exchange information among different video frames. The attention takes place in the time axis and models the temporal dependency of each grid-wise feature. However, directly adapting bidirectional temporal attention here as [8, 42, 105, 124] can hardly acquire the ability of causal reasoning, since the predictions will be inevitably dependent on the subsequent frames instead of past conditions. Therefore, we restrict the attention direction by adding a causal attention mask, as shown in the last row of Fig. 4(c), to encourage the model to fully exploit knowledge from past observations and faithfully reason about the future as if in real-world driving. We empirically found that the causality constraint greatly regularizes the predicted frames to be coherent with past frames. Following common practice, we also add temporal bias implemented as relative position embeddings on the time axis [94] to distinguish different frames of a sequence for temporal attention.

**Decoupled Spatial Attention.** As driving videos feature fast perspective changes, features in a specific grid could vary greatly in different timesteps and are hard to correlate and learn by temporal attention, which suffers from a limited receptive field. In light of this, we introduce spa-

tial attention to propagate each grid feature in spatial axes to aid in gathering information for temporal attention. We implement a decoupled variant of self-attention for its efficiency with linear computational complexity, compared to quadratic full self-attention. As shown in Fig. 4(c), the two decoupled attention layers propagate features in horizontal and vertical axes, respectively.

**Deep Interaction.** Intuitively, the spatial blocks fine-tuned in stage one refine features of each frame independently towards photorealism, whereas the temporal blocks introduced in stage two align features of all video frames towards coherency and consistency. To further boost the spatiotemporal feature interaction, we interleave the proposed temporal reasoning blocks with the original Transformer blocks in SDXL, *i.e.*, spatial attention, cross attention, and feedforward network, as shown in Fig. 4(b).

**Zero Initialization.** Similar to the previous practices [2, 126], for each block that is newly introduced in stage two, we initialized all parameters of its final layer as zero. This avoids disrupting the prior knowledge of the well-trained image generation model in the beginning and stabilizes the training process.

**Training.** GenAD is trained to predict the future by jointly denoising from the noisy latents with the guidance of past frames and text conditions. We first project $T$ consecutive frames of a video clip into a batch of latents $\mathbf{v} = \{\mathbf{v}^m, \mathbf{v}^n\}$, where the leading $m$ frame latents $\mathbf{v}^m$ are clean, representing historical observations, and other $n = T - m$ frame latents $\mathbf{v}^n$ indicate the future to be predicted. $\mathbf{v}^n$ are then corrupted to $\mathbf{v}_t^n$ by the forward diffusion process, where $t$ indexes a randomly sampled noise scale. The model is trained to predict the noise of $\mathbf{v}_t^n$ conditioned on observations $\mathbf{v}^m$ and text $\mathbf{c}$. The learning objective of the video prediction model is formulated as follows:

$$\mathcal{L}_{\text{vid}} := \mathbb{E}_{\mathbf{v}, \boldsymbol{\epsilon} \sim \mathcal{N}(0,1), \mathbf{c}, t}\left[\|\boldsymbol{\epsilon} - \mathbf{f}_{\theta, \phi}(\mathbf{v}_t^n; \mathbf{v}^m, \mathbf{c}, t)\|_2^2\right], \quad (2)$$

where $\theta$ denotes the inherited stage-one model and $\phi$ represents the newly inserted temporal reasoning blocks. Following [8], we freeze $\theta$ and only train the temporal reasoning blocks to avoid perturbing the generation ability of the image generation model and focus on learning temporal dependencies in videos. Note that only the outputs from the corrupted frames $\mathbf{v}_t^n$ contribute to the training loss while those from condition frames $\mathbf{v}^m$ are ignored. Our training recipe is also readily applicable to video interpolation with minor modifications, *i.e.*, switching the indices of condition frames.

### 3.3. Extensions

Relying on the well-trained video prediction capability in driving scenarios, we further exploit the potential of the pre-trained model in action-controlled prediction and planning,

which are important for real-world driving systems. Here, we explore the downstream tasks on nuScenes [12] which provides recorded poses.

**Action-conditioned Prediction.** To make our predictive model controllable with exact ego actions and act as a simulator [52], we fine-tune the model with the paired future trajectory as an additional condition. Specifically, we map the raw trajectory to a high-dimensional feature with Fourier embeddings [100]. After further projection by a linear layer, it is added to the original conditions. Thus, the ego actions are injected into the network through the conditional cross-attention layer in Fig. 4(b).

**Planning.** By learning to predict the future, GenAD acquires strong representations of complex driving scenes, which can be further exploited for planning. Specifically, we extract spatiotemporal features of two historical frames through the UNet encoder of the *frozen* GenAD, which is nearly half the size of the entire model, and feed them to a multi-layer perceptron (MLP) to predict future waypoints. With the frozen GenAD encoder and a learnable MLP layer, the training process of our planner can be sped up by 3400 times compared to an end-to-end planning model UniAD [47], validating the effectiveness of the learned spatiotemporal feature of GenAD.

## 4. Experiments

### 4.1. Setup and Protocols

GenAD is learned in two stages on OpenDV-2K but with different learning objectives (in Sec. 3) and input formats. In stage one, the model takes input (image, text) pairs and is trained on text-to-image generation. We broadcast the command annotation, which is labeled for each 4s video sequence, to all frames included. The model is trained for 300K iterations on 32 NVIDIA Tesla A100 GPUs with a total batch size of 256. In the second stage, GenAD is trained to jointly denoise future latents conditioned on past latents and texts. Its inputs are (video clip, text) pairs where each video clip is 4s at 2Hz. The current version of GenAD is trained on 64 GPUs for 112.5K iterations with a total batch size of 64. The input frames are resized to $256 \times 448$ for training in both stages, and the text condition $\mathbf{c}$ is dropped at a probability of $p = 0.1$ to enable classifier-free guidance [40] in sampling, which is commonly used in diffusion models to improve sample quality. More training and sampling details are in Appendix D.

### 4.2. Results of Video Prediction Pre-training

**Comparison to Recent Video Generation Approaches.** We compare GenAD to recent advances on an unseen set with geofencing from OpenDV-YouTube, Waymo [98], KITTI [30], and Cityscapes [21] in a *zero-shot* generation

Figure 5. **Task on zero-shot video prediction for unseen scenarios**. We show the generation results (in blue boxes) of different models given the same starting frames. GenAD makes more robust, realistic, and reasonable future predictions on unseen datasets (scenarios). More comparisons (Fig. 14) and visualizations (Fig. 13) are shown in Appendix.

| Method | Training Dataset | Pred. | nuScenes FID (↓) | FVD (↓) |
|---|---|---|---|---|
| DriveGAN [52] | nuScenes | ✓ | 73.4 | 502 |
| DriveDreamer* [104] | | ✓ | 52.6 | 452 |
| DrivingDiffuion* [64] | | ✗ | 15.8 | 332 |
| GenAD-nus (Ours) | nuScenes | ✓ | **15.4** | 244 |
| GenAD (Ours) | OpenDV-2K | ✓ | **15.4** | **184** |

Table 2. **Video generation quality compared to state-of-the-arts trained on nuScenes**. "Pred.": evaluation by future prediction. *: requiring 3D layout inputs.

manner. Fig. 5 depicts the qualitative results. Image-to-video models I2VGen-XL [129] and VideoCrafter1 [15] can not strictly follow the given frames to make predictions, yielding poor consistency between the predicted frames and past frames. The video prediction model DMVFN [46] that is trained on Cityscapes suffers from the unfavorable shape distortions in its predictions, especially on the three unseen datasets. In contrast, GenAD exhibits remarkable zero-shot generalization ability and visual quality although *none* of these sets are included in the training.

**Comparison to nuScenes Experts.** We also compare GenAD with the most recent available driving video generation models which are exclusively trained for nuScenes.



Figure 6. **Task on langauge-conditioned prediction**. Given two frames of a rainy scenario in the intersection and three high-level text conditions, GenAD simulates reasonable futures accordingly.

Tab. 2 shows that GenAD surpasses all previous methods in both image fidelity (FID) and video coherence (FVD). Specifically, GenAD significantly reduces FVD by **44.5%** compared to DrivingDiffusion [64], without taking 3D future layouts as additional inputs. For fair comparisons, we train a model variant (GenAD-nus) on nuScenes dataset only. We find that although GenAD-nus performs on par with GenAD on nuScenes, it struggles to generalize to unseen datasets like Waymo, where the generation degrades

Figure 7. **Case study for model designs**. All components help alleviate artifacts and improve the consistency of future predictions.

| Method | YouTube | | |
| --- | --- | --- | --- |
| | FID ($\downarrow$) | FVD ($\downarrow$) | CLIPSIM ($\uparrow$) |
| Baseline | 18.32 | 244.44 | 0.8405 |
| + Deep Interaction | 17.96 | 201.69 | 0.8409 |
| + Temporal Causality | **16.54** | 207.45 | 0.8550 |
| + Decoupled Spatial Attn. | 17.67 | **189.54** | **0.8652** |

Table 3. **Ablation on model designs in GenAD**. All proposed designs contribute to the final performance.

to the nuScenes visual pattern. In contrast, GenAD trained on OpenDV-2K exhibits strong generalization ability across datasets as shown in Fig. 5.

We provide language-conditioned prediction samples on nuScenes in Fig. 6, where GenAD simulates various futures from the same start following different textual instructions. The impressive generation quality is exhibited in the intricate details of the environment, and the natural transition of ego motion.

**Ablation Study.** We perform ablations by training each variant on a subset of OpenDV-2K for 75K steps. Starting from the baseline with plain temporal attentions [8, 42], we gradually introduce our proposed components. Notably, by interleaving the temporal blocks with the spatial blocks, the FVD significantly improves (-17%) due to more sufficient spatiotemporal interactions. Both temporal causality and decoupled spatial attention contribute to better CLIP-SIM, improving the temporal consistency between future predictions and the condition frames. To be clear, the slight increase in FID and FVD, shown in fourth and third rows of Tab. 3 respectively, does not faithfully reflect a decline in generation quality as discussed in [8, 10, 79]. The effectiveness of each design is shown in Fig. 7.



Figure 8. **Task on action-conditioned prediction (simulation)**. Given the same starting frames and different future trajectories (shown in yellow dots in the first column), GenAD-act can simulate diverse futures following different ego intentions. More visualizations are in Appendix Fig. 15.

| Method | Condition | nuScenes Action Prediction Error ($\downarrow$) |
| --- | --- | --- |
| Ground truth | - | 0.90 |
| GenAD | text | 2.54 |
| GenAD-act | text + traj. | **2.02** |

Table 4. **Task on action-conditioned prediction**. Compared to GenAD with text conditions only, GenAD-act enables more precise future predictions that follow the action condition.

