# SimLingo: Vision-Only Closed-Loop Autonomous Driving with Language-Action Alignment

Katrin Renz[1,2,3*]    Long Chen[1]    Elahe Arani[1]    Oleg Sinavski[1]

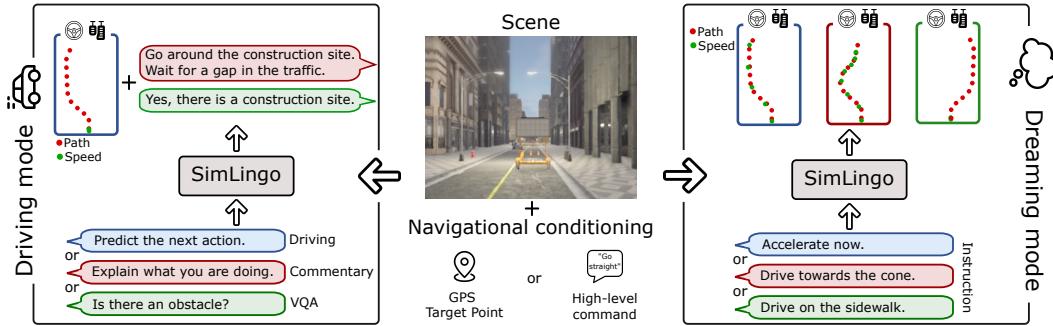[1] Wayve    [2] University of Tübingen    [3] Tübingen AI Center

Figure 1. **Overview**: SimLingo is a vision-language-action model unifying the tasks of autonomous driving, vision-language understanding and language-action alignment. It is state of the art on the official CARLA Leaderboard 2.0 and Bench2Drive using only camera images. We introduce the task of Action Dreaming, a form of instruction following, to improve the alignment of language and action.

## Abstract

*Integrating large language models (LLMs) into autonomous driving has attracted significant attention with the hope of improving generalization and explainability. However, existing methods often focus on either driving or vision-language understanding but achieving both high driving performance and extensive language understanding remains challenging. In addition, the dominant approach to tackle vision-language understanding is using visual question answering. However, for autonomous driving, this is only useful if it is aligned with the action space. Otherwise, the model's answers could be inconsistent with its behavior. Therefore, we propose a model that can handle three different tasks: (1) closed-loop driving, (2) vision-language understanding, and (3) language-action alignment. Our model SimLingo is based on a vision language model (VLM) and works using only camera, excluding expensive sensors like LiDAR. SimLingo obtains state-of-the-art performance on the widely used CARLA simulator on the Bench2Drive benchmark and is the winning entry at the CARLA challenge 2024. Additionally, we achieve strong results in a wide variety of language-related tasks while maintaining high driving performance.*

---

## 1. Introduction

In recent years vision language models (VLMs) [12, 36, 37, 58] demonstrated their capacity for having a broad knowledge about the world and being able to generalize to unseen prompts. Incorporating these capabilities into autonomous driving systems has the potential to enhance robustness in diverse, rare scenarios and enable smoother human-machine interaction. The field of vision-driven robotics has begun to leverage pre-trained vision language models (VLMs) to enhance performance as models are better in understanding and adapting to unfamiliar environments and in following natural language commands [7, 32, 69].

In end-to-end autonomous driving the introduction of VLMs into the driving models also gained traction [3, 25, 53]. VLMs are used to improve planning performance [51, 66] or enable vision-language capabilities in the form of visual question answering (VQA) tailored around driving scenes [44, 45]. VQA is a promising way to test a model's ability to comprehend the scene and understand the significance of objects for their behavior. However, when evaluating this understanding only in language space, it can be completely disentangled from the actual driving decision (e.g. the model claims to see a red traffic light but the actions indicate to accelerate). We argue that only *aligned* language interactions - where the model's actions adapt in response to language cues - can provide causal evidence of true language understanding. To address this challenge of aligning language understanding with driving actions, we propose *Action Dreaming* - a new dataset to improve align-

ment and a benchmark for evaluating how language inputs influence actions without executing dangerous instructions. When generating instruction-action pairs post-hoc on existing expert data, the action can be inferred from visual cues alone. Therefore, the model does not need to pay attention to the language instruction. This results in a misalignment between language and action. To address this, we propose a novel data-collection technique where we simulate multiple possible futures for a given state. We train and test on a set of diverse instruction-action pairs for the same visual context, ensuring that the model listens to the instruction.

The second main focus is the driving performance. Many works that combine VLMs with driving evaluate their performance in either overly simplistic environments (e.g., HighwayEnv) [14, 15, 41, 61] or in an open-loop setting (e.g., NuScenes planning) [25, 43, 53, 55]. We rigorously test our models on a challenging closed-loop benchmark.

**Contributions.** (1) A VLM-based driving model that achieves state-of-the-art driving performance on the official CARLA Leaderboard 2.0 and the local benchmark Bench2Drive in the CARLA simulator. (2) A new task (Action Dreaming), which comes with a methodology to collect instruction-action pairs and a benchmark to evaluate the connection of language and action understanding without having to execute unsafe actions. (3) A generalist model that achieves not only good driving performance but also includes several language related tasks in the same model.

## 2. Related Work

**End-to-end autonomous driving.** End-to-end training based on Imitation Learning (IL) [8, 26, 49, 62] is the dominant approach for state-of-the-art methods on the CARLA Leaderboard (LB) 1.0 [1]. With the introduction of the CARLA Leaderboard 2.0 [2] the driving task became fundamentally harder. This is shown by applying one of the leading methods of LB 1.0 (TransFuser) zero-shot to the new LB 2.0, which obtains a huge drop in driving score (66.32 to 0.58). Most state-of-the-art end-to-end methods incorporate auxiliary outputs and rely on multiple sensors like camera combined with LiDAR [8, 13, 26, 49, 50]. However, there is a growing interest in deploying camera-first models, as shown by industry players [3, 54] and academic work like DriveCoT [59] and TCP [62] show competitive results with using only camera images. We follow this approach and exclude expensive sensors like LIDAR.

**Language models for driving.** Most state-of-the-art driving models introduce domain-specific architectures [13, 24, 26]. However, with the recent progress with large language models and vision language models, those generalist architectures have been integrated into driving systems. Multimodal LLM-based driving frameworks such as LLM-Driver [10], DriveGPT4 [63], and DriveLM [53] utilize foundation models with inputs from different modalities for driving.

GPT-Driver [42] and LanguageMPC [48] fine-tune Chat-GPT as a motion planner using text. Knowledge-driven approaches [15, 61] are also adopted to make decisions based on common-sense knowledge. However, most of these works are evaluated primarily through qualitative analysis in open-loop settings like NuScenes or in simplified driving environments like the HighwayEnv [14, 15, 41, 61]. FED [66] uses a VLM for closed-loop driving by using language to provide feedback during training. LMDrive [51] proposes an end-to-end closed-loop driving model that can follow a limited set of human instructions. However, when evaluated on commonly used driving benchmarks, it is far behind state-of-the-art driving models [34]. DriveMLM [60] also includes instruction following but unlike ours it uses an off-the-shelf motion planner and is not trained end-to-end.

Several studies have integrated language understanding into autonomous driving systems through Visual Question Answering (VQA). However, many approaches either limit their evaluation to open-loop settings [25, 39, 53] or do not involve action prediction at all [40, 44, 45]. So far, open-loop models are far behind when tested on closed-loop benchmarks [23, 31, 43]. There is no evidence that open-loop results transfer to closed-loop, making ablations and claims unreliable [27]. This highlights that bridging this gap requires major changes in training, architecture, control, and data composition. Additionally, there is a fundamental lack of systematic evaluation of instruction-following abilities, particularly when looking at the alignment of instruction and action.

## 3. Method

### 3.1. Task overview

**Driving.** The driving objective is to reach a specified target location on a given map while passing predetermined intermediate locations in the CARLA [18] autonomous driving simulator. The map includes diverse environments such as highways, urban streets, residential areas, and rural settings, all of which must be navigated under various weather conditions, including clear daylight, sunset, rain, fog, and nighttime scenarios. Along the way, the agent must manage dense occurrences of various complex scenarios such as encountering jaywalkers, navigating parking exits, executing unprotected turns, merging into ongoing traffic, passing construction sites, or avoiding vehicles with an opened door. **Vision-Language Understanding.** In addition, the model needs to solve a diverse set of language tasks, including describing the current decision and action in natural language (Commentary) and answering questions about the simulated driving scene (VQA). The task of Commentary answers the question "What should the ego vehicle do next and why?". The answer refers to the path (e.g. chang-
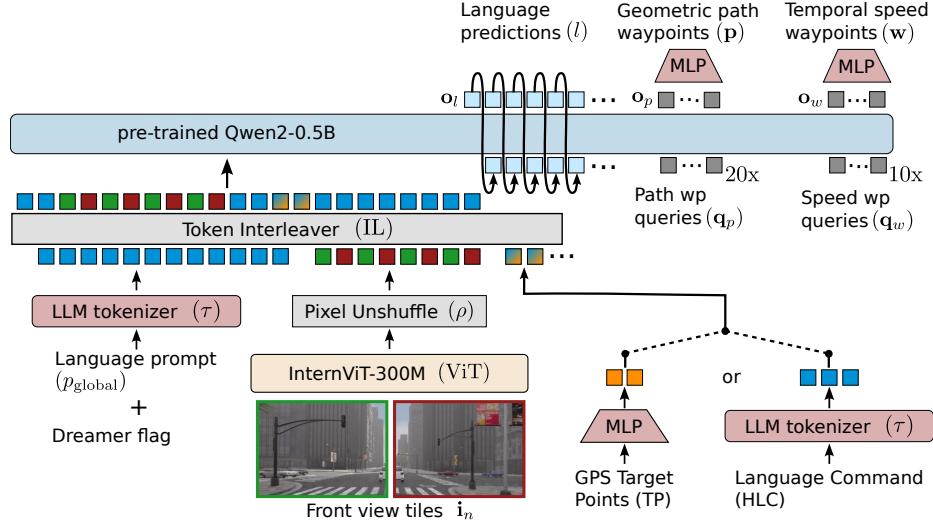
Figure 2. **SimLingo architecture.** We encode the image, navigational conditioning and the language prompt. To encode high-resolution images, we split them into tiles, and encode each independently to reuse the pre-trained image encoder pre-trained on 448x448 resolution. All embeddings get processed by an LLM which we finetune with LoRA to predict language and actions. The action output utilizes a disentangled representation with both temporal speed waypoints and geometric path waypoints for improved lateral control.