### 4.3. Results of Extensions

**Action-conditioned Prediction.** We further showcase the performance of the action-conditioned model fine-tuned on nuScenes, GenAD-act, in Fig. 8 and Tab. 4. Given two starting frames and a trajectory **w** composed of 6 future waypoints, GenAD-act imagines 6 future frames following the trajectory sequence. To evaluate the consistency between the input trajectory **w** and predicted frames, we establish an inverse dynamics model (IDM) on nuScenes as the evaluator, which projects a video sequence into a corresponding ego trajectory. We leverage the IDM to translate predicted frames into the trajectory $\hat{\mathbf{w}}$, and calculate the L2 distance between **w** and $\hat{\mathbf{w}}$ as the Action Prediction Error. Specifically, GenAD-act substantially reduces the Action Prediction Error by 20.4% compared to GenAD with text condition, allowing for more accurate future simulations.

**Planning Results.** Tab. 5 depicts the planning results on nuScenes where ground truth poses for the ego vehicle are available. By freezing GenAD encoder and only optimizing an additional MLP on top of it, the model can effectively learn to plan. Notably, by pre-extracting image features through the UNet encoder of GenAD, the entire learning process for planning adaptation takes only 10 minutes on a single NVIDIA Tesla V100 device, which is 3400 times more efficient than the training of the UniAD planner [47].

| Method | # Trainable Params. | nuScenes | |
|---|---|---|---|
| | | ADE (↓) | FDE (↓) |
| ST-P3* [45] | 10.9M | 2.65 | 3.73 |
| UniAD* [47] | 58.8M | 1.03 | 1.65 |
| GenAD (Ours) | 0.8M | 1.23 | 2.31 |

Table 5. **Task on open-loop planning**. A lightweight MLP with *frozen* GenAD gets competitive planning results with $73\times$ fewer trainable parameters and front-view image alone. *: multi-view inputs. Evaluation protocols are aligned with UniAD [47].

## 5. Limitations and Discussion

We study the system-level development of GenAD, a large-scale generalized video predictive model for autonomous driving. We also validate the adaptation of the learned representation of GenAD to driving tasks, *i.e.*, learning a "world model" and motion planning. Although we obtain improved generalization to open domains, the increased model capacity poses challenges in both training efficiency and real-time deployment. We envision the unified video prediction task will serve as a scalable objective for future research on representation learning and policy learning. Another interesting direction involves distilling the encoded knowledge for a wider range of downstream tasks [60].

## Acknowledgements

## References

[1] Anurag Ajay, Seungwook Han, Yilun Du, Shaung Li, Abhi Gupta, Tommi Jaakkola, Josh Tenenbaum, Leslie Kaelbling, Akash Srivastava, and Pulkit Agrawal. Compositional foundation models for hierarchical planning. In *NeurIPS*, 2023. 17

[2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning. In *NeurIPS*, 2022. 2, 6

[3] Mohammadhossein Bahari, Saeed Saadatnejad, Ahmad Rahimi, Mohammad Shaverdikondori, Amir Hossein Shahidzadeh, Seyed-Mohsen Moosavi-Dezfooli, and Alexandre Alahi. Vehicle trajectory prediction works, but not everywhere. In *CVPR*, 2022. 2

[4] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *ICCV*, 2021. 18

[5] Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampedro, and Jeff Clune. Video pretraining (VPT): Learning to act by watching unlabeled online videos. In *NeurIPS*, 2022. 17

[6] Fan Bao, Chongxuan Li, Jiacheng Sun, and Jun Zhu. Why are conditional generative models better than unconditional ones? *arXiv preprint arXiv:2212.00362*, 2022. 4

[7] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, Wesam Manassra, Prafulla Dhariwal, Casey Chu, Yunxin Jiao, and Aditya Ramesh. Improving image generation with better captions, 2023. 17

[8] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *CVPR*, 2023. 2, 5, 6, 8, 16, 18

[9] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. RT-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022. 18

[10] Tim Brooks, Janne Hellsten, Miika Aittala, Ting-Chun Wang, Timo Aila, Jaakko Lehtinen, Ming-Yu Liu, Alexei Efros, and Tero Karras. Generating long videos of dynamic scenes. In *NeurIPS*, 2022. 8, 16

[11] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *NeurIPS*, 2020. 18

[12] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yuxin Pan,

Giancarlo Baldan, and Oscar Beijbom. nuScenes: A multi-modal dataset for autonomous driving. In *CVPR*, 2020. 3, 4, 6, 18, 24, 25, 28

[13] Holger Caesar, Juraj Kabzan, Kok Seang Tan, Whye Kit Fong, Eric M. Wolff, Alex Lang, Luke Fletcher, Oscar Beijbom, and Sammy Omari. nuPlan: A closed-loop ml-based planning benchmark for autonomous vehicles. In *CVPR Workshops*, 2021. 3, 4, 24, 25

[14] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*, 2017. 28

[15] Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, Chao Weng, and Ying Shan. VideoCrafter1: Open diffusion models for high-quality video generation. *arXiv preprint arXiv:2310.19512*, 2023. 7, 16

[16] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart-α: Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023. 17

[17] Kai Chen, Enze Xie, Zhe Chen, Yibo Wang, Lanqing Hong, Zhenguo Li, and Dit-Yan Yeung. GeoDiffusion: Text-prompted geometric control for object detection data generation. *arXiv preprint arXiv:2306.04607*, 2023. 17

[18] Li Chen, Penghao Wu, Kashyap Chitta, Bernhard Jaeger, Andreas Geiger, and Hongyang Li. End-to-end autonomous driving: Challenges and frontiers. *arXiv preprint arXiv:2306.16927*, 2023. 2, 18

[19] Yun Chen, Frieda Rong, Shivam Duggal, Shenlong Wang, Xinchen Yan, Sivabalan Manivasagam, Shangjie Xue, Ersin Yumer, and Raquel Urtasun. GeoSim: Realistic video simulation via geometry-aware composition for self-driving. In *CVPR*, 2021. 17

[20] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *CVPR*, 2023. 26

[21] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The Cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016. 3, 6, 18, 28, 30

[22] Yilun Dai, Mengjiao Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Josh Tenenbaum, Dale Schuurmans, and Pieter Abbeel. Learning universal policies via text-guided video generation. In *NeurIPS*, 2023. 2, 17

[23] Pradipto Das, Chenliang Xu, Richard F Doell, and Jason J Corso. A thousand frames in just a few words: Lingual description of videos through latent topics and sparse object stitching. In *CVPR*, 2013. 18

[24] Thierry Deruyttere, Simon Vandenhende, Dusan Grujicic, Luc Van Gool, and Marie-Francine Moens. Talk2Car: Taking control of your self-driving car. In *EMNLP*, 2019. 3, 4, 25

[25] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *CVPR*, 2021. 4, 5

[26] David F Fouhey, Wei-cheng Kuo, Alexei A Efros, and Jitendra Malik. From lifestyle vlogs to everyday interactions. In *CVPR*, 2018. 18

[27] Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, Hongsheng Li, and Yu Qiao. LLaMA-Adapter V2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010*, 2023. 17

[28] Ruiyuan Gao, Kai Chen, Enze Xie, Lanqing Hong, Zhenguo Li, Dit-Yan Yeung, and Qiang Xu. MagicDrive: Street view generation with diverse 3d geometry control. *arXiv preprint arXiv:2310.02601*, 2023. 17

[29] Zeyu Gao, Yao Mu, Ruoyan Shen, Chen Chen, Yangang Ren, Jianyu Chen, Shengbo Eben Li, Ping Luo, and Yanfeng Lu. Enhance sample efficiency and robustness of end-to-end urban autonomous driving via semantic masked world model. *arXiv preprint arXiv:2210.04017*, 2022. 16, 17

[30] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The KITTI dataset. *IJRR*, 2013. 3, 6, 18, 27, 30

[31] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014. 18

[32] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, Miguel Martin, Tushar Nagarajan, Ilija Radosavovic, Santhosh Kumar Ramakrishnan, Fiona Ryan, Jayant Sharma, Michael Wray, Mengmeng Xu, Eric Zhongcong Xu, Chen Zhao, Siddhant Bansal, Dhruv Batra, Vincent Cartillier, Sean Crane, Tien Do, Morrie Doulaty, Akshay Erapalli, Christoph Feichtenhofer, Adriano Fragomeni, Qichen Fu, Abrham Gebreselasie, Cristina González, James Hillis, Xuhua Huang, Yifei Huang, Wenqi Jia, Weslie Khoo, Jáchym Kolář, Satwik Kottur, Anurag Kumar, Federico Landini, Chao Li, Yanghao Li, Zhenqiang Li, Karttikeya Mangalam, Raghava Modhugu, Jonathan Munro, Tullie Murrell, Takumi Nishiyasu, Will Price, Paola Ruiz, Merey Ramazanova, Leda Sari, Kiran Somasundaram, Audrey Southerland, Yusuke Sugano, Ruijie Tao, Minh Vo, Yuchen Wang, Xindi Wu, Takuma Yagi, Ziwei Zhao, Yunyi Zhu, Pablo Arbeláez, David Crandall, Dima Damen, Giovanni Maria Farinella, Christian Fuegen, Bernard Ghanem, Vamsi Krishna Ithapu, C. V. Jawahar, Hanbyul Joo, Kris Kitani, Haizhou Li, Richard Newcombe, Aude Oliva, Hyun Soo Park, James M. Rehg, Yoichi Sato, Jianbo Shi, Mike Zheng Shou, Antonio Torralba, Lorenzo Torresani, Mingfei Yan, and Jitendra Malik. Ego4D: Around the world in 3,000 hours of egocentric video. In *CVPR*, 2022. 18

[33] Jiaxi Gu, Shicong Wang, Haoyu Zhao, Tianyi Lu, Xing Zhang, Zuxuan Wu, Songcen Xu, Wei Zhang, Yu-Gang Jiang, and Hang Xu. Reuse and diffuse: Iterative

denoising for text-to-video generation. *arXiv preprint arXiv:2309.03549*, 2023. 18

[34] Agrim Gupta, Stephen Tian, Yunzhi Zhang, Jiajun Wu, Roberto Martín-Martín, and Li Fei-Fei. MaskViT: Masked visual pre-training for video prediction. In *ICLR*, 2023. 2, 17, 18

[35] Niklas Hanselmann, Katrin Renz, Kashyap Chitta, Apratim Bhattacharyya, and Andreas Geiger. KING: Generating safety-critical driving scenarios for robust imitation via kinematics gradients. In *ECCV*, 2022. 2

[36] William Harvey, Saeid Naderiparizi, Vaden Masrani, Christian Weilbach, and Frank Wood. Flexible diffusion modeling of long videos. In *NeurIPS*, 2022. 18

[37] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 19

[38] Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion models for high-fidelity long video generation. *arXiv preprint arXiv:2211.13221*, 2022. 18

[39] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*, 2017. 28

[40] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS Workshops*, 2021. 6, 27

[41] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 18

[42] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. In *NeurIPS*, 2022. 2, 5, 8

[43] Anthony Hu, Gianluca Corrado, Nicolas Griffiths, Zachary Murez, Corina Gurau, Hudson Yeo, Alex Kendall, Roberto Cipolla, and Jamie Shotton. Model-based imitation learning for urban driving. In *NeurIPS*, 2022. 16, 17

[44] Anthony Hu, Lloyd Russell, Hudson Yeo, Zak Murez, George Fedoseev, Alex Kendall, Jamie Shotton, and Gianluca Corrado. GAIA-1: A generative world model for autonomous driving. *arXiv preprint arXiv:2309.17080*, 2023. 2, 17, 18

[45] Shengchao Hu, Li Chen, Penghao Wu, Hongyang Li, Junchi Yan, and Dacheng Tao. ST-P3: End-to-end vision-based autonomous driving via spatial-temporal feature learning. In *ECCV*, 2022. 9, 27, 28

[46] Xiaotao Hu, Zhewei Huang, Ailin Huang, Jun Xu, and Shuchang Zhou. A dynamic multi-scale voxel flow network for video prediction. In *CVPR*, 2023. 7, 18

[47] Yihan Hu, Jiazhi Yang, Li Chen, Keyu Li, Chonghao Sima, Xizhou Zhu, Siqi Chai, Senyao Du, Tianwei Lin, Wenhai Wang, Lewei Lu, Xiaosong Jia, Qiang Liu, Jifeng Dai, Yu Qiao, and Hongyang Li. Planning-oriented autonomous driving. In *CVPR*, 2023. 6, 8, 9, 27, 28

[48] Fan Jia, Weixin Mao, Yingfei Liu, Yucheng Zhao, Yuqing Wen, Chi Zhang, Xiangyu Zhang, and Tiancai Wang. Adriver-i: A general world model for autonomous driving. *arXiv preprint arXiv:2311.13549*, 2023. 2

[49] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. 18

[50] Tarasha Khurana, Peiyun Hu, David Held, and Deva Ramanan. Point cloud forecasting as a proxy for 4d occupancy forecasting. In *CVPR*, 2023. 16

[51] Jinkyu Kim, Teruhisa Misu, Yi-Ting Chen, Ashish Tawari, and John Canny. Grounding human-to-vehicle advice for self-driving vehicles. In *CVPR*, 2019. 3, 4, 24, 25

[52] Seung Wook Kim, Jonah Philion, Antonio Torralba, and Sanja Fidler. DriveGAN: Towards a controllable high-quality neural simulation. In *CVPR*, 2021. 2, 6, 7, 17

[53] Seung Wook Kim, Bradley Brown, Kangxue Yin, Karsten Kreis, Katja Schwarz, Daiqing Li, Robin Rombach, Antonio Torralba, and Sanja Fidler. NeuralField-LDM: Scene generation with hierarchical latent diffusion models. In *CVPR*, 2023. 17

[54] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 18

[55] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollar, and Ross Girshick. Segment anything. In *ICCV*, 2023. 2

[56] Lei Lai, Zhongkai Shangguan, Jimuyang Zhang, and Eshed Ohn-Bar. XVO: Generalized visual odometry via cross-modal self-training. In *ICCV*, 2023. 28

[57] Yann LeCun. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27. *Open Review*, 62, 2022. 2, 17

[58] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 2015. 2

[59] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A multi-modal model with in-context instruction tuning. *arXiv preprint arXiv:2305.03726*, 2023. 17

[60] Daiqing Li, Huan Ling, Amlan Kar, David Acuna, Seung Wook Kim, Karsten Kreis, Antonio Torralba, and Sanja Fidler. DreamTeacher: Pretraining image backbones with deep generative models. In *ICCV*, 2023. 9, 17

[61] Hongyang Li, Yang Li, Huijie Wang, Jia Zeng, Pinlong Cai, Dahua Lin, Junchi Yan, Feng Xu, Lu Xiong, Jingdong Wang, Futang Zhu, Kai Yan, Chunjing Xu, Tiancai Wang, Beipeng Mu, Shaoqing Ren, Zhihui Peng, and Yu Qiao. Open-sourced data ecosystem in autonomous driving: the present and future. 2023. 18

[62] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, 2023. 2, 3, 4, 17, 19

[63] Xin Li, Wenqing Chu, Ye Wu, Weihang Yuan, Fanglong Liu, Qi Zhang, Fu Li, Haocheng Feng, Errui Ding, and Jingdong Wang. VideoGen: A reference-guided latent diffusion approach for high definition text-to-video generation. *arXiv preprint arXiv:2309.00398*, 2023. 18

[64] Xiaofan Li, Yifu Zhang, and Xiaoqing Ye. DrivingDiffusion: Layout-guided multi-view driving scene video generation with latent diffusion model. *arXiv preprint arXiv:2310.07771*, 2023. 7, 17

[65] Zhenyu Li, Zehui Chen, Ang Li, Liangji Fang, Qinhong Jiang, Xianming Liu, and Junjun Jiang. Unsupervised domain adaptation for monocular 3d object detection via self-training. In *ECCV*, 2022. 3

[66] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014. 17

[67] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023. 17

[68] Yaofang Liu, Xiaodong Cun, Xuebo Liu, Xintao Wang, Yong Zhang, Haoxin Chen, Yang Liu, Tieyong Zeng, Raymond Chan, and Ying Shan. Evalcrafter: Benchmarking and evaluating large video generation models. *arXiv preprint arXiv:2310.11440*, 2023. 16

[69] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2018. 27

[70] Haoyu Lu, Guoxing Yang, Nanyi Fei, Yuqi Huo, Zhiwu Lu, Ping Luo, and Mingyu Ding. VDT: An empirical study on video diffusion with transformers. *arXiv preprint arXiv:2305.13311*, 2023. 18

[71] Jiachen Lu, Ze Huang, Jiahui Zhang, Zeyu Yang, and Li Zhang. Wovogen: World volume-aware diffusion for controllable multi-camera driving scene generation. *arXiv preprint arXiv:2312.02934*, 2023. 2

[72] Jiageng Mao, Minzhe Niu, Chenhan Jiang, Hanxue Liang, Xiaodan Liang, Yamin Li, Chao Ye, Wei Zhang, Zhenguo Li, Jie Yu, Hang Xu, and Chunjing Xu. One million scenes for autonomous driving: ONCE dataset. In *NeurIPS Datasets and Benchmarks*, 2021. 3, 4, 24, 25

[73] Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In *CVPR*, 2023. 17

[74] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. HowTo100M: Learning a text-video embedding by watching hundred million narrated video clips. In *ICCV*, 2019. 18

[75] Arsha Nagrani, Paul Hongsuck Seo, Bryan Seybold, Anja Hauth, Santiago Manen, Chen Sun, and Cordelia Schmid. Learning audio-video modalities from image captions. In *ECCV*, 2022. 18

[76] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. GLIDE: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 18

[77] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 2, 17

[78] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022. 3, 4, 22, 23

[79] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 2, 4, 8, 17, 26

[80] Xiaojuan Qi, Zhengzhe Liu, Qifeng Chen, and Jiaya Jia. 3D motion decomposition for rgbd future dynamic scene synthesis. In *CVPR*, 2019. 2

[81] Shengyi Qian, Linyi Jin, Chris Rockwell, Siyi Chen, and David F Fouhey. Understanding 3d object articulation in internet videos. In *CVPR*, 2022. 18

[82] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 2, 5, 18, 26, 28

[83] Ilija Radosavovic, Tete Xiao, Stephen James, Pieter Abbeel, Jitendra Malik, and Trevor Darrell. Real-world robot learning with masked visual pre-training. In *CoRL*, 2023. 17

[84] Vasili Ramanishka, Yi-Ting Chen, Teruhisa Misu, and Kate Saenko. Toward driving scene understanding: A dataset for learning driver behavior and causal reasoning. In *CVPR*, 2018. 3, 4, 19, 25, 26

[85] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 18

[86] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *ICML*, 2015. 18

[87] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 2, 4, 17, 18

[88] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Raphael Gontijo-Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. In *NeurIPS*, 2022. 18

[89] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022. 17

[90] Ramon Sanabria, Ozan Caglayan, Shruti Palaskar, Desmond Elliott, Loïc Barrault, Lucia Specia, and Florian Metze. How2: a large-scale dataset for multimodal language understanding. *arXiv preprint arXiv:1811.00347*, 2018. 18

[91] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and

Jenia Jitsev. LAION-5B: An open large-scale dataset for training next generation image-text models. In *NeurIPS Datasets and Benchmarks*, 2022. 18

[92] Dhruv Shah, Ajay Sridhar, Nitish Dashora, Kyle Stachowicz, Kevin Black, Noriaki Hirose, and Sergey Levine. ViNT: A foundation model for visual navigation. In *CoRL*, 2023. 17

[93] Dandan Shan, Jiaqi Geng, Michelle Shu, and David F Fouhey. Understanding human hands in contact at internet scale. In *CVPR*, 2020. 18

[94] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. *arXiv preprint arXiv:1803.02155*, 2018. 5

[95] Chonghao Sima, Katrin Renz, Kashyap Chitta, Li Chen, Hanxue Zhang, Chengen Xie, Ping Luo, Andreas Geiger, and Hongyang Li. Drivelm: Driving with graph visual question answering. *arXiv preprint arXiv:2312.14150*, 2023. 18

[96] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 18, 27

[97] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012. 18

[98] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott M. Ettinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. Scalability in perception for autonomous driving: Waymo open dataset. In *CVPR*, 2020. 3, 6, 27, 30

[99] Alexander Swerdlow, Runsheng Xu, and Bolei Zhou. Street-view image generation from a bird's-eye view layout. *arXiv preprint arXiv:2301.04634*, 2023. 17

[100] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. In *NeurIPS*, 2020. 6, 26

[101] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges. *arXiv preprint arXiv:1812.01717*, 2018. 16, 28

[102] Vikram Voleti, Alexia Jolicoeur-Martineau, and Chris Pal. MCVD-masked conditional video diffusion for prediction, generation, and interpolation. In *NeurIPS*, 2022. 18

[103] Jinpeng Wang, Yixiao Ge, Rui Yan, Yuying Ge, Kevin Qinghong Lin, Satoshi Tsutsui, Xudong Lin, Guanyu Cai, Jianping Wu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. All in one: Exploring unified video-language pre-training. In *CVPR*, 2023. 17

[104] Xiaofeng Wang, Zheng Zhu, Guan Huang, Xinze Chen, and Jiwen Lu. DriveDreamer: Towards real-world-driven world models for autonomous driving. *arXiv preprint arXiv:2309.09777*, 2023. 2, 7, 17, 18

[105] Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan He, Jiashuo Yu, Peiqing Yang, Yuwei Guo, Tianxing Wu, Chenyang Si, Yuming Jiang, Cunjian Chen, Chen Change Loy, Bo Dai, Dahua Lin, Yu Qiao, and Ziwei Liu. LAVIE: High-quality video generation with cascaded latent diffusion models. *arXiv preprint arXiv:2309.15103*, 2023. 5, 16, 18

[106] Yuqi Wang, Jiawei He, Lue Fan, Hongxin Li, Yuntao Chen, and Zhaoxiang Zhang. Driving into the future: Multiview visual forecasting and planning with world model for autonomous driving. *arXiv preprint arXiv:2311.17918*, 2023. 2