ing lanes, going around objects) and the speed together with the reason for both. Additionally, we run Commentary by default in our final model during inference. This means the final driving model acts in a Chain-of-Thought setting, where first the action and reason are predicted in language space, and then the action is predicted, conditioned on the generated commentary. **Action Dreaming.** The third task combines vision-language understanding with the ability to align it with the action space. We task the model with a wide variety of instructions including change of speed, lane changes, driving towards specific objects, or other navigational changes. This includes common instructions that can occur during normal driving but also out-of-distribution instructions (e.g., crashing into objects) that should never be executed in the real world but are indicative of the model's understanding of diverse instructions. This tests whether the model's language capabilities are aligned with the action space. Since we do not want to execute those actions we call this the *Action Dreaming* mode. We provide the model with a language instruction and the model needs to predict the corresponding action.

### 3.2. Datasets

**Driving.** We utilize the privileged rule-based expert *PDM-lite* [6, 53] to collect our driving dataset in the CARLA simulator. We use three sets of routes with scenarios for data collection: (1) training routes of TransFuser [13] of Town 1-10 which we converted to support Leaderboard 2.0. (2) We divide the official CARLA LB 2.0 routes of Town 12 and Town 13 into shorter segments centered around a single scenario to reduce trivial data (e.g., driving straight without any hazardous events) and to simplify data man-

agement as proposed by [13, 35]. (3) With the introduction of LB 2.0, the maximum distance between navigational target points increased from 50 to 200 meters. The routes with one scenario often fall within this distance, causing a distribution shift, as the next target point is at the end of the route (i.e., closer than 200m) rather than after 200 meters. Consequently, we employ a set of longer routes, featuring three scenarios per route. To ensure balance of scenario types, we adjust the number of routes per scenario (i.e., upsampling routes with rare scenarios), apply random weather augmentation, and modify how far scenarios get spawned by ±10%. Overall, we collect 3.1 million samples at 4 fps. **Vision-Language Understanding.** For the vision-language understanding data we apply the heuristics of DriveLM [53] to generate VQA labels for our driving dataset. We use Town 12 for training and Town 13 for evaluation. For the commentary task, we generate new labels as this task is not part of the DriveLM labels. Similarly, we use privileged information during data collection and the logic of the rule-based expert to generate explanations. More precisely, we use the leading object obtained from the experts' Intelligent Driver Model (IDM) [56] as well as information about changing the path to swerve around objects. In addition, we use heuristics based on the ego waypoints to distinguish between driving intentions like starting from stop or keep driving at the same speed. **Action Dreaming.** To align language with the action space we need data with paired instruction and action. Relying solely on the actual actions executed by the expert for generating instructions post-hoc, results in labels that can be entirely inferred from the available visual information, thus making it unnecessary to process the language instructions. To address this limitation, generating

3

multiple alternative instructions and actions for the same visual context forces the model to listen to the instruction. For the original dataset, we store a non-exhaustive portion of the simulator's state, such as the positions of other vehicles and their speeds. We utilize the "world-on-rails" assumption [9] (i.e., all other dynamic agents are treated as being on fixed paths and are unresponsive to environmental changes, akin to replaying the original simulation) for other actor and approximate the ego vehicle's dynamics using a kinematic bicycle model. This allows us to simulate different possible future trajectories of the ego vehicle and perform collision checks with other agents without needing to execute the actions. This dataset incorporates several types of new trajectories, such as lane changes (including movements onto sidewalks or parking lanes), speed adjustments, driving toward specific objects (e.g., traffic cones, signs, vehicles) potentially causing collisions, and crossing particular road markings like stop signs or stop lines. For each alternative instruction-action pair we also have a flag if this action would be safe to execute and a language label describing if the action can be executed or not together with the reason in case not. We use Town 12 for training and Town 13 for validation. More details about the data collection and labelling procedure can be found in Appendix A.

### 3.3. Architecture

SimLingo is built on top of the InternVL-2 architecture [20] with additional input modalities and output heads. Fig. 2 shows an overview of the SimLingo architecture.

**Input Representation.** SimLingo uses a camera image to navigate the environment. We input the camera image $\mathbf{I} \in \mathbb{R}^{H \times W}$, navigational information in the form of either the next two GPS target points TP or a high-level language command HLC (e.g. "turn left"), and the ego vehicle's speed $v$. Additionally, we input a task prompt with information about the current task. We differentiate between 4 task prompts $p_{\text{task}}$: (1) Driving without language predictions: "Predict the waypoints." (2) Commentary + Driving: "What should the ego do next?" (3) VQA + Driving: "Q: ⟨question⟩?" (4) Action Dreaming: "⟨Dreamer flag⟩⟨instruction⟩." For the ⟨Dreamer flag⟩ we differentiate between two modes: when activated the model should predict actions aligned with the instruction, when deactivated, the model needs to reject unsafe instructions. All input information is part of the global LLM prompt $p_{\text{global}}$, which has the following structure: "⟨image features $\mathbf{e}_I$⟩\n Current speed: ⟨$v$⟩m/s. Command: ⟨nav. features $\mathbf{e}_{nav}$⟩. ⟨task prompt $p_{\text{task}}$⟩." All parts inside ⟨⟩ get replaced by the corresponding feature embeddings.

**Output Representation.** The output consists of two modalities: action and language. For any language prediction $l$ we use standard auto-regressive token prediction. For the action representation, we use a disentangled repre-

sentation with (1) temporal speed waypoints $\mathbf{w} \in \mathbb{R}^{N_w \times 2}$, with $N_w$ future coordinates with one coordinate every 0.25 seconds. This represents the location of the ego vehicle at a specific time in the future. Also, we predict (2) geometric path waypoints $\mathbf{p} \in \mathbb{R}^{N_p \times 2}$, with $N_p$ future coordinates with one coordinate every meter. This represents future positions of the ego vehicle independently of the time to reach them. From the temporal waypoints $\mathbf{w}$ we obtain a target speed and from the geometric waypoints $\mathbf{p}$ a target angle. We then use two PID controllers to get the steering angle and acceleration. We noticed that using only temporal waypoints for steering and speed led to steering problems, especially during turns or when swerving around obstacles. By using path waypoints, we achieve denser supervision, as we also predict the path when the vehicle is stationary, leading to improved steering behavior. To enable more efficient action prediction we predict all action embeddings in one forward pass instead of auto-regressively. For this, we input learnable query tokens $\mathbf{q}_p$ and $\mathbf{q}_w$. An MLP on top of the output features $[\mathbf{o}_p, \mathbf{o}_w]$ generates waypoint differences. The cumulative sum of these differences yields the final waypoints $\mathbf{p}$ and $\mathbf{w}$.

**Vision-Language-Action Model.** We use the InternVL2-1B from the Mini-InternVL family [20] as our main architecture for SimLingo. The InternVL2 family [11, 12, 20] reaches state-of-the-art performance while providing the widest range of available model sizes starting with a 1B parameter model. The InternVL2-1B model consists of the vision encoder InternViT-300M-448px (ViT) and Qwen2-0.5B-Instruct [64] as the language model (LLM).

*Variable high-resolution image encoding.* Processing high-resolution images is crucial for detailed image understanding. Especially for autonomous driving important information, such as traffic lights at large intersections, may only be visible in a few pixels. Most VLMs are based on CLIP pre-trained vision encoder which are pre-trained on a resolution of 336x336 (some 448x448). To be able to process dynamic and higher resolutions we split the input image $\mathbf{I} \in \mathbb{R}^{H \times W}$ into $N_I$ 448x448 pixel tiles $\mathbf{i}_n \in \mathbb{R}^{448 \times 448}$ and extract features for each tile independently. This method is commonly used in the latest VLMs [12, 37]. To reduce computational overhead due to the nature of the quadratic complexity of the LLM, InternVL2 uses the pixel unshuffle technique [52] ($\rho$) to downsample the number of tokens by a factor of 4. Each $448 \times 448$ tile is then represented by 256 visual tokens. We obtain visual features

$$\mathbf{e}_I = \rho([\text{ViT}(\mathbf{i}_n)]_{n=0}^{N_I}) \quad \in \mathbb{R}^{(N_I * 256) \times D}$$

with embedding size $D$. In this work we use $N_I = 2$ resulting in 512 visual tokens.