[107] Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinyuan Chen, Yaohui Wang, Ping Luo, Ziwei Liu, Yali Wang, Limin Wang, and Yu Qiao. InternVid: A large-scale video-text dataset for multimodal understanding and generation. *arXiv preprint arXiv:2307.06942*, 2023. 18

[108] Zhou Wang, Eero P Simoncelli, and Alan C Bovik. Multiscale structural similarity for image quality assessment. In *ACSSC*, 2003. 16

[109] Xinshuo Weng, Jianren Wang, Sergey Levine, Kris Kitani, and Nicholas Rhinehart. Inverting the pose forecasting pipeline with SPF2: Sequential pointcloud forecasting for sequential pose forecasting. In *CoRL*, 2021. 16

[110] Benjamin Wilson, William Qi, Tanmay Agarwal, John Lambert, Jagjeet Singh, Siddhesh Khandelwal, Bowen Pan, Ratnesh Kumar, Andrew Hartnett, Jhony Kaesemodel Pontes, Deva Ramanan, Peter Carr, and James Hays. Argoverse 2: Next generation datasets for self-driving perception and forecasting. In *NeurIPS Datasets and Benchmarks*, 2021. 3

[111] Daniel M Wolpert, Zoubin Ghahramani, and Michael I Jordan. An internal model for sensorimotor integration. *Science*, 1995. 2

[112] Penghao Wu, Li Chen, Hongyang Li, Xiaosong Jia, Junchi Yan, and Yu Qiao. Policy pre-training for autonomous driving via self-supervised geometric modeling. In *ICLR*, 2023. 18

[113] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui. The rise and potential of large language model based agents: A survey. *arXiv preprint arXiv:2309.07864*, 2023. 2

[114] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezatofighi, and Dacheng Tao. GMFlow: Learning optical flow via global matching. In *CVPR*, 2022. 19

[115] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezatofighi, Fisher Yu, Dacheng Tao, and Andreas Geiger. Unifying flow, stereo and depth estimation. *IEEE TPAMI*, 2023. 19

[116] Yinghao Xu, Menglei Chai, Zifan Shi, Sida Peng, Ivan Skorokhodov, Aliaksandr Siarohin, Ceyuan Yang, Yujun Shen,

Hsin-Ying Lee, Bolei Zhou, and Sergey Tulyakov. Dis-CoScene: Spatially disentangled generative radiance fields for controllable 3d-aware scene synthesis. In *CVPR*, 2023. 17

[117] Hongwei Xue, Tiankai Hang, Yanhong Zeng, Yuchong Sun, Bei Liu, Huan Yang, Jianlong Fu, and Baining Guo. Advancing high-resolution video-language representation with large-scale video transcriptions. In *CVPR*, 2022. 18

[118] Kairui Yang, Enhui Ma, Jibin Peng, Qing Guo, Di Lin, and Kaicheng Yu. BEVControl: Accurately controlling street-view elements with multi-perspective consistency via bev sketch layout. *arXiv preprint arXiv:2308.01661*, 2023. 17

[119] Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Dale Schuurmans, and Pieter Abbeel. Learning interactive real-world simulators. *arXiv preprint arXiv:2310.06114*, 2023. 17

[120] Ze Yang, Yun Chen, Jingkang Wang, Sivabalan Manivasagam, Wei-Chiu Ma, Anqi Joyce Yang, and Raquel Urtasun. UniSim: A neural closed-loop sensor simulator. In *CVPR*, 2023. 17

[121] Zetong Yang, Li Chen, Yanan Sun, and Hongyang Li. Visual point cloud forecasting enables scalable autonomous driving. In *CVPR*, 2024. 17

[122] Lijun Yu, Yong Cheng, Kihyuk Sohn, José Lezama, Han Zhang, Huiwen Chang, Alexander G. Hauptmann, Ming-Hsuan Yang, Yuan Hao, Irfan Essa, and Lu Jiang. MAGVIT: Masked generative video transformer. In *CVPR*, 2023. 18

[123] Rowan Zellers, Ximing Lu, Jack Hessel, Youngjae Yu, Jae Sung Park, Jize Cao, Ali Farhadi, and Yejin Choi. MERLOT: Multimodal neural script knowledge models. In *NeurIPS*, 2021. 18

[124] David Junhao Zhang, Jay Zhangjie Wu, Jia-Wei Liu, Rui Zhao, Lingmin Ran, Yuchao Gu, Difei Gao, and Mike Zheng Shou. Show-1: Marrying pixel and latent diffusion models for text-to-video generation. *arXiv preprint arXiv:2309.15818*, 2023. 5, 16

[125] Jimuyang Zhang, Ruizhao Zhu, and Eshed Ohn-Bar. SelfD: self-learning large-scale driving policies from the web. In *CVPR*, 2022. 18

[126] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, 2023. 6

[127] Qihang Zhang, Zhenghao Peng, and Bolei Zhou. Learning to drive by watching youtube videos: Action-conditioned contrastive policy pretraining. In *ECCV*, 2022. 18

[128] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 16

[129] Shiwei Zhang, Jiayu Wang, Yingya Zhang, Kang Zhao, Hangjie Yuan, Zhiwu Qin, Xiang Wang, Deli Zhao, and Jingren Zhou. I2VGen-XL: High-quality image-to-video synthesis via cascaded diffusion models. *arXiv preprint arXiv:2311.04145*, 2023. 7, 16

[130] Wenliang Zhao, Lujia Bai, Yongming Rao, Jie Zhou, and Jiwen Lu. UniPC: A unified predictor-corrector framework for fast sampling of diffusion models. In *NeurIPS*, 2023. 17

[131] Wenliang Zhao, Yongming Rao, Zuyan Liu, Benlin Liu, Jie Zhou, and Jiwen Lu. Unleashing text-to-image diffusion models for visual perception. In *ICCV*, 2023. 17

[132] Jilai Zheng, Chao Ma, Houwen Peng, and Xiaokang Yang. Learning to track objects from unlabeled videos. In *ICCV*, 2021. 19

[133] Luowei Zhou, Chenliang Xu, and Jason Corso. Towards automatic learning of procedures from web instructional videos. In *AAAI*, 2018. 18

[134] Yunsong Zhou, Linyan Huang, Qingwen Bu, Jia Zeng, Tianyu Li, Hang Qiu, Hongzi Zhu, Minyi Guo, Yu Qiao, and Hongyang Li. Embodied understanding of driving scenarios. *arXiv preprint arXiv:2403.04593*, 2024. 18

[135] Ruizhao Zhu, Peng Huang, Eshed Ohn-Bar, and Venkatesh Saligrama. Learning to drive anywhere. In *CoRL*, 2023. 2

# Appendix

# A. Discussions

To assist a better understanding of our work, we supplement discussions on intuitive questions that one may raise.

**Q1.** *Why do you propose the training target of the predictive model as videos?*

Video is a particularly universal and scalable target given a wealth of uncalibrated driving videos. Different from BEV representations [29, 43] that require camera extrinsic parameters and point clouds [50, 109] that are restricted by different LiDAR configurations, video prediction can be performed in a pose-agnostic manner. This characteristic offers significant advantages in scalability to more diverse data sources, which are key for the generalization ability of the learned model.

**Q2.** *Why do you predict multiple frames simultaneously with historical frames as input? How about alternatively using an auto-regressive design, i.e., predicting future frames one by one?*

Indeed, auto-regressive prediction can further stabilize the prediction process by leveraging conditional dependencies on previously generated frames, thereby enhancing consistency. Nevertheless, we still choose to employ a joint denoising procedure for two primary reasons. To start, diffusion models are typically computationally expensive, and our model is no exception. For videos comprised of multiple frames, predicting them auto-regressively would multiply their computational intensity, making it inefficient for implementation and deployment.

Moreover, conducting auto-regressive predictions makes it challenging to effectively apply conditions that require significant changes. Consider the scenario of a driver making a turn, which typically takes several seconds and involves a long sequence of frames. If the prediction duration is too short, the model may struggle to follow the given instructions, as it is impossible to achieve substantial changes within a single frame. Instead, it might simply continue the tendency of previously determined frames and completely disregard the provided instructions. Therefore, joint prediction also allows us to effectively apply complex controls and facilitate more coherent action generation.

**Q3.** *What is the criterion to prove good generalization ability of your model? How much data do we need to guarantee generalization?*

Currently, it is hard to define a specific criterion to assess the generalization ability of our predictive models for the reason that the quality judgment is subjective [68] and it is impossible to find an aligned method that is available to compare. However, through our exhaustive exploitation of public data, we have discovered that increasing the scale of the data is advantageous for zero-shot generation on existing datasets. It is also important to note that our method is easily scalable, offering opportunities to continuously enhance its generalization ability by leveraging vast amounts of unlabeled data.

**Q4.** *Why not evaluate models using typical video prediction metrics? What are the appropriate metrics to evaluate the performance of the driving video prediction model with multiple conditions?*

Common practices in the task of video prediction use Structural Similarity Index Measure (SSIM) [108] and a perceptual metric LPIPS [128] for quantitative evaluation. These two metrics calculate frame-wise similarities between predicted frames and corresponding ground truth frames. They are designed to assess the model's ability to *exactly* follow the recorded events. Consequently, models optimized for these metrics tend to copy certain patterns and could overfit the small datasets adopted, thereby restricting their potential for diverse future generations. This limitation is particularly problematic for predictive models in driving scenarios, where multiple futures may occur and proactive preparation is essential for each of these cases.

We sought to use the distribution-based metrics, including FVD [101] and CLIPSIM which are widely adopted by diffusion-based generation approaches [8, 105, 124]. However, for the image-to-video generation models [15, 129] in our comparison, they do not directly compose the input image as any specific video frames they generate, mainly preserving semantics and contents from input images. Thus, it becomes challenging to align the comparison settings with ours for metrics like FVD, which measures the distribution distance of consecutive frames, or CLIPSIM, which can be used to evaluate the semantic similarity between the conditional frame and generated frames. Moreover, these metrics are not perfect. For instance, FVD could be blind to unrealistic repetition and prefer small-scale motion, as discussed in [8, 10].

In short, from the existing metrics, it is hard to quantitatively evaluate the prediction abilities of a generalized model for real-world driving, which encompasses multi-modal conditions and requires temporal consistency. There still needs to be effort made to design an appropriate metric that can effectively evaluate such models.

**Q5. Broader impact.** *What are potential applications and future directions with the provided large-scale OpenDV-2K data and the GenAD model, for both academia and industry?*

To the best of our knowledge, OpenDV-2K is the largest available data corpus that we can collect from public sources. It significantly enhances the quantity and diversity of driving video footage in multiple dimensions, providing the research community with a massive high-quality resource for exploring open avenues in autonomous driving. In addition to video prediction, we hope our dataset can also benefit the community to enable broader applications [22, 83, 92, 103, 119].

In this work, we have demonstrated that the strong representation of GenAD can be beneficial for planning. Similarly, it is also promising to adapt it to a broader range of downstream tasks such as perception [131]. To improve the flexibility and efficiency for deployment, transferring the knowledge of generative models via distillation [60] is also worth investigating. Except for its powerful representations, the prediction futures conditioned on actions also open the opportunities for model-predictive control [34, 57] and inverse dynamics model [1, 5, 22] to enable trajectory planning, which are beyond the scope of this paper. Note that our model will be made publicly available to benefit the community and it is flexible to further fine-tune it on in-house data for the industry.

**Q6. Limitations.** *What are the issues with current designs, and corresponding preliminary solutions?*

It is known that the captions used for training have a great impact on generation quality [7, 16]. Currently, the context description of OpenDV-2K is automatically annotated by BLIP-2 [62]. However, we empirically find that the generated captions have two main limitations. First, the BLIP-2 captions tend to be short and plain, lacking enough details about the complicated driving scenes and becoming indistinguishable from one another. Second, the alignment between the image and caption still needs to be improved. The BLIP-2 captions are mostly centric on a single object and thus fail to include the majority of important content in the scene. In addition, affected by its fine-tuning samples [66], BLIP-2 is unaware of the state of the image observer itself. Hence, it fails to infer ego intentions, which may lead to conflicts with high-level commands. To overcome these limitations, it is promising to utilize more advanced vision-language models that have a more comprehensive understanding and text-rich description of the whole scene [27, 67, 77], and have temporal awareness [59].

We opt for SDXL [79] as our starting point to inherit its merits in high quality of visual details, large capacity of model size, and better rendering abilities of text encoders. On the other hand, we have noticed that SDXL is slow to sample and computationally expensive. Our model does suffer from that as well. However, as a pioneering work exploring how to build a generalized predictive model on internet-scale driving data, the main focus of this work is the generalization ability to diverse unseen driving scenarios instead of computation overhead. Future works may include trying faster sampling methods [73, 89, 130] and transferring our general recipe to more efficient diffusion models [87].

While there is not a silver bullet yet we hope that future work takes a deep and grounded look at these discussions, identifying what more downstream applications could be applied - and more importantly, why they work or fail. Our hope is that GenAD serves as a starting point, as the main paper argues, a generalized video pre-trained paradigm that is built on top of the largest available driving videos and excels at a wide spectrum of autonomous driving tasks.

## B. Related Work

Related work is introduced below due to the limited space in the main paper.

### B.1. Driving Scene Generation

Over the past few years, scene generation has gained increasing popularity due to its importance for safety-critical domains like autonomous driving. One family of works [19, 116, 120, 121] perform 3D-aware rendering for sensor simulation. In particular, GeoSim [19] augments the existing images by borrowing objects from other scenes and rendering them at novel poses. UniSim [120] creates digital twins of driving logs with manipulable foreground objects to enable close-loop simulation. However, these methods can only manipulate objects from the collected assets, and novel objects cannot be created unless further collected. Recent advancements [17, 28, 53, 64, 99, 104, 118] use diffusion models to synthesize scenes with novel content beyond the collected data. As a dual task of perception, several works simulate realistic sensor data controlled by input layouts such as 2D bounding boxes [17] and bird's-eye view (BEV) segmentation maps [53, 99, 118]. More recent works [28, 64, 104] choose 3D bounding boxes for better geometry control. These methods also have potential to serve as data engine [17, 28, 64, 99, 118, 120], *i.e.*, the simulated data can be further adopted as augmented samples to boost the performance of existing perception models. However, their control abilities are acquired from manually annotated datasets, preventing them from scaling to more unlabelled data and increasing both diversity and generalization.

Besides layout-controlled sensor simulation, another thread of progresses [29, 43, 44, 52, 104] focuses on simulating the temporal dynamics of the driving scenarios. Specifically, MILE [43] firstly introduces a model of the world incorporating the BEV representation. By imagining the world within the designed space, the world dynamics can be implicitly encoded

17

and the behaviors of vehicles can be interpretably decoded. This opens the opportunity for executing planning policies without having access to real observations. Differently, inspired by the advances in video generation, DriveDreamer [104] and GAIA-1 [44] propose to build a realistic world model in the form of video frames. Particularly, GAIA-1 is scaled up to about 10B model parameters on 4700 hours of in-house videos, showing highly appealing results. However, the diversity of their generation is still limited by the datasets they adopt. To be specific, the nuScenes [12] used by DriveDreamer is collected in Singapore and Boston, while GAIA-1's driving logs are recorded within London. Both of them use fixed or similar camera settings. The distribution of their data sources limits their generalization abilities to unseen scenarios, different camera poses, and other settings. Moreover, how to utilize the learned knowledge for downstream applications, *e.g.*, planning, is still rarely mentioned and explored.

## B.2. Video Generation and Prediction

Video generation and prediction are effective ways to model the real world. Several practices [34, 46, 122] have been made to synthesize future driving videos. With the renaissance of diffusion models [41, 96], recent progresses [76, 85, 87, 88] have demonstrated that diffusion models show a great advantage over other generative methods [31, 54, 86] in both fidelity and diversity. These advantages have also been extended to the temporal domain by numerous works in video generation [8, 33, 38, 63, 105]. Among them, many works [8, 36, 70, 102] include public driving datasets [12, 21, 30] as touchstones for their evaluation. However, none of these methods have proposed effective designs that are specialized for driving scenarios, which are known to be more complex and challenging [34] as we discussed in the main paper. In addition, due to their exclusive training strategy, the model capability is greatly limited by each small and simple dataset [12, 21, 30], hindering the generalization ability to diverse driving scenes in the real world. In contrast, we explore the first practice of building a generalized prediction model via training on large-scale driving videos in a joint manner.

## B.3. Learning from Web Driving Videos

Learning the general capabilities from large-scale data has been well studied in the field of both vision and language [11, 82, 91]. It is also promising to exploit the internet-scale videos for autonomous driving. However, due to the unlabeled nature of the web data, there exist great challenges and there are only a few methods that leverage this idea to driving tasks for different purposes. SelfD [125] learns driving policies via semi-supervised learning on YouTube videos. The policy network is pre-trained with pseudo trajectories and then transferred to the target datasets via fine-tuning. Instead of directly pre-training the policy, ACO [127] introduces an action contrastive learning method to obtain action-related representations for downstream tasks. However, both SelfD and ACO rely on pseudo-labeling of trajectories or actions on vast amounts of driving videos. This could be highly sensitive to domain changes, thus compromising their reliability. More recently, PPGeo [112] proposes a fully self-supervised learning pipeline to learn a motion-aware encoder through geometric reconstruction. The encoder can be further fine-tuned to benefit downstream tasks. However, their pipeline requires separating each component into different training stages. Instead, our method directly conducts self-supervised learning via future prediction, which is more intuitive and flexible. This allows us to easily apply it to such massive and diverse uncalibrated driving videos for the first time. In addition, our predictions generate interpretable visual outputs that implicitly perform the planning process and seamlessly serve as a real-world driving simulator.

## B.4. Video Datasets from the Internet

Large-scale datasets have been proven to be a core component for generalizable foundation models [82, 91]. For video tasks, collecting data in laboratories or through crowd-sourcing is a common strategy for specific tasks, such as robotics [9] and ego-centric perception [32]. However, the collection and annotation process is costly and hard to scale. Therefore, researchers have sought YouTube or similar websites as video sources as they cover diverse topics and environments, and support academic usage licenses. For example, some pioneering works manually annotate YouTube videos for action classifications [49, 81, 97], action descriptions or captions [23, 133], and hand-object intersections [26, 93]. Recently, researchers have begun to leverage alt-text [4], automatic speech recognition [74, 117, 123], original image captions [75], or paired subtitles [74, 90] to enlarge the annotation scale for video captions. With the development of foundation models, Wang *et al.* [107] employ image captioning models and language models to generate video captions. These video-text pairs have demonstrated great help for general-domain video-language pre-training. ACO [127] and SelfD [125] are the only two that collect 120 and 100 hours of driving videos from YouTube, respectively, to pre-train an encoder for policy learning (Details in Appendix B.3). In contrast, we exhaustively mine driving videos from YouTube and construct the largest driving video datasets publicly available, accumulating over 1700 hours. Besides, our videos are paired with descriptions and command labels which can be used for broader applications such as language-guided autonomous driving [18, 61, 95, 134].

## C. OpenDV-2K Dataset

Our data suite, OpenDV-2K, the *largest* public driving dataset to date, contains 2059 hours of driving video along with diverse text conditions, including *contexts* and *commands*. In this section, we detail the YouTube video collection process (Appendix C.1.1), language annotation method (Appendix C.1.2 for OpenDV-YouTube and Appendix C.2 for other public datasets), more examples and analysis to illustrate the diversity of OpenDV-2K (Appendix C.1.3 and Appendix C.1.4).

### C.1. OpenDV-YouTube

#### C.1.1  Data Collection

**Data Acquisition.**  We first search for videos of driving tours on YouTube and select 43 video uploaders worldwide, *i.e.*, YouTubers, who continuously post high-quality driving videos. We further check the quality of videos from these YouTubers in terms of resolution, frame rate, scene transition frequency, *etc.*, resulting in 2139 high-quality front-view driving videos. We take all videos from 3 selected YouTubers as the validation set, including `Pete Drives USA`, `KenoVelicanstveni`, and `Driving Experience`, while the other videos are used for training. We illustrate the diversity of the OpenDV-YouTube in Fig. 9.

**Format Conversion.**  To simplify the data usage for training both image and video models, We pre-process all videos into sets of consecutive frames in image format using `decord` and `opencv` packages. We sample videos with resolutions no less than 720p (*e.g.*, 1280×720 for 16 : 9 videos) at 10Hz.

**Data Cleaning.**  To ensure the quality of our dataset, we exclude non-driving frames which are commonly shown in each video and introduce unwanted noise. Specifically, we discard the first 90 seconds and the last 30 seconds for most videos to remove the channel introduction at the beginning and the subscription reminder at the end. For YouTubers with longer video introductions, we discard the first 180 or 300 seconds from their videos. We further detect and remove black frames and transition frames with the help of vision-language models. We first search for frames with phrases like `words`, `watermark`, `dark night`, `dark street`, and `blur` in their BLIP-2 [62] -generated contexts, followed by the manual quality check to determine their removal. For details on BLIP-2 descriptions, please refer to Appendix C.1.2.

#### C.1.2  Language Annotation

Our OpenDV-YouTube possesses two types of annotations, frame descriptions (contexts) and ego-driver commands. The context aims to benefit text-to-image learning, helping the model understand the concepts of open-world objects and scenarios, whereas the command is designed to correlate the future predictions with ego actions and further enables the language as control signals. We show some examples in Fig. 10 and introduce the annotation method below.

**Frame Descriptions (Contexts).**  We leverage the established BLIP-2 [62] to describe the main objects or scenarios in each frame with the following prompt. The language annotations are also used in data cleaning, as mentioned in Appendix C.1.1.