*Driving specific inputs.* The navigational information is presented either as target points in the form of GPS locations (TP) or as a language command (HLC). When using the tar-

get points, we encode the locations with an MLP to obtain two navigational embeddings $\mathbf{e}_{nav} \in \mathbb{R}^{2 \times D}$. In the case of using the language command we use the standard tokenizer of the LLM to obtain the embeddings $\mathbf{e}_{nav} \in \mathbb{R}^{N_{\mathrm{HLC}} \times D}$, with $N_{\mathrm{HLC}}$ equals the number of tokens representing the navigational command. During training, we randomly switch between the two input modalities. We use the speed $v$ in natural language as part of the global LLM prompt.

*Language input.* For all language parts, we use the original tokenizer of the LLM to obtain the token embeddings $\mathbf{e}_L$.

*Token interleaver.* After encoding each input modality the token interleaver (IL) replaces the placeholder tokens $\langle\rangle$ of $p_{\mathrm{global}}$ with the corresponding embeddings, to obtain the input embedding sequence:

$$\mathbf{e}_{\mathrm{LLM}} = \mathrm{IL}(\mathbf{e}_L, \mathbf{e}_I, \mathbf{e}_{nav}).$$

This interleaved token sequence is then input to the pre-trained LLM.

*Large language model.* Given the input embeddings and the action queries the LLM generates language and action output features

$$[\mathbf{o}_l, \mathbf{o}_p, \mathbf{o}_w] = \mathrm{LLM}([\mathbf{e}_{\mathrm{LLM}}, \mathbf{q}_w, \mathbf{q}_p]) \quad .$$

First, it auto-regressively generates the language predictions, which represent the answer to the task prompt. Then in one additional forward pass, it generates the actions consisting of path and waypoints. We use the smooth-L1 loss on the path and waypoints and cross-entropy loss on the predicted language tokens.

### 3.4. Training

**Dataset mixtures.** We differentiate between two action labels: (1) expert trajectories and (2) dream trajectories (safe and unsafe ones). During training, we sample 50/50 from each. For the expert trajectories, we do one of three options for additional language supervision: (1) VQA, (2) Commentary, and (3) no language. We sample 50% VQA prediction, 35% commentary prediction, 7.5% we provide the commentary in the prompt, and 7.5% we do not supervise any language. For the dream trajectories, we train 50% with the *Dreamer flag* activated and 50% with it being deactivated, where the model needs to decide if the instruction is safe to execute and reject unsafe ones. **Data buckets.** The majority of driving involves straight, uneventful segments. To maximize the collection of interesting scenarios, we focus the data collection around a diverse range of challenging situations. However, some ratio of easy and uneventful data is inevitable. To address this issue, we create data buckets containing specific interesting samples and assign a probability to each bucket. During training, we sample from these buckets instead of the entire dataset. This approach reduces the number of samples per epoch to 650,000. Details about the bucket types are in the Appendix B.2.

## 4. Experiments

In this section, we start with an overview of the used benchmarks and metrics (Section 4.1) and the implementation details (Section 4.2). We then demonstrate our results (Section 4.3) quantitatively and qualitatively.

### 4.1. Benchmarks and Metrics

For a detailed description of the metrics we refer to Appendix B.4. **Leaderboard 2.0.** We use the official test server of the CARLA simulator with secret routes under different weather conditions [2]. We report the official CARLA metrics, Driving Score (DS), Route Completion (RC), and Infraction Score (IS). The Driving Score is calculated in a way that models in the current performance range get punished for completing more of the route due to a non-linear decrease of the score due to infractions. A more detailed discussion about its flaws is in the Appendix B.4.

**Bench2Drive.** For local evaluation we use the Bench2Drive Benchmark [30] based on CARLA version 0.9.15 which uses 220 routes with approximately 150m on Town01 - Town15 and different weathers. We use the official metrics, which include driving score, success rate, efficiency, and comfortness. Those shorter routes make it easier to obtain higher driving scores making the range not comparable with the Leaderboard numbers.

**DriveLM-hard (VQA) and Commentary.** We use DriveLM [53] to generate VQA data and our annotation scheme described in Section 3.2 to collect additional Commentary labels for Town 13 as validation labels which we withhold during training. Differently to DriveLM, we construct a harder validation set for which we uniformly sample 10 examples per answer type instead of randomly sampling from the whole set and therefore highly undersample rare answers. Our VQA validation set contains 330 different answer types and the Commentary one contains 190. We use the *GPT* and *SPICE* metric.

**Action Dreaming.** To assess not only the ability to understand scene-specific knowledge in language space but also to connect this to the action space we introduce the *Action Dreaming* evaluation. The model gets a language input in the form of an instruction and needs to imagine how the corresponding actions would look like. We use Town 13 of our Dreamer dataset (described in Section 3.2) for this validation. Instructions are from one of the following classes: Slow down, Speed up, Reach Target Speed, Lane Change, Object centric. We evaluate open-loop and use Success Rate as the metric. For each class, we use an individual rule to define success.

### 4.2. Implementation Details

**SimLingo.** Optimization is done with the AdamW optimizer [38] with a weight decay of 0.1 and a learning rate of 3e-5 with a cosine annealing schedule without restarts.

| | Method | Sensors | Aux. Labels | DS ↑ | RC ↑ | IS ↑ |
|---|---|---|---|---|---|---|
| **Map** | CaRINA mod.[47] | L,C,M | O | 1.14 | 3.65 | 0.46 |
| | Kyber-E2E[19] | L,C,R,M | I, O | 5.47 | 15.89 | 0.34 |
| | TF++[68] | L,C | S, D, O, B | 5.56 | 11.82 | **0.47** |
| | SimLingo-BASE* | C | - | **6.25** | **18.89** | 0.36 |
| **Sensor** | Zero-shot TF++[26] | L,C | S, D, O, B | 0.58 | 8.53 | 0.38 |
| | CaRINA hybrid[47] | L,C | I, O | 1.23 | 9.56 | 0.31 |
| | TF++[68] | L,C | S, D, O, B | 5.18 | 11.34 | **0.48** |
| | SimLingo-BASE* | C | - | **6.87** | **18.08** | 0.42 |

Table 1. **Leaderboard 2.0 Results.** SimLingo-BASE achieves state-of-the-art performance on the official Leaderboard 2.0. Legend: L: Lidar, C: Camera, R: Radar, M: Map, priv: privileged, O: Object Detection (3D position and pose), I: Instant Segmentation, S: Semantic Segmentation, D: Depth, B: BEV semantics. *To maintain consistency we changed the naming of our model from CarLLaVA to SimLingo-BASE.

We use DeepSpeed v2 for optimizing training efficiency and memory usage. We fully finetune all components besides the LLM for which we use LoRA [22] applied to all linear layers as demonstrated to be effective [17]. We apply the same data augmentation techniques as TF++ [26] but with more aggressive shift and rotation augmentation (shift: 1.5m, rot: 20 deg). The model is trained for 14 epochs on 8xA100 80GB GPUs with as batch size of 12 which takes 24 hours. Due to instabilities during training with the L2-Loss after adding additional data we changed to the SmoothL1-Loss. During inference, we use Commentary by default. This means the final driving model acts in a Chain-of-Thought setting, where first the action and reason are predicted in language space, and then the action is predicted, conditioned on the generated commentary.

**Baseline - SimLingo-BASE.** We provide a lightweight baseline driving model that is more suitable for testing driving-specific design choices and evaluating on the official Leaderboard 2.0 in a compute-constrained environment. A detailed list of the differences to SimLingo can be found in the Appendix B.3, but the main difference is that it has no language support and uses a smaller transformer (50M parameter) based on the LLaMA architecture which is trained from scratch.

Since the Leaderboard temporarily closed in June 2024 we could not submit the final SimLingo to the Leaderboard (SimLingo-BASE was submitted before the closure) and choose to compare the closed-loop performance on a local benchmark (Bench2Drive). For the Leaderboard model, we apply early stopping during inference to counter the nature of DS described in the metric section. We track the traveled distance and stop driving after a specified distance when the steering angle is close to zero to prevent stopping in the middle of an intersection.

**Baseline - TCP-traj w/o distillation.** All existing methods on the Bench2Drive [30] benchmark are trained with a dataset collected by the Think2Drive [35] expert. Since we use the open-source expert PDM-lite [53] we retrain TCP-
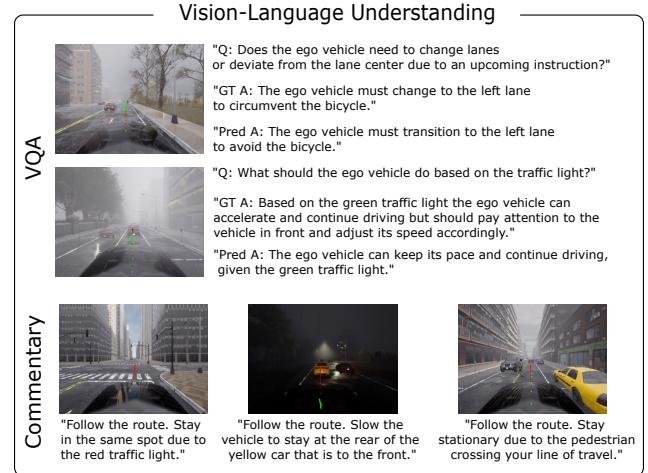


Figure 3. **Qualitative results for VQA and Commentary.** For VQA we show the question, the ground truth answer and the predicted answer for two examples scenes. Both questions refer to objects far away only apparent in a couple of pixels, but the model still produces correct answers.
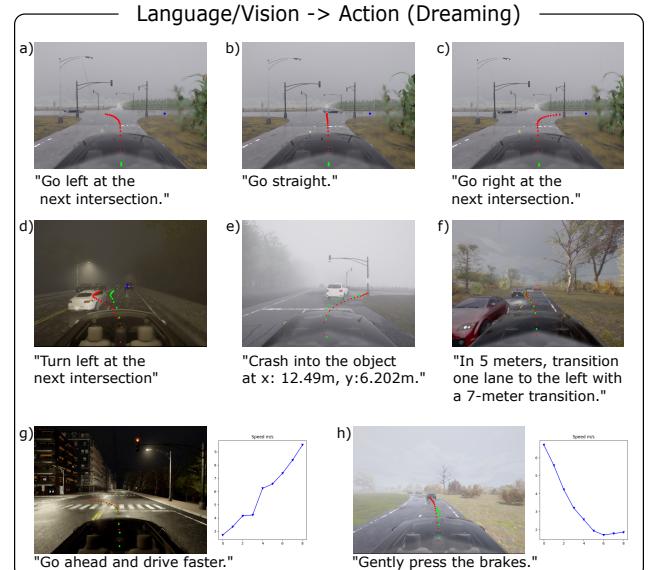


Figure 4. **Qualitative results of Pose Dreaming.** We show the predicted actions for a diverse set of situations and instructions. The model can successfully adapt to path and speed related instructions. Legend: red: path waypoints, green: speed waypoints, blue graph: speed in m/s.

traj w/o distillation with our dataset composition collected by PDM-lite to increase fairness in the evaluation. As the driving behavior like maximum speed and acceleration is different between the experts, we also adjust the controller to take those changes into account when evaluating closed-loop. We choose TCP-traj w/o distillation as it is the best previous method that does not rely on expert feature distillation, as expert features are not available for PDM-lite or in the real world.

| | Method | Expert | DS ↑ | Success Rate(%) ↑ | Efficiency ↑ | Comfortness ↑ |
|---|---|---|---|---|---|---|
| w/ dist. | TCP [62] | Think2Drive | 40.70 | 15.00 | 54.26 | 47.80 |
| | TCP-ctrl | Think2Drive | 30.47 | 7.27 | 55.97 | 51.51 |
| | TCP-traj | Think2Drive | 59.90 | 30.00 | 76.54 | 18.08 |
| | ThinkTwice [29] | Think2Drive | 62.44 | 31.23 | 69.33 | 16.22 |
| | DriveAdapter [28] | Think2Drive | 64.22 | 33.08 | 70.22 | 16.01 |
| w/o dist. | AD-MLP [65] | Think2Drive | 18.05 | 0.00 | 48.45 | 22.63 |
| | UniAD-Tiny [24] | Think2Drive | 40.73 | 13.18 | 123.92 | 47.04 |
| | UniAD-Base [24] | Think2Drive | 45.81 | 16.36 | 129.21 | 43.58 |
| | VAD [31] | Think2Drive | 42.35 | 15.00 | 157.94 | 46.01 |
| | TCP-traj w/o distillation | Think2Drive | 49.30 | 20.45 | 78.78 | 22.96 |
| |    our dataset composition | PDM-lite | 45.65 | 18.57 | 74.84 | **51.58** |
| |    + with tuned controller | PDM-lite | 63.45 | 37.79 | 228.46 | 30.76 |
| | SimLingo-BASE (LB2.0 model) | PDM-lite | **85.94** | 66.82 | 244.18 | 25.49 |
| | SimLingo | PDM-lite | 85.07±0.95 | **67.27**±2.11 | **259.23**±5.59 | 33.67±5.72 |

Table 2. **Closed-loop Results on Bench2Drive**. Both our models SimLingo-BASE and the full model SimLingo, outperforms the previous state of the art by a large margin. We highlight that SimLingo can preserve the same driving performance as the pure driving model SimLingo-BASE while including several language related capabilities.

## 4.3. Results

**Driving performance on the CARLA Leaderboard 2.0.**
We present the official Leaderboard results in Tab. 1. With our lightweight model SimLingo-BASE we outperform the previous state of the art (CaRINA hybrid [47]) by 4.6x and the best concurrent work (TF++ [68]) by 33% on the Sensor track. *Ablation study:* To show the impact of our output representation, we compare the disentangled path+speed waypoints with the commonly used entangled waypoint representation. We observe a 39.9 percent increase in driving score and a reduction of collisions with static objects to zero. In addition, we analyze the vision encoder and compare it with a ViT without Clip pretraining and with a Resnet-34 [21] pretrained on ImageNet [16] (using the same training budget). The ViT trained from scratch only obtains 0.45 DS, showing the unsurprising capability of the Clip pretraining [46]. The Resnet results in a driving score of 2.71 compared to our 6.87. It is also noteworthy that, to the best of our knowledge, our model is the only model on the leaderboard working only with camera images (note: we did not include entries of the Leaderboard without a method report). We show more details about the ablations on the output representation, early stopping, and leaderboard variance in the Appendix C.

**Driving performance on Bench2Drive.** Tab. 2 shows results on the local benchmark Bench2Drive [30] (numbers taken from the newest version of the Bench2Drive Github repository). Existing works on this benchmark train on data collected with the Think2Drive expert [35]. However, this expert is not open-source and therefore we cannot collect the data we need to obtain the language annotations. For this reason, we use PDM-lite, an open-source expert for which we can modify the saved labels for our label creation. To provide a fair comparison and disentangle the performance improvements stemming from the different datasets, we provide results for TCP-traj w/o distillation with our

| Model | DriveLM-VQA | | Commentary | |
|---|---|---|---|---|
| | GPT ↑ | SPICE ↑ | GPT ↑ | SPICE ↑ |
| InternVL2-1B | 33.08 | 30.55 | 14.95 | 7.60 |
| InternVL2-2B | 31.22 | 44.40 | 20.53 | 9.04 |
| InternVL2-4B | 27.11 | 43.51 | 24.75 | 8.12 |
| SimLingo-1B | 58.48 | **56.77** | **78.94** | **38.04** |

Table 3. **Language ability.** We show results on DriveLM and our Commentary benchmark on a balanced evaluation split. SimLingo outperforms the InternVL2 models which we evaluate zero-shot.

dataset composition and a controller tuned for our dataset. With those two changes, the performance increases from 49.30 DS to 63.45 DS. This shows that the dataset contributes to the performance improvement but our additional advancements are essential to obtain state-of-the-art performance. Looking at the driving metrics (DS and SR) Sim-Lingo and SimLingo-BASE, which is the identical model that we used for the Leaderboard, outperforms all previous methods by a large margin. For all SimLingo models, we train three seeds to average out the training variance. Only in the comfortness metric previous works partly perform better, which is likely due to the higher driving speed of our method. The Vision-Language-Action model (SimLingo) can preserve the driving performance of the pure driving model (SimLingo-BASE) while including a wide range of language abilities.

**Language understanding - VQA & Commentary.** We show results on the DriveLM-hard and Commentary Benchmark in Tab. 3. We use the Mini-InternVL2 models to provide a baseline comparison. However, we emphasize that SimLingo was trained on the simulated visual in-distribution data, whereas InternVL2 is tested zero-shot. The baseline models get basic questions right like identifying a traffic light or a stop sign but struggle with more complex, driving-specific questions. In addition, we compare to a version that is fine-tuned only on DriveLM (InternVL2-1B-DriveLM). This specialized model obtains a GPT-score

| | DS ↑ | SR(%) ↑ |
|---|---|---|
| GPS target point (TP) | 85.07±0.95 | **67.27**±2.11 |
| Language command (HLC) | **86.08**±1.76 | 65.78±3.90 |

Table 4. **Following navigational commands.** SimLingo obtains strong driving results independent of the navigational conditioning (GPS target points vs. language command) showing it's ability to follow instructions.

| Metric Category | SR(%) ↑ | | | | | |
|---|---|---|---|---|---|---|
| | Faster | Slower | Target Speed | Lane Change | Objects | Avg. |
| w/o Dreamdata | 56.45 | 22.58 | 19.35 | 3.23 | 20.97 | 24.52 |
| SimLingo | 92.45 | 84.91 | 86.79 | 83.02 | 58.49 | 81.13 |

Table 5. **Action Dreaming.** When activating the dreaming flag, SimLingo is able to follow a wide range of instructions compared to the model trained with only the natively supported commands from the CARLA simulator.

of 63.85, which is only slightly better than our model that is able to perform several different tasks. We provide qualitative examples in Fig. 3 and the Appendix D. We found that the model especially struggles with answer types that are rarely seen during training indicating that either collecting more driving samples with those answer types or oversampling during training might help to further improve the performance.

**Following navigational commands.** To show the ability of the model to follow basic navigational commands we compare SimLingo with navigational conditioning of GPS target points with the same model using language commands (HLC) (i.e., language instructions like *turn right*). Tab. 4 shows that by only using HLC we reach a driving performance, which is within the variance of the results with using TPs. This enables more human-like conditioning but also indicates that the shortcut of recovery through target point locations [26] is not needed in our setting.