`Prompt` = "Question: Describe the image of a driving scenario concisely. Answer: "

Table 6. **BLIP-2 Prompt for generating context of each frame**.

**Driver Commands.**  Similar to the conventional behavior planning approach [84], we classify the commands for ego vehicle into 13 categories, *i.e.*, {`forward, intersection passing, left turn, right turn, left lane change, right lane change, left lane branch, right lane branch, crosswalk passing, rail passing, merge, U-Turn, stop/decelerate, deviate`}. We train an action model based on optical flow to annotate the command for the unlabeled YouTube dataset. Specifically, we leverage the pre-trained GMFlow [114, 115] to extract optical flow between adjacent frames of a driving video sequence. Taking as input both the optical flow and its distance map [132], we train a ResNet-18 [37] to classify the action of each 4s video clip. The training is conducted on the merged dataset of Honda-HDD-Action and Honda-HDD-Cause [84], which provides specified action annotations. For each type of action, we match it with multiple expressions to enrich language understanding. During training, we randomly select one text from the matched caption set for each action. The dictionary for paraphrasing is shown in Tab. 8 below.

Figure 9. **Diverse video samples in OpenDV-YouTube**. We only showcase certain frames from videos due to space limits. OpenDV-YouTube covers a wide spectrum of diversity in multiple axes, including geographic locations, traffic scenarios, time periods, weather conditions, *etc*. We strictly construct the Train/Val split from different YouTubers for zero-shot evaluation.

Figure 10. **Examples of language annotations for different data sources in OpenDV-2K**. We unify the paired text as command and context for all data sources after careful pre-processing. The command represents the action of the ego vehicle, whereas the context covers various aspects of information in the driving scenario. For details on how to merge public driving datasets, please refer to Appendix C.2.

```
command_caption_dict = {

    0: [    "Move forward.",           "Move steady.",                "Go forward.",
            "Go straight.",            "Proceed.",                    "Drive forward.",
            "Drive straight.",         "Drive steady.",               "Keep the direction.",
            "Maintain the direction."                                                            ],

    1: [    "Pass the intersection.",          "Cross the intersection.",         "Traverse the intersection.",
            "Drive through the intersection.", "Move past the intersection.",     "Pass the junction.",
            "Cross the junction.",             "Traverse the junction.",          "Drive through the junction.",
            "Move past the junction.",         "Pass the crossroad.",             "Cross the crossroad.",
            "Traverse the crossroad.",         "Drive through the crossroad.",    "Move past the crossroad."        ],

    2: [    "Turn left.",              "Turn to the left.",           "Make a left turn.",
            "Take a left turn.",       "Turn to the left.",           "Left turn.",
            "Steer left.",             "Steer to the left."                                                         ],

    3: [    "Turn right.",             "Turn to the right.",          "Make a right turn.",
            "Take a right turn.",      "Turn to the right.",          "Right turn.",
            "Steer right.",            "Steer to the right."                                                        ],

    4: [    "Make a left lane change.",        "Change to the left lane.",        "Switch to the left lane.",
            "Shift to the left lane.",         "Move to the left lane."                                            ],

    5: [    "Make a right lane change.",       "Change to the right lane.",       "Switch to the right lane.",
```

```
                    "Shift to the right lane.",        "Move to the right lane."                              ],

        6:  [    "Go to the left lane branch.",      "Take the left lane branch.",      "Move into the left lane branch.",
                 "Follow the left lane branch.",     "Follow the left side road."                              ],

        7:  [    "Go to the right lane branch.",     Take the right lane branch.,       "Move into the right lane branch.",
                 "Follow the right lane branch.",    "Follow the right side road."                             ],

        8:  [    "Pass the crosswalk.",              "Cross the crosswalk.",            "Traverse the crosswalk.",
                 "Drive through the crosswalk.",     "Move past the crosswalk.",        "Pass the crossing area.",
                 "Cross the crossing area.",         "Traverse the crossing area.",     "Drive through the crossing area.",
                 "Move past the crossing area."                                                                ],

        9:  [    "Pass the railroad.",               "Cross the railroad.",             "Traverse the railroad.",
                 "Drive through the railroad.",      "Move past the railroad.",         "Pass the railway.",
                 "Cross the railway.",               "Traverse the railway.",           "Drive through the railway.",
                 "Move past the railway."                                                                      ],

        10: [    "Merge.",                           "Merge traffic.",                  "Merge into traffic.",
                 "Merge into the traffic.",          "Join the traffic.",               "Merge into the traffic flow.",
                 "Join the traffic flow.",           "Merge into the traffic stream.",  "Join the traffic stream.",
                 "Merge into the lane.",                                                                       ],

        11: [    "Make a U-turn.",                   "Make a 180-degree turn.",         "Turn 180 degree.",
                 "Turn around.",                     "Drive in a U-turn."                                      ],

        12: [    "Stop.",                            "Halt.",                           "Decelerate.",
                 "Slow down.",                       "Brake."                                                  ],

        13: [    "Deviate.",                         "Deviate from the path.",          "Deviate from the lane.",
                 "Change the direction.",            "Shift the direction."                                    ]

}
```

Table 8. **Paraphrasing dictionary for command generation**. Each index corresponds to one of the 13 actions inferred by the classifier.

### C.1.3 Analyses Methods

In this section, we elaborate on the means of data analysis for OpenDV-YouTube. The analysis results are reported in the main paper and Appendix C.1.4.

**Geographic Diversity Analysis.** We take GPT-3.5-turbo [78] to infer the geographic information of each video from its title. We also apply handmade rules to post-process the results from GPT-3.5-turbo to deal with multiple aliases of one city or one country. The prompts are shown in Tab. 9 where `title` denotes the video title to be inferred. For simplicity, we assume that all clips of a video are taken in the same place. For videos with multiple inferred locations, we assume that all clips included in that video are uniformly distributed in these locations. For a video composed of $M$ clips with $N$ inferred locations, we assume there are $\frac{M}{N}$ clips taken in each site.

```
Messages = [

{ "role":  "system", "content":  f""" You are a helpful assistant, who is a geography expert and is also good at recog-
nizing different languages. """ },

{ "role":  "user", "content":  f""" Try to infer in which city or state a video is taken from its title. Please answer the
city name, the state name, and the country name in English respectively and briefly, in the following form: \

"Country: {{the name of the country}}
State: {{the name of the state or the province}}
```

City: {{the name of the city}}". \

If something cannot be inferred, fill the corresponding blank with "N/A". If there is more than one city in the video, first check if all the answers are valid, i.e. the name of cities, instead of the names of districts or towns. If there are multiple cities after checking the validity, use "," to separate different cities. \

You should also try to infer the state or province where the cities belong and fill the answer into the blank of "State". Note that you must infer the country where the video might be taken. Moreover, please discard meaningless words like "city", "country", "province" or "state" when filling in the blanks. \

The title of the video is as follows: {`title`}"""}]

Table 9. **Prompt for geographic inference of videos**.

**Scenario Diversity Analysis.**  For scene analysis, we visualize the frequency of different scenes in frame descriptions generated in Appendix C.1.2. For analyses on weather and time period, we observe that some language hints such as "foggy" and "night" are often present in videos' titles, thus we prompt GPT-3.5-turbo [78] to infer the weather and photographed period of the video from its title. The prompt is shown in Tab. 10 where `title` denotes the video title to be inferred.

`Messages` = [

{ `"role"`:  `"system"`, `"content"`:  f" You are a helpful assistant, who has a good command of multiple languages. "},

{ `"role"`:  `"user"`, `"content"`:  f""" Try to infer in which weather and period a video is taken from its title. Please answer the weather and period in English respectively and briefly, in the following form: \

"Weather: {{the weather}}
Period: {{the period}}". \

If something cannot be inferred, fill the corresponding blank with "N/A". The weather must be one of the following: "sunny", "rainy", "foggy", "snowy", "cloudy", "storm". The period must be one of the following: "daytime", "dusk", "dawn", "nighttime". \

The title of the video is as follows: {`title`}"""}]

Table 10. **Prompt for weather and time period inference of videos**.

### C.1.4 Diversity Highlights

**Geographic Distribution.**  As indicated by the human-refined GPT inference results, YouTube videos are taken from over 244 cities in more than 40 countries, covering considerably more areas than any existing public driving datasets, as shown in Tab. 1 and Fig. 2 in the main paper. Note that the result is still *underestimated* since the geographic information may not be included in the title for some videos and cannot be inferred. Taking the two most popular areas as an example, OpenDV-YouTube contains 36.4M clips in the US, covering 40 out of 50 states, and 12.9M clips in China, covering 26 out of 34 provinces. Moreover, to test the zero-shot performance of a model in its unseen locations, our YouTube-Val subset contains videos from 3 countries that are not included in YouTube-Train, *i.e.*, `Bosna i Hercegovina`, `Denmark,` and `Hungary`. There are also videos from 1 state of the US unseen in YouTube-Train, *i.e.*, `Maine`.

**Camera Settings.**  Considering that the online videos are sourced from different YouTubers around the globe, our dataset enjoys high diversity in photography equipment, leading to plentiful color settings, camera intrinsic parameters, and camera poses. For instance, a front-view video on a double-deck bus (see the second left picture in the last row of Fig. 9) is provided in our YouTube-Val subset while no similar cases are included in the YouTube-Train subset.

Figure 11. **Word cloud of frame descriptions for OpenDV-2K**. Only the top 500 most frequently mentioned objects, agents, or scenarios are included in the word cloud.

| Driver Action | Forward | Stop | Left Turn | Right Turn | U-Turn | Lane Change | Intersection Passing |
|---|---|---|---|---|---|---|---|
| Estimated Proportion | 81.39% | 8.85% | 1.89% | 1.81% | 0.27% | 0.30% | 5.49% |

Table 11. **Driver action distribution of OpenDV-YouTube**.

| Period | Daytime | Dawn | Dusk | Nighttime |
|---|---|---|---|---|
| Estimated Frame Count | 54M | 425K | 2M | 4M |

Table 12. **Time period distribution of OpenDV-YouTube**.

| Weather | Normal | Rainy | Cloudy | Foggy | Snowy | Storm |
|---|---|---|---|---|---|---|
| Estimated Frame Count | 58M | 690K | 503K | 284K | 503K | 117K |

Table 13. **Weather distribution of OpenDV-YouTube**.

**Scenarios.** We claim that there is sufficient data of diverse driver actions, weather conditions, photographed periods, and scenes in our OpenDV-YouTube. Results are shown in Tab. 11, Tab. 13, Tab. 12, and Fig. 11, respectively. Note that according to the analysis process in Appendix C.1.3, the diversity of scenarios in our dataset is *estimated* values since not all videos provide weather and filming periods in their titles.

**Corner Cases.** YouTube videos also contain corner cases and safety-critical cases. Several special cases from OpenDV-YouTube are given in Fig. 9, *e.g.*, dark tunnels with limited lighting (the leftmost and the rightmost in the 2nd row from bottom), intersections crowded with numerous pedestrians during nighttime (the 2nd right in the 3rd row from bottom), beaches at sunset (the 2nd left in the 5th row from bottom), rooftop (the leftmost in the last row), and videos captured with raindrops on the camera lens (the rightmost in the 5th row from top).