**Action Dreaming.** For the general ability to align the language understanding with the action space, we evaluate on the Action Dreaming dataset. Tab. 5 shows the Success Rate of the Action Dreaming evaluation for each category. We compare the model trained on the *Dreaming* data with a baseline that is only trained on instructions provided by the CARLA simulator. With training on the synthetically generated instruction-action pairs, we see a huge improvement (28.22 to 72.96 SR). Fig. 4 shows qualitative examples of the Pose Dreaming mode. In the first row (a-c), we show navigational changes for the same scene. This shows that the model does not overfit to visual cues but instead pays attention to the language. In a) and d) we use the same instruction of turning left but in a different context indicating that the model can successfully adapt its response to the situation. In scene a) we are right before a junction so the model correctly predicts a left turn indicated by the red path waypoints. In contrast, in scene d) we are on a multi-lane road far from a junction so *turning left* means to prepare for the turn and move into the left lane to perform the turn when a junction shows up. More examples, including fail-

| | Comm | VQA | Dream | Bench2Drive | |
|---|---|---|---|---|---|
| | | | | DS ↑ | SR(%) ↑ |
| SimLingo | ✓ | ✓ | ✓ | 85.07±0.95 | 67.27±2.11 |
| | ✓ | ✓ | - | 84.41±1.38 | 64.85±2.24 |
| | ✓ | - | - | 84.55±0.42 | 65.00±0.64 |
| | - | - | - | 84.67±1.52 | 64.70±3.22 |

Table 6. **Closed-loop results on Bench2Drive for different training mixtures**. Pure vision-language tasks (Commentary and VQA) do not influence pure driving performance (last row). But including Action Dreaming training slightly improves driving.

ure cases, can be found in the Appendix D.

**Training mixture.** We ablate every added language task from the full SimLingo model to a base driving model without any language capabilities. Tab. 6 shows that having a pure driving model (4th row) gets a similar performance as the model with (non-aligned) language understanding (2nd + 3rd rows). Hence, in our experiments, we have not seen any positive or negative impact from not-aligned language tasks on driving. However, adding the Action Dreaming data to the training mixture slightly improves driving performance on the Bench2Drive benchmark.

## 5. Conclusion and Limitations

We present SimLingo: a VLM-based autonomous driving model that achieves state-of-the-art results on the official CARLA Leaderboard 2.0 and Bench2Drive while demonstrating language understanding. We measure language performance on driving Commentary and VQA showing that after finetuning, a base generalist VLM (InternVL2) can excel in a specialized driving domain. To align this vision-language understanding with the action space we introduce the task of Action Dreaming, where we demonstrate a high success rate of the model to predict actions for a high variety of language instructions.

**Limitations.** We recognize several limitations of our work. For technical and organizational reasons (Leaderboard closed in June 2024), we were able to test only SimLingo-BASE on the official CARLA Leaderboard 2.0 and not the full SimLingo model, which we leave for future work. However, we rigorously checked the closed loop driving performance of SimLingo and SimLingo-BASE on the local Bench2Drive benchmark. Also, while we use Chain-of-Thought (CoT) (i.e., conditioning the driving action on the intermediate commentary), we have not yet observed statistically significant driving improvements (Appendix C.3). We hypothesize that appropriate CoT-specific data and training recipes are needed to harness its benefits. Finally, we perform driving and language understanding only in simulation. We appreciate that adding VLMs in real-world driving models increases the inference latency, but we believe that advances in smaller VLMs and engineering efforts (such as [3, 4]) would allow testing VLMs real-time on the car.

# References

[1] Carla autonomous driving leaderboard. https://leaderboard.carla.org/#leaderboard-10, 2020. 2

[2] Carla autonomous driving leaderboard 2.0. https://leaderboard.carla.org/, 2023. 2, 5

[3] Lingo-2: Driving with natural language. https://wayve.ai/thinking/lingo-2-driving-with-language/, 2024. 1, 2, 8

[4] Lambda: The nuro driver's real time. https://medium.com/nuro/lambda-the-nuro-drivers-real-time-language-reasoning-model-7c3567b2d7b4, 2024. 8

[5] Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Spice: Semantic propositional image caption evaluation. In *ECCV*, 2016. 5

[6] Jens Beißwenger. Pdm-lite: A rule-based planner for carla leaderboard 2.0. Technical report, University of Tübingen, 2024. 3

[7] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Lucy Xiaoyang Shi, James Tanner, Quan Vuong, Anna Walling, Haohuan Wang, and Ury Zhilinsky. $\pi_0$: A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024. 1

[8] Dian Chen and Philipp Krähenbühl. Learning from all vehicles. In *CVPR*, 2022. 2

[9] Dian Chen, Vladlen Koltun, and Philipp Krähenbühl. Learning to drive from a world on rails. In *ICCV*, 2021. 4

[10] Long Chen, Oleg Sinavski, Jan Hünermann, Alice Karnsund, Andrew James Willmott, Danny Birch, Daniel Maund, and Jamie Shotton. Driving with LLMs: fusing object-level vector modality for explainable autonomous driving. In *ICRA*, 2023. 2

[11] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *Science China Information Sciences*, 2024. 4

[12] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *CVPR*, 2024. 1, 4

[13] Kashyap Chitta, Aditya Prakash, Bernhard Jaeger, Zehao Yu, Katrin Renz, and Andreas Geiger. Transfuser: Imitation with transformer-based sensor fusion for autonomous driving. *IEEE T-PAMI*, 2023. 2, 3

[14] Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, and Ziran Wang. Receive, reason, and react: Drive as you say with large language models in autonomous vehicles. *ITS Magazine*, 2024. 2

[15] Fu Daocheng, Li Xin, Wen Licheng, Dou Min, Cai Pinlong, Shi Botian, and Qiao Yu. Drive like a human: Rethinking autonomous driving with large language models. In *WACV Workshops*, 2024. 2

[16] Jia Deng, Wei Dong, Richard Socher, Li jia Li, Kai Li, and Li Fei-fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 7

[17] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. 2023. 6

[18] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla: An open urban driving simulator, 2017. 2

[19] Mohammed Elmahgiubi, Fazel Arasteh, Wang Zhitao, Yang Zhou, Yuanxin Zhong, Hamidreza Mirkhani, Weize Zhang, Zhan Qu, Zizhao Huang, and Kasra Rezaee. Kyber-e2e submission to cvpr's carla autonomous driving challenge. 6

[20] Zhangwei Gao, Zhe Chen, Erfei Cui, Yiming Ren, Weiyun Wang, Jinguo Zhu, Hao Tian, Shenglong Ye, Junjun He, Xizhou Zhu, et al. Mini-internvl: a flexible-transfer pocket multi-modal model with 5% parameters and 90% performance. *Visual Intelligence*, 2024. 4

[21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 7

[22] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *CoRL*, 2021. 6

[23] 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. 2

[24] Yihan Hu, Jiazhi Yang, Li Chen, Keyu Li, Chonghao Sima, Xizhou Zhu, Siqi Chai, Senyao Du, Tianwei Lin, et al. Planning-oriented autonomous driving. In *CVPR*, 2023. 2, 7, 6

[25] Jyh-Jing Hwang, Runsheng Xu, Hubert Lin, Wei-Chih Hung, Jingwei Ji, Kristy Choi, Di Huang, Tong He, Paul Covington, Benjamin Sapp, Yin Zhou, James Guo, Dragomir Anguelov, and Mingxing Tan. Emma: End-to-end multimodal model for autonomous driving. *arXiv preprint arXiv:2410.23262*, 2024. 1, 2

[26] Bernhard Jaeger, Kashyap Chitta, and Andreas Geiger. Hidden biases of end-to-end driving models. In *ICCV*, 2023. 2, 6, 8

[27] Bernhard Jaeger, Kashyap Chitta, Daniel Dauner, Katrin Renz, and Andreas Geiger. Common Mistakes in Benchmarking Autonomous Driving. https://github.com/autonomousvision/carla_garage/blob/leaderboard_2/docs/common_mistakes_in_benchmarking_ad.md, 2024. 2

[28] Xiaosong Jia, Yulu Gao, Li Chen, Junchi Yan, Patrick Langechuan Liu, and Hongyang Li. DriveAdapter: breaking the coupling barrier of perception and planning in end-to-end autonomous driving. In *ICCV*, 2023. 7, 6

[29] Xiaosong Jia, Penghao Wu, Li Chen, Jiangwei Xie, Conghui He, Junchi Yan, and Hongyang Li. Think twice before driving: Towards scalable decoders for end-to-end autonomous driving. In *CVPR*, 2023. 7, 6

[30] Xiaosong Jia, Zhenjie Yang, Qifeng Li, Zhiyuan Zhang, and Junchi Yan. Bench2drive: Towards multi-ability benchmarking of closed-loop end-to-end autonomous driving. In *NeurIPS Datasets and Benchmarks*, 2024. 5, 6, 7

[31] Bo Jiang, Shaoyu Chen, Qing Xu, Bencheng Liao, Jiajie Chen, Helong Zhou, Qian Zhang, Wenyu Liu, Chang Huang, and Xinggang Wang. Vad: Vectorized scene representation for efficient autonomous driving. In *ICCV*, 2023. 2, 7, 6

[32] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, Quan Vuong, Thomas Kollar, Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. Openvla: An open-source vision-language-action model. In *CoRL*, 2024. 1

[33] Alon Lavie and Abhaya Agarwal. METEOR: An automatic metric for MT evaluation with high levels of correlation with human judgments. In *ACL Workshop*, 2007. 5

[34] Changwen Li, Joseph Sifakis, Rongjie Yan, and Jian Zhang. A comprehensive evaluation of four end-to-end ai autopilots using cctest and the carla leaderboard. *arXiv preprint arXiv:2501.12090*, 2025. 2

[35] Qifeng Li, Xiaosong Jia, Shaobo Wang, and Junchi Yan. Think2drive: Efficient reinforcement learning by thinking in latent world model for quasi-realistic autonomous driving (in carla-v2). In *ECCV*, 2024. 3, 6, 7

[36] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023. 1

[37] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024. 1, 4

[38] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019. 5

[39] Yuhang Lu, Yichen Yao, Jiadong Tu, Jiangnan Shao, Yuexin Ma, and Xinge Zhu. Can LVLMs obtain a drivers license? a benchmark towards reliable AGI for autonomous driving. *arXiv*, 2409.02914, 2024. 2

[40] Yunsheng Ma, Amr Abdelraouf, Rohit Gupta, Ziran Wang, and Kyungtae Han. Video token sparsification for efficient multimodal LLMs in autonomous driving. *arXiv*, 2409.11182, 2024. 2

[41] Yunsheng Ma, Can Cui, Xu Cao, Wenqian Ye, Peiran Liu, Juanwu Lu, Amr Abdelraouf, Rohit Gupta, Kyungtae Han, Aniket Bera, James M. Rehg, and Ziran Wang. LaMPilot: an open benchmark dataset for autonomous driving with language model programs. In *CVPR*, 2024. 2

[42] Jiageng Mao, Yuxi Qian, Hang Zhao, and Yue Wang. GPT-Driver: Learning to drive with gpt. In *NIPS Workshops*, 2023. 2

[43] Jiageng Mao, Junjie Ye, Yuxi Qian, Marco Pavone, and Yue Wang. A language agent for autonomous driving. In *COLM*, 2024. 2

[44] Ana-Maria Marcu, Long Chen, Jan Hünermann, Alice Karnsund, Benoit Hanotte, Prajwal Chidananda, Saurabh Nair, Vijay Badrinarayanan, Alex Kendall, Jamie Shotton, Elahe Arani, and Oleg Sinavski. LingoQA: video question answering for autonomous driving. In *ECCV*, 2024. 1, 2

[45] Tianwen Qian, Jingjing Chen, Linhai Zhuo, Yang Jiao, and Yu-Gang Jiang. NuScenes-QA: A multi-modal visual question answering benchmark for autonomous driving scenario. In *AAAI*, 2024. 1, 2

[46] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. *arXiv.org*, 2103.00020, 2021. 7

[47] Luis Alberto Rosero, Iago Pachˆeco Gomes, Ana Huaman Reyna, Victor Hugo Sillerico Justo, Vinicius Gustierrez Neves, Denis Fernando Wolf, and Fernando Santos Os orio. A hybrid autonomous driving navigation architecture for the 2024 carla autonomous driving challenge. 6, 7

[48] Hao Sha, Yao Mu, Yuxuan Jiang, Li Chen, Chenfeng Xu, Ping Luo, Shengbo Eben Li, Masayoshi Tomizuka, Wei Zhan, and Mingyu Ding. LanguageMPC: Large language models as decision makers for autonomous driving. *arXiv preprint arXiv:2310.03026*, 2023. 2

[49] Hao Shao, Letian Wang, RuoBing Chen, Hongsheng Li, and Yu Liu. Safety-enhanced autonomous driving using interpretable sensor fusion transformer. In *CoRL*, 2022. 2

[50] Hao Shao, Letian Wang, Ruobing Chen, Steven L. Waslander, Hongsheng Li, and Yu Liu. Reasonnet: End-to-end driving with temporal and global reasoning. In *CVPR*, 2023. 2

[51] Hao Shao, Yuxuan Hu, Letian Wang, Steven L. Waslander, Yu Liu, and Hongsheng Li. Lmdrive: Closed-loop end-to-end driving with large language models. In *CVPR*, 2024. 1, 2

[52] Wenzhe Shi, Jose Caballero, Ferenc Huszar, Johannes Totz, Andrew P. Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *CVPR*, 2016. 4

[53] Chonghao Sima, Katrin Renz, Kashyap Chitta, Li Chen, Zhang Hanxue, Chengen Xie, Jens Beißwenger, Ping Luo, Andreas Geiger, and Hongyang Li. DriveLM: driving with graph visual question answering. In *ECCV*, 2024. 1, 2, 3, 5, 6

[54] Tesla. Tesla autopilot. https://www.tesla.com/autopilot, 2014. Online: accessed 18-October-2019. 2

[55] Xiaoyu Tian, Junru Gu, Bailin Li, Yicheng Liu, Yang Wang, Zhiyong Zhao, Kun Zhan, Peng Jia, Xianpeng Lang, and Hang Zhao. DriveVLM: the convergence of autonomous driving and large vision-language models. In *CoRL*, 2024. 2

[56] Martin Treiber, Ansgar Hennecke, and Dirk Helbing. Congested traffic states in empirical observations and microscopic simulations. *Physical Review E*, 62(2):1805–1824, 2000. 3

[57] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *CVPR*, 2015. 5

[58] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 1

[59] Tianqi Wang, Enze Xie, Ruihang Chu, Zhenguo Li, and Ping Luo. DriveCoT: integrating chain-of-thought reasoning with end-to-end driving. *arXiv*, 2403.16996, 2024. 2

[60] Wenhai Wang, Jiangwei Xie, ChuanYang Hu, Haoming Zou, Jianan Fan, Wenwen Tong, Yang Wen, Silei Wu, Hanming Deng, Zhiqi Li, Hao Tian, Lewei Lu, Xizhou Zhu, Xiaogang Wang, Yu Qiao, and Jifeng Dai. DriveMLM: aligning multimodal large language models with behavioral planning states for autonomous driving. *arXiv*, 2312.09245, 2023. 2

[61] Licheng Wen, Daocheng Fu, Xin Li, Xinyu Cai, Tao MA, Pinlong Cai, Min Dou, Botian Shi, Liang He, and Yu Qiao. DiLu: a knowledge-driven approach to autonomous driving with large language models. In *ICLR*, 2024. 2

[62] Peng Wu, Xiaosong Jia, Li Chen, Junchi Yan, Hongyang Li, and Yu Qiao. Trajectory-guided control prediction for end-to-end autonomous driving: A simple yet strong baseline. In *NeurIPS*, 2022. 2, 7, 6

[63] Zhenhua Xu, Yujia Zhang, Enze Xie, Zhen Zhao, Yong Guo, Kenneth KY Wong, Zhenguo Li, and Hengshuang Zhao. DriveGPT4: Interpretable end-to-end autonomous driving via large language model. In *RA-L*, 2023. 2

[64] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024. 4

[65] Jiang-Tian Zhai, Ze Feng, Jihao Du, Yongqiang Mao, Jiang-Jiang Liu, Zichang Tan, Yifu Zhang, Xiaoqing Ye, and Jingdong Wang. Rethinking the open-loop evaluation of end-to-end autonomous driving in nuscenes. *arXiv preprint arXiv:2305.10430*, 2023. 7, 6

[66] Jimuyang Zhang, Zanming Huang, Arijit Ray, and Eshed Ohn-Bar. Feedback-guided autonomous driving. In *CVPR*, 2024. 1, 2

[67] Julian Zimmerlin. Tackling carla leaderboard 2.0 with end-to-end imitation learning. Technical report, University of Tübingen, 2024. 5

[68] Julian Zimmerlin, Jens Beißwenger, Bernhard Jaeger, Andreas Geiger, and Kashyap Chitta. Hidden biases of end-to-end driving datasets. In *CVPR Workshops*, 2024. 6, 7

[69] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, et al. RT-2: Vision-language-action models transfer web knowledge to robotic control. In *CoRL*, 2023. 1

## A. Datasets

We provide a more detailed description of the collected datasets and how we generate the labels for each language-related dataset.

### A.1. Driving dataset - Scenarios

In each Town, we collect data containing different scenarios, which we detail in the following (descriptions are taken from https://leaderboard.carla.org/scenarios/):