## C.2. Merged Public Datasets

Though the annotations in the OpenDV-YouTube are on a large scale, annotations are subject to limited patterns. *Contexts* from BLIP-2 follow certain syntax while *commands* are generated by the paraphrase dictionary. To provide more diverse expressions of contexts and commands, we merge annotations and sensor data from existing public datasets after converting their labels into complete sentences with correct grammar and format.

### C.2.1 Contexts Generation

**nuScenes & nuPlan.** Contexts are directly inherited from the scenario description of its belonging scenario in nuScenes [12] or nuPlan [13].

**ONCE.** In the metadata of ONCE [72], weather condition and filming time period are provided. These annotations are directly inherited in OpenDV-ONCE as contexts.

**Honda-HAD.** Diverse contexts are generated by refining and paraphrasing driving events provided by Honda-HAD [51]. The prompt for refinement and paraphrasing is as follows.

```
Messages = [

{ "role": "system", "content": f""" You are a helpful assistant. """ },

{ "role": "user", "content": f"""Generate {NUM_GEN} descriptions with exactly the same meanings as the following
reference sentence, REF: {current_caption}. \
```

Please write these sentences concisely in diverse ways, and try to use common and simple words if possible. \

Each generated sentence denotes a short description of noteworthy elements (e.g. pedestrians, traffic lights, cars) in this driving scenario. \

> There might be some typos, grammar errors, or unnatural expressions in the REF sentence, and you might need to correct these issues in the generated sentences. Each generated sentence should be correct in grammar and spelling, easy to understand, in natural and smooth expression. All sentences have the same meaning with the reference sentence REF, and the only difference is the wording. \
>
> Your complete response is only a python list including {NUM_GEN} strings (No other text needed), each one is an example sentence with an identifier `'\n'` in the end. """}]

Table 14. **Prompt for paraphrasing contexts in Honda-HAD dataset**.

### C.2.2 Commands Annotation

**nuScenes & nuPlan.** Since vehicle trajectory is given in nuScenes [12] and nuPlan [13], ego-vehicle commands can be easily calculated from trajectories by mathematical methods. After commands are generated, we can refer to Tab. 8 to provide diverse expressions of driver commands.

**Talk2Car.** Talk2Car [24] provides texts of possible human intentions for each scene in nuScenes. These annotations are inherited after they are refined by GPT-3.5-turbo to be grammatically correct and in appropriate formats. The prompt used for refinement is as follows.

> ```
> Messages = [
> ```
> { "role": "system", "content": f""" You are a helpful assistant. """ },
>
> { "role": "user", "content": f""" Please correct the capitalization and punctuation issues in this sentence: "{current_caption}". The original characters and words should be exactly the same without any changes. Do not add quotation marks. """}]

Table 15. **Prompt for refining texts in Talk2Car dataset**.

**ONCE.** Behaviours of the ego-vehicle can be obtained from the change in camera pose provided in ONCE [72]. They are further converted to natural language using Tab. 8.

**Honda-HAD.** Since driver behaviours are not directly provided in Honda-HAD [51], we implement the video classifier trained in Appendix C.1.2 and refer to Tab. 8 to generate ego-vehicle behaviours. Moreover, Honda-HAD does provide sufficient driving advice for each scene. We use GPT-3.5-turbo to refine and paraphrase these annotations so that diverse expressions are contained in our OpenDV-2K. Prompts for GPT-3.5-turbo are as follows.

> ```
> Messages = [
> ```
> { "role": "system", "content": f""" You are a helpful assistant. """ },
>
> { "role": "user", "content": f"""Generate {NUM_GEN} driving commands with exactly the same meanings as the following sentence: {current_caption}. \
>
> Please write these sentences concisely in diverse ways, and try to use common and simple words if possible. Remember all sentences have the same meaning, which is an instruction or intention for the planning of the ego vehicle. \
>
> Your complete response is only a python list including {NUM_GEN} strings (No other text needed), each one is an example sentence with an identifier `'\n'` in the end. """}]

Table 16. **Prompt for paraphrasing command annotations from driving advice in Honda-HAD dataset**.

**Honda-HDD-Action.** Honda-HDD-Action [84] contains 104 hours of videos with corresponding labels of driving commands. Since some clips begin with transitions from a completely green frame, we remove the first 30 frames from all

clips. Moreover, driving events in videos with a duration too long might be inconsistent with human-annotated behaviors. Therefore, we have to discard all video clips longer than 20 seconds. Meanwhile, since in the training stage, our model takes videos no shorter than 4 seconds as input, we also remove all videos shorter than 4 seconds. Only 32 hours of videos are left after this cleaning process. For the remaining clips, we directly use the labels as driver commands and use Tab. 8 to generate command texts.

**Honda-HDD-Cause.** There are 12 hours of videos in Honda-HDD-Cause [84], as well as corresponding human-annotated driving behaviors and human explanations. Similar to Honda-HDD-Action, we apply the same cleaning process on Honda-HDD-Cause, with about 1 hour of cleaned driving videos preserved. To align with OpenDV-YouTube, we convert these videos into frame sets by sampling the sensor videos at 10Hz. For command annotations, causal explanations in the form of phrases in the original dataset are inherited after refining and paraphrasing by GPT-3.5-turbo. The prompts used are as follows.

```
elements = "sign, congestion, traffic light, pedestrian, parked car"

Messages = [

{ "role": "system", "content": f""" You are a helpful AI driving assistant, who gives commands to the ego vehicle in natural language for safe driving. \

You are provided with one of the following elements of the driving scenario, namely, {elements}. Based on the given element, produce a driving command indicating either 'stop' or 'deviate' to the ego vehicle. Specifically, sign, congestion, traffic light, crossing vehicle, and pedestrian simulates a 'stop' command, and only the parked car leads to a 'deviate' command. \

You should write {NUM_GEN} fluent, concise, and diverse sentences for each command, using common and simple words. Half of these sentences are descriptions of the action of the ego vehicle (or driver), and the other half should be imperative sentences. All sentences should have the same meaning, and the only difference is the wording. """ },

{ "role": "user", "content": current_caption }]
```

Table 17. **Prompt for refining driving commands in Honda-HDD-Cause dataset**.

# D. Implementation Details of GenAD

## D.1. Model Design

### D.1.1 GenAD

GenAD is built upon 2.7B SDXL [79], which is a large-scale text-to-image generation model. We first fine-tune it in the first stage to transfer its domain knowledge to driving view synthesis. After that, we freeze the original blocks in the denoising UNet, and interleave them with our proposed temporal reasoning blocks, in total 2.5B, to allow for modeling on video sequences in video prediction pre-training. Following the original SDXL, the language conditions are encoded by two frozen CLIP variants with 817M parameters, namely, CLIP ViT-L [82] and OpenCLIP ViT-bigG [20], and the projection between the pixel space and latent space is performed by a pre-trained autoencoder with 83.7M parameters. As a result, GenAD has 5.9B parameters in total. The computational complexity is 5.27 TFLOPs.

### D.1.2 Extension on Action-condition Prediction

Besides the text conditions and past-frame conditions, we introduce the future trajectory of the ego vehicle as an additional condition signal to guide the denoising and therefore control the future imagination. We implement it by transforming the low-dimensional future waypoints into high-dimensional continuous embeddings [100], then projecting it with a zero-initialized linear layer into the same dimension with the text conditions $\mathbf{c}$. With zero initialization, the knowledge of future trajectory could be gradually injected into the model through the conditional cross-attention layer alongside $\mathbf{c}$, avoiding disturbing the learned prior on other conditions in the first place. It further controls the future simulation to be consistent with the ego intentions. Here the conditional future trajectory includes 6 waypoints at 2Hz.

### D.1.3 Extension on Planning

Since GenAD is capable of predicting reasonable futures given past observations, it encodes past frames in a meaningful way to guide the denoising of future frames. Therefore, we take the pre-trained GenAD as a strong feature extractor to obtain the spatiotemporal representations from past frames for downstream policy learning. We only utilize the encoder part of GenAD's denoising UNet to extract intermediate semantic features rather than acquiring the noise from the decoder part. Specifically, given the past two frames and a high-level command generated in the same way as in [45, 47], the frozen GenAD encoder extracts the spatiotemporal features, which are passed to a randomly initialized multi-layer perceptron (MLP) to project them into the future trajectory of ego vehicle. The MLP is composed of 6 linear layers and 5 ReLU activations, containing only 0.8M parameters in total. The first two linear layers downsample the features channel-wise, then the features of two frames are concatenated in channel dimension and further downsampled by the third linear layer. After that, the features are average-pooled in spatial dimensions, and the resulting vector is projected to the future trajectory, which is composed of 6 waypoints at 2Hz (3s), by the last two linear layers of the MLP.

## D.2. Training Details

GenAD is trained in two phases, *i.e.*, image domain transfer and video prediction pre-training. In the first stage, we fine-tune the pre-trained SDXL on per-image denoising with 2.7B trainable parameters of its denoising UNet. It is trained on 65.1M image-text pairs of OpenDV-2K. Each text condition is unified as "command, context". For some commands and contexts that are originally labeled for video sequences rather than static images, we simply associate them with all image frames included in that video sequence. We train the model for 300K iterations on 32 GPUs with a total batch size of 256 with AdamW [69]. We linearly warm up the learning rate for $10^4$ steps in the beginning then keep it constant at $1.25 \times 10^{-6}$. The default GPUs in most of our experiments are NVIDIA Tesla A100 devices unless otherwise specified.

In the second stage, we train the model on video-level denoising using video-text pairs lifting it to predict the future iteratively during inference. For compute efficiency, we freeze all blocks of the fine-tuned image model and only optimize our introduced temporal reasoning blocks, resulting in 2.5B trainable parameters in this stage. To maximize the data efficiency for constructing video clips, we take each frame of a 10Hz YouTube video as a starting frame to form a 4s training sequence at 2Hz, resulting in 65M video sequences for training. For each sequence with 8 frames at 2Hz, we randomly take the leading $m \in \{1, 2\}$ frames as conditional frames and the remaining $n \in \{7, 6\}$ frames to be corrupted for video denoising, with probabilities $p \in \{0.1, 0.9\}$, respectively. We do not add noise on conditional frames since there is no need to generate *past observations*. The text condition is structured in the same way as the first stage, and we acquire the context from the middle frame of the sequence. GenAD is trained on 64 GPUs for 112.5K iterations with a total batch size of 64. The learning rate is set as $1.25 \times 10^{-5}$ after $10^4$ warm-up steps.

In both stages, the input frames are resized to $256 \times 448$, and the text condition **c** is dropped at a probability of $p = 0.1$ to enable classifier-free guidance [40] in sampling. Both CLIP text encoders and the autoencoder are kept frozen throughout our experiments.