- **Control Loss without Previous Action:** The ego-vehicle loses control due to poor road conditions and must recover.
- **Unprotected Left Turn at Intersection with Oncoming Traffic:** The ego-vehicle performs an unprotected left turn at an intersection (can occur at both signalized and unsignalized intersections).
- **Right Turn at Intersection with Crossing Traffic:** The ego-vehicle makes a right turn at an intersection while yielding to crossing traffic (signalized and unsignalized intersections).
- **Crossing Negotiation at Unsignalized Intersection:** The ego-vehicle navigates an unsignalized intersection by negotiating with other vehicles. The assumption is that the vehicle entering the intersection first has priority.
- **Crossing Traffic Running a Red Light at an Intersection:** While the ego-vehicle is traveling straight through an intersection, a crossing vehicle runs a red light (signalized and unsignalized intersections).
- **Crossing with Oncoming Bicycles:** The ego-vehicle must turn at an intersection while yielding to bicycles crossing the intersection.
- **Highway Merge from On-Ramp:** The ego-vehicle merges into moving traffic on a highway.
- **Highway Cut-In from On-Ramp:** A vehicle merges into the ego-vehicle's lane from an on-ramp. The ego-vehicle must decelerate, brake, or change lanes to avoid a collision.
- **Static Cut-In:** Another vehicle cuts into the ego lane from a queue of stationary traffic. It must decelerate, brake, or change lanes to avoid a collision.
- **Highway Exit:** To exit the highway the ego-vehicle needs to cross a lane of moving traffic.
- **Yield to Emergency Vehicle:** An emergency vehicle is approaching from behind. The ego must create space for it to pass safely.
- **Obstacle in Lane - Same Direction:** An obstacle (e.g., a construction zone, an accident, or a parked vehicle) is blocking the ego lane. The ego vehicle must change lanes into traffic moving in the same direction to bypass the obstacle.
- **Obstacle in Lane - Opposite Direction:** An obstacle (e.g., construction zone, an accident, or a parked vehicle) is blocking the ego lane. The ego vehicle must change lanes into traffic moving in the opposite direction to bypass the obstacle.
- **Door obstacle:** The ego-vehicle needs to avoid a parked vehicle with its door opening into the lane.
- **Slow moving hazard at lane edge:** A slow-moving hazard (e.g. bicycle) partially obstructs the ego vehicle's lane. The ego-vehicle must either brake or carefully bypass the hazard (bypassing on lane with traffic in the same or opposite direction).
- **Vehicle invading lane on bend:** On a bend, an oncoming vehicle invades the ego vehicle's lane to avoid an obstacle. The ego-vehicle must brake or move to the side of the road to safely navigate past the oncoming vehicle.
- **Longitudinal control after leading vehicle's brake:** The leading vehicle in front of the ego-vehicle brakes suddenly to avoid an obstacle. The ego-vehicle must execute an emergency brake or avoidance maneuver to prevent a collision.
- **Obstacle avoidance without prior action:** The ego-vehicle encounters an unexpected obstacle or entity on the road. It must perform an emergency brake or avoidance maneuver.
- **Pedestrian emerging from behind parked vehicle:** A pedestrian suddenly emerges from behind a parked vehicle and enters the lane. The ego-vehicle must brake or take evasive action to avoid hitting the pedestrian.
- **Obstacle avoidance with prior action - pedestrian or bicycle:** While in the middle of a turn, the ego-vehicle encounters an obstacle such as a pedestrian or a bicycle crossing the road or a stopped vehicle in the road and must perform an emergency brake or avoidance maneuver.
- **Parking Cut-in:** A parked vehicle exits a parallel parking space into the ego-vehicle's path. The ego-vehicle must slow down to allow the parked vehicle to merge into traffic.
- **Parking Exit:** The ego-vehicle must exit a parallel parking space and merge into moving traffic.

### A.2. VQA - DriveLM

For the VQA data, we use the DriveLM-Carla [53] data generation method. Since we generate a new driving dataset, we extract questions and answers for our dataset instead of using the original dataset. For the training set, we generate in total 28M QA-pairs for 1M frames of Town 12. For the evaluation set, we use the keyframe extraction of DriveLM to evaluate on more interesting and less redundant frames. In addition, we balance the validation set to capture the same amount of samples for each answer type in the dataset.

Since the labels are auto-generated with a heuristic-based

procedure, the QAs follow the same sentence structures and wordings. To avoid overfitting to specific phrases we include data augmentation. For this, we prompt GPT-4 to generate 20 alternative sentences for each question and answer, from which we sample during data loading.

### A.3. Commentary

We generate language labels for the *Commentary* task based on a subset of the simulator state, which we save during the data collection. The structure of the *Commentary* labels is as follows: The first sentence describes the action according to the route with its justification (e.g., staying on the current lane, doing a lane change to go around an obstacle). It is followed by a description of the speed-related action (e.g., accelerate, keep the speed, stop), and the reason (e.g., because of the pedestrian, to follow the vehicle in front, to drive closer to the stop sign).

In the following, we detail the steps to obtain each part of the *Commentary* labels.

**Route action.** The default description is "*Follow the route.*". Only in special cases, we change the description. For this, we check if any scenario is active in the given frame and get the scenario type. We only change the default description for the scenarios requiring a deviation from the center lane of the original route (e.g., obstacle in lane, vehicle invading the lane, door obstacle). Since the ego action differs depending on the location relative to the scenario, we extract the relative positioning from the simulator information. Those locations are grouped into three phases: (1) before the lane deviation, (2) during the deviation, and (3) end of the deviation. *Before the deviation*, we differentiate between the scenario types and use a template sentence for each, for instance:

- Overtake the bikes on your lane.
- Go around the vehicle with the open door.
- Give way to the emergency vehicle.
- Go around the accident in your lane.
- Go around the construction site.
- Move slightly to the right to circumvent the oncoming cars entering your lane because of the construction cones.

*During the deviation*, describes the phase in which the ego already shifted lanes. We reuse the templates from (1) but add "Stay on your current lane to" before the templates (e.g., Stay on your current lane to overtake the bikes on your lane.)

*End of the deviation* is the phase where the ego vehicle needs to shift back to its original lane. Since our model is based on only front-view cameras we use a generic sentence for this (i.e., "Return to your original route after avoiding the obstacle.") as often the type of the obstacle or scenario is not visible anymore. This template can be easily changed for a model supporting multi-view inputs.

**Speed action.** We generate a high-level description of the ego action based on the current speed, the desired target speed based on the expert decision, and the current speed limit. We differentiate between the following types:

- Remain stopped
- Come to a stop now
- Maintain your current speed
- Maintain the reduced speed
- Increase your speed
- Slow down

For the scenarios that are used for the *route action* description (i.e., where the expert needs to deviate from the route), we use a different sentence template. This is only the case for the situation when the ego vehicle is *before the deviation* and is stopped and remains stopped for the next two seconds. In this case, we use the template "Wait for a gap in the traffic before changing lanes".

**Speed Reason.** Next, we leverage the Inteligent-Driver-Model (IDM) features, which the expert algorithm is based on. IDM identifies leading objects and calculates the optimal target speed for the ego vehicle based on the distance to the leading object. Leading objects include dynamic objects like other vehicles or pedestrians and static objects like traffic lights, stop signs, or construction sites. With this, we know for any sample in the dataset which object the main reason is for the given target speed and therefore the *speed action*. Based on the type of the leading object, we construct a language description of the object. For vehicles, this consists of the color of the vehicle, the type (e.g., SUV, police car), and a rough position relative to the ego (e.g., to the front right). For static objects, it is the name of the object (e.g. traffic light) and in case the object has a state this is also included (e.g., *red* traffic light). For pedestrians, we differentiate between children and adults. Based on the type of the leading object and the *speed action* we construct different sentences, for instance:

- since you've already stopped at the stop sign
- to avoid a collision with the *object description*
- due to the pedestrian crossing in front of you
- to remain behind the red SUV that is slowing down because of a red traffic light.
- to reach the target speed
- because the traffic light is green

If we are right before a junction we add another notice label regarding the positioning of the other vehicles in the junction. With this, the model needs to reason about whether it is safe to enter a junction or not. We collect the position and driving direction of each vehicle close to the junction and summarize the situation based on one of the following sentences:

- the other vehicles are stopped at the junction and the junction is clear
- the other vehicles are stopped at the junction and the vehicle in the junction is moving away

- pay attention to the vehicles coming towards the junction
- pay attention to the vehicle in the junction

### A.4. Action Dreamer

We construct an offline, non-reactive simulation based on the collected dataset to generate alternative ego trajectories and evaluate their feasibility in terms of collision avoidance and adherence to traffic rules. For this purpose, we utilize the Kinematic Bicycle Model in combination with the PID controllers from the PDM-lite expert algorithm.

The core functionality of the Action Dreamer simulation is the *ego forecasting*, which predicts future ego vehicle poses based on the ego actions in each timestep. There are several approaches to generate these ego actions, allowing for modification to obtain the alternative trajectories. One approach involves perturbing the ground truth actions to produce slightly modified trajectories. Another approach uses the PID controllers to compute actions based on pre-defined path waypoints and target speeds. In this case, the lateral PID controller generates steering angles, while the longitudinal PID determines acceleration and braking values based on the desired target speed. Using these actions, the Kinematic Bicycle Model calculates the next vehicle pose. This process can be iteratively unrolled over multiple time steps to derive a complete trajectory.

We start with obtaining the current state of the simulator from the saved dataset. For each dynamic object, we also get the states for the following 10 timesteps. With this we can get the non-reactive trajectories for each object and perform collision checks with the ego vehicle. As default, we use the ground truth actions of the ego vehicle. In addition, we use the ground truth path provided by the experts' path planner, which includes waypoints spaced every 10 cm, as the default route to be followed. We then change those default values to obtain alternative trajectories for the modes: objects (collision), faster, slower, target speed, and lane changes.

The following steps describe how we obtain the necessary information (e.g., actions, updated paths with desired speeds, or a combination of both) for the ego forecasting method for each of the different modes we have.

- **Objects (Collision) Mode:** Filter all dynamic and static objects, retaining only those within a 15-meter radius of the ego route and at least 3 meters ahead of the ego vehicle. For each object, calculate its position, distance to the ego vehicle, and whether the ego vehicle (given its current speed) could reach the object within a given future timestep. Objects that are too far and unreachable are discarded. For reachable objects, we calculate the target speed required for the ego vehicle to precisely reach the object's location in the required time. For the path, we adjust the route waypoints to ensure they intersect with the

| Epochs | 14 |
|---|---|
| Learning Rate | 3e-5 |
| Batch Size | 96 |
| Optimizer | AdamW |
| Weight decay | 0.1 |
| Betas | (0.9, 0.999) |
| LR schduler | One cycle cosine |
| Warmup steps | 5% of total steps |
| LoRA alpha | 64 |
| LoRA r | 32 |
| LoRA dropout | 0.1 |

Table 7. **Hyperparameter** choices to train SimLingo.

target object's position. This modified route, along with the computed target speed, is then passed to the ego forecasting process.
- **Faster Mode:** We use the original path and actions for steering but set the acceleration to a random value above 50%.
- **Slower Mode:** Similar to the faster mode but we set the acceleration to zero and activate the brake.
- **Random Target Speed:** We assign a random target speed between 0 and 35 m/s and directly pass this together with the default path to the forecasting method.
- **Lane Changes:** We exclude frames where the vehicle is already performing a lane change or is in a junction. To obtain the number of possible options, we extract information on the number and types of lanes (e.g., driving lanes, parking lanes, or sidewalks) in both the same and opposite directions. For each of the options, we change the default path so that it reaches the specified lane. We randomize the starting distance of the lane change and the length of the transition phase. Those parameters are conditioned on the current ego speed.