For extensions on action-conditioned prediction, we fine-tune the pre-trained GenAD as well as the linear projection layer for trajectory conditions on nuScenes. We conduct training on 16 GPUs for 100K steps with a total batch size of 16. Other training protocols such as the learning rate are the same with video prediction pre-training. For extensions on planning, we adapt a lightweight MLP to project the spatiotemporal features from frozen GenAD to future trajectory. We only optimize the MLP with 0.8M trainable parameters to adapt to planning. The MLP is trained for 12 epochs with a batch size of 16 and a learning rate of $5 \times 10^{-4}$, taking only 10 minutes to converge on a single NVIDIA Tesla V100 device.

## D.3. Sampling Details

Given two types of conditions including the past two frames and text, GenAD simulates 6 future frames accordingly via iteratively denoising its input latent, which starts from random Gaussian noises. The image resolution is $256 \times 448$ and the video sequence is at 2Hz. The sampling process is performed by Denoising Diffusion Implicit Models (DDIM) [96]. We use 100 sampling steps and set the scale of classifier-free guidance to 7.5. The sampling speed is 539.41 ms/step.

# E. Experimental Setup

## E.1. Data Preparation

We conduct extensive experiments on multiple datasets to evaluate the performance of our method. Specifically, the experiments of zero-shot transfer (Appendix F.2) are conducted on OpenDV-YouTube, Waymo [98], KITTI [30] and

Cityscapes [21]. Experiments of action-condition prediction (Appendix F.3) and motion planning (Main Sec. 4.3) are established on nuScenes [12]. The results of text-to-image generation (Appendix F.1) are shown in OpenDV-YouTube. As for failure case studies (Appendix F.4, Fig. 16), there are three cases in OpenDV-YouTube (a, b, d) and one case in Waymo (c). All results are reported in the validation set, which is completely unseen in the training of GenAD. All images and video frames are resized to $256 \times 448$ before being fed into GenAD. For tasks based on video prediction, we construct 2 frames in 1s at 2Hz as conditional frames. Each video sequence is paired with text conditions composed of command and context. For zero-shot datasets, the command and context are generated by the BLIP-2 model and video classifier respectively, following the preparation of training data. For nuScenes, we generate the command from logged trajectory following [45, 47] and map them to language using dictionary in Tab. 8, and we take the scenario descriptions as the context, which are officially provided in the dataset.

## E.2. Metrics

We use various metrics in multiple aspects for quantitative evaluation. These metrics include Fréchet Inception Distance (FID) [39] , Fréchet Video Distance (FVD) [101], CLIP-Similarity (CLIPSIM), Action Prediction Error, Average Displacement Error (ADE) and Final Displacement Error (FDE). For video prediction tasks, all predicted future frames are at 2Hz. We refer readers for discussions on metrics in Appendix A (Q4).

**FID:** It evaluates the generation quality of images, which are video frames in our experiments, by measuring the distribution distance of features between the predictions and original frames in the dataset. The features are extracted by a pre-trained Inception model. For quantitative comparison on nuScenes, FID is evaluated on 6019 generated frames and ground-truth frames. For experiments on YouTube, FID is calculated on 18000 frames from both generation and the dataset.

**FVD:** It measures the semantic similarity between real and synthesized videos with a pre-trained I3D action classification model [14] as the feature extractor. We evaluate 4369 video clips for the nuScenes comparison experiment, and 3000 video clips for YouTube.

**CLIPSIM:** We use the CLIP ViT-L/14 [82] to evaluate the consistency and coherence of the predicted video by computing the average similarity score of CLIP features between 6 generated frames and the first conditional frame. We take 3000 video sequences for evaluation.

**Action Prediction Error:** For experiments of action-condition prediction on nuScenes, it measures the consistency between the input trajectory $\mathbf{w}$ and predicted future frames of GenAD. We transform the future frames into trajectory $\hat{\mathbf{w}}$ using an inverse dynamics model (IDM), which is trained on nuScenes to project a video sequence into a trajectory following the design in [56]. This metric is then calculated as the mean L2 distance between all corresponding waypoints of $\mathbf{w}$ and $\hat{\mathbf{w}}$. Here both $\mathbf{w}$ and $\hat{\mathbf{w}}$ include 6 waypoints in 2 Hz, and $\mathbf{w}$ is generated from the logged trajectory in ego coordinate.

**ADE/FDE:** To evaluate the performance of planning on nuScenes, we calculate the ADE and FDE between the predicted trajectory and ground-truth trajectory in an open-loop setting. Here, ADE is the mean L2 distance between all waypoints of these two trajectories, and FDE is the L2 distance between the final waypoints of them.

## F. More Visualizations

### F.1. Image Generation in Driving Domain

After image domain transferring, the fine-tuned image model now focuses on synthesizing images in realistic driving views. Given text prompts in Tab. 18, the corresponding generated images are shown in Fig. 12 where the generated samples greatly reflect the abundant visual details in complex and driving scenes. The ability of high-quality driving-view generation laid the foundation for simulating a realistic futuristic driving world, which is learned through video prediction pre-training.

Figure 12. **Generated images by the fine-tuned image model.** Corresponding text prompts are listed in Tab. 18.

```
1.  Take a left turn.  A city at night with a lot of lights.
2.  Move steady.  A car driving down a highway with a view of the sky.
3.  Move steady.  A car driving through a tunnel.
4.  Drive steady.  A city street at night with cars and taxis.
5.  Keep the direction.  A city street with a crosswalk and tall buildings.
6.  Go straight.  A car driving down a mountain road.
7.  Maintain the direction.  A city street with parked cars.
8.  Turn to the left.  A car driving down a city street.
9.  Steer right.  A car driving on a mountain road.
10.  Make a right turn.  A car driving down a mountain road.
11.  Drive steady.  A car driving down a city street.
12.  Move steady.  A car driving down a road in a small village.
13.  Proceed.  A car driving on a highway with a sun in the sky.
14.  Drive steady.  A car driving down a snowy road.
15.  Take a left turn.  A car driving down a hill with houses on the side.
16.  Drive through the junction.  A red car is driving down a street in Boston.
17.  Brake.  A city street with cars and tall buildings.
18.  Proceed.  A car is driving down a hill with parked cars on the side.
19.  Decelerate.  A car is driving on a highway with cars behind it.
20.  Drive straight.
21.  Move forward.
22.  Move forward.  A car driving on a mountain road.
23.  Keep the direction.  A red double decker bus driving down a city street.
24.  Stop.  A car driving on a busy street.
25.  Drive straight.  A tram on a street at night.
26.  Drive forward.  A city street with a crosswalk.
27.  Proceed.  A green light on a street with cars and pedestrians.
28.  Move steady.  A view of a highway with a city in the background.
29.  Maintain the direction.  A view of a city street with buildings and mountains in the background.
30.  Drive straight.  A city street with a lot of cars and buildings.
31.  Steer right.  A car driving down a cobblestone street in a city.
32.  Drive forward.  A car driving on a dark road at night.
33.  Move forward.  A car driving down a busy street at night.
34.  Proceed.  A car driving through a tunnel.
35.  Drive straight.  A white van driving down a city street.
36.  Maintain the direction.  A city street at night with a ferris wheel.
```

Table 18. **Prompts for image generation in Fig. 12**, in the sequential order (from left to right and top to bottom).

## F.2. Zero-shot Transfer

With a strong capability on video prediction, the pre-trained GenAD can generalize to multiple unseen datasets in a zero-shot manner. In Fig. 13, we showcase multiple zero-shot video prediction results on OpenDV-YouTube. In Fig. 14, we illustrate the superiority of our method by comparing it to the previous state-of-the-arts on 4 datasets, including OpenDV-YouTube, Waymo [98], Cityscapes [21] and KITTI [30].

## F.3. Action-conditioned Prediction

By introducing an additional trajectory condition, the fine-tuned GenAD-act can be controlled to simulate different futures according to the input trajectory. We show four groups of action-conditioned prediction in Fig. 15. Both the input trajectory conditions (shown in the left bird's-eye view map) and imagined future frames are in 3s at 2Hz.

## F.4. Failure Cases

We showcase four failure cases generated by our model in Fig. 16. The model is sometimes disturbed by misleading contexts and is not strong enough to produce high-quality human details, as discussed in the Appendix A Q6. In some cases, the motion is not smooth enough. Meanwhile, the model fails to keep up with out-of-distribution camera height for 3s, even though succeeds in the first 2 seconds. These cases are worth future explorations.

Figure 13. **Zero-shot video prediction on OpenDV-YouTube** (the YouTube-Val subset from different YouTubers with strict geofence). The corresponding text conditions from top to bottom are as follows. 1. "Move steady. A car driving down a highway with cars behind it.", 2. "Turn to the left. A car is driving on a roof.", 3. "Maintain the direction. A taxi driving on a city street at night.", 4. "Drive forward. A car driving down a city street.", 5. "Proceed. A car driving on a bridge at night.", 6. "Steer left. A car driving on a road with trees and a blue sky.", 7. "Slow down. A street in a city with buildings and cars.", 8. "Go straight. A blue car driving down a city street.", 9. "Decelerate. A car driving on a city street.", 10. "Keep the direction. A view of a city street from the driver's seat.", 11. "Brake. A car driving down a street with trees and buses.", 12. "Proceed. A car driving on a city street at night.", 13. "Move forward. A car driving down a road near a river.", 14. "Drive straight. A van is driving down a highway with tall buildings in the background.".

Figure 14. **Zero-shot video prediction on public datasets compared with state-of-the-art video generation/prediction models**. Videos generated by I2VGen-XL are inconsistent with the condition frame. VideoCrafter1 appears to generate static scenarios. DMVFN suffers from huge image distortions. Meanwhile, all the other 3 models fail to generate videos when the ego vehicle should turn to the left and follow the lane (see the rightmost case in the last row). Our model manages to succeed in predictive video generation with great consistency with the conditional frames. We only show the first, third, and fifth frames from 6 predicted frames of our model due to space limits.

Figure 15. **Action-conditioned prediction on nuScenes**. We show four groups of video predictions for comparison, where each group is conditioned on the same two starting frames and different trajectories. In each group, the results in the first and second row are conditioned on the blue and green trajectories shown in the leftmost bird's-eye view, respectively.



Keep the direction. A city at night with a ferris wheel.

(a) Negative context effects

Move forward. A car is parked at a gas station.

(b) Sudden speed-up

Move steady. A woman is walking in a parking lot.

(c) Losses of human details

Pass the intersection. A yellow bus is driving down a street in a city.

(d) Failure in long-term maintainence of O.O.D. camera settings

Figure 16. **Examples of failure cases**. Examples (a, b, d) are from OpenDV-YouTube, and example (c) is from Waymo. We notice that sometimes contexts exert negative impacts on generated videos since the model tends to sacrifice temporal consistency to explicitly generate the object in the context under some circumstances (see example (a)). In examples (b) and (c), the model faces challenges in generating smooth motion and human details, respectively. In example (d), the model succeeds in holding on to the out-of-distribution camera setting, *i.e.*, on a double-deck bus, for the first 4 frames. But the camera height gradually falls down as normal in the last 2 frames.