## B. Implementation details

### B.1. Hyperparameter

Table 7 shows the hyperparameter we use to train Sim-Lingo.

### B.2. Training buckets

The majority of driving involves straight, uneventful segments. To maximize the collection of interesting scenarios during data collection, we focus on capturing a diverse range of challenging situations. However, some ratio of easy and uneventful data is inevitable. Training models on the entire dataset revealed that straight driving without hazards is effectively learned in the early epochs, resulting in wasted computation in later epochs as the models continue to train on these uninteresting samples. To address this issue, we create data buckets containing only the interesting
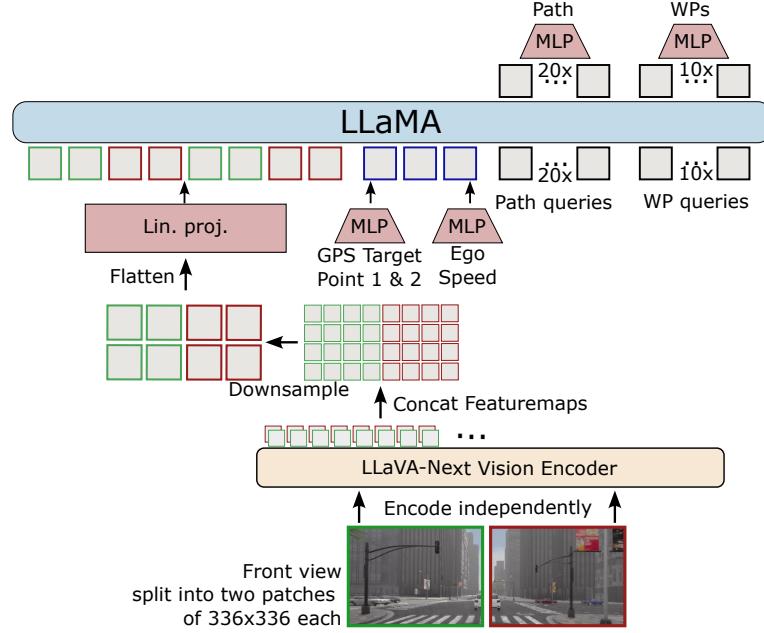
Figure 5. **SimLingo-BASE architecture.** The images are split in two, and each split is independently encoded and then concatenated, downsampled, and projected before feeding it into a transformer decoder which is based on the LLaMA architecture. The output utilizes the same output representation as SimLingo.

samples and sample from these buckets during training instead of the entire dataset. We use (1) five buckets for different amounts of acceleration and deceleration with one specifically for starting from stop, excluding samples with an acceleration between -1 and 1, (2) two buckets for steering, excluding samples for going straight, (3) three buckets for vehicle hazard with vehicles coming from different directions, (4) one for a stop sign, red light, and walker hazards each, (5) one bucket for swerving around obstacles, (6) one bucket for "old" Towns commonly used in Leaderboard 1.0 models (Town 01-10) and (7) one bucket that samples from the whole dataset to keep a small portion of uneventful data such as driving straight.

### B.3. SimLingo vs. SimLingo-BASE

Our base model SimLingo-BASE was designed as a lightweight model to research driving-specific design choices. With adding language capabilities, we also changed some of the settings to better match the added requirements. We note that because of the closure of the CARLA Leaderboard and because of computational overhead we did not repeat the experiments of the SimLingo-BASE with the new settings. We believe that the claims and results can be expected to still hold. Fig. 5 shows the architecture of SimLingo-BASE. We detail the exact differences between SimLingo-BASE and SimLingo:

1. Number of epochs: SimLingo-BASE is trained for 30 epochs. For SimLingo we reduce the number of epochs to 14.

2. Image encoder: SimLingo-BASE uses a Clip-ViT that is used as default in the LLaVA VLM. SimLingo uses the original image encoder of the InternVL2-1B. In both cases, we do full finetuning.

3. Language model: SimLingo-BASE uses a 50M parameter transformer decoder based on the LLaMA architecture which we train from scratch. SimLingo uses the default pretrained LLM from the InternVL2-1B model which we finetune with LoRA.

4. Loss function: SimLingo-BASE uses L2-Loss. After adding the additional Action Dreaming data we observed instabilities during training with the L2-Loss so we changed to the SmoothL1-Loss.

### B.4. Metric descriptions

#### Leaderboard 2.0.

The Leaderboard uses the official CARLA metrics: Driving Score, Route Completion, and Infraction Score. Each metric is calculated for each route independently. After all routes are completed, the final metrics are derived by taking the arithmetic mean of the metrics across all routes. The overall driving score, calculated using the global values, is the primary metric for ranking methods.

*Driving Score.* The primary evaluation criterion is the *Driving Score*, denoted as:

$$DS_i = RC_i \cdot IS_i,$$

where $RC_i$ represents the percentage of the $i$-th route completed, and $IS_i$ is a penalty factor accounting for infractions

incurred during the route.

*Route Completion.* This metric quantifies the proportion of the route successfully completed by the agent, expressed as a percentage.

*Infraction Penalty.* The penalty due to infractions, $IS_i$, is calculated as a product of all infractions:

$$IS_i = \prod_{j=1}^{N_I} (p_j)^{\#\text{infractions}_j},$$

where $p_j$ denotes the penalty coefficient for the $j$-th type of infraction out of a total of $N_I$ infraction types, which we specify int following. $\#\text{infractions}_j$ is the number of times this infraction occurred. The calculation begins with a base score of 1.0, which decreases with each infraction.

Infractions are categorized by severity, each associated with a penalty coefficient that reduces the driving score. Key infractions include:

- **Collisions with pedestrians:** $p_j = 0.50$.
- **Collisions with vehicles:** $p_j = 0.60$.
- **Collisions with static objects:** $p_j = 0.65$.
- **Running a red light:** $p_j = 0.70$.
- **Ignoring a stop sign:** $p_j = 0.80$.
- **Failure to yield to emergency vehicles:** $p_j = 0.7$.
- **Failure to maintain minimum speed:** Up to $p_j = 0.7$.
- **Off-road driving:** Reduces route completion score proportionally.

When one of the following events occurs, the route stops immediately:

- Route deviation (more than 30 meters off route).
- Blocked agent (more than 180 simulation seconds without action).
- Communication timeout (more than 60 seconds).
- Route timeout (exceeding allowed simulation time).

**Leaderboard 2.0 metric discussion.**

The Driving Score is calculated in a way, that it can be advantageous not completing the whole route. This is the case if the infractions incurred during a segment of the route reduce the driving score more than the potential gain from continuing the route. In this case stopping early to avoid further penalties leads to an overall higher driving score. This tradeoff only occurs for long routes. We refer to [67] for a mathematical description of this tradeoff, a detailed discussion and a proposal for a better metric calculation.

**Bench2Drive.**

*Driving Score.* The Driving Score is calculated similarly to the Leaderboard 2.0. The only difference to the original Driving Score is that the penalty "Failure to maintain minimum speed" is ignored. To take the driving speed into account the authors of the benchmark introduced the "Efficiency" metric. Since Bench2Drive uses short routes the discussed trade-off does not occur on this benchmark.

*Success Rate.* The Success Rate measures the percentage of completed routed without any infractions (ignoring the minimum speed penalty).

*Efficiency.* This uses the ratio of the ego vehicle speed to the speed of the surrounding actors. Since there was no penalty in Leaderboard 1.0 for low speeds most models used a very low speed, which makes driving and reacting to other dynamic actors much easier. The higher this efficiency metric the faster the model drives, making the driving task harder.

*Comfortness.* The comfortness metric takes the jerk magnitude, lateral and longitudinal acceleration, yaw acceleration, longitudinal jerk, and the yaw rate into account. If the mean of the ego vehicles measurement over the full route falls into the following thresholds it is treated as success.

- **Jerk Magnitude:** Maximum absolute magnitude of jerk is $8.37$ m/s$^3$.
- **Lateral Acceleration:** Maximum absolute lateral acceleration is $4.89$ m/s$^2$.
- **Longitudinal Acceleration:**
  - Maximum longitudinal acceleration is $2.40$ m/s$^2$.
  - Minimum longitudinal acceleration is $-4.05$ m/s$^2$.
- **Yaw Acceleration:** Maximum absolute yaw acceleration is $1.93$ rad/s$^2$.
- **Longitudinal Jerk:** Maximum absolute longitudinal jerk is $4.13$ m/s$^3$.
- **Yaw Rate:** Maximum absolute yaw rate is $0.95$ rad/s.

The thresholds are taken from the Bench2Drive repository.

**Vision-Language Understanding.**

For the tasks of VQA and Commentary, we use SPICE [5] and DriveLM's GPT Score [53] as metrics. SPICE is a metric used for image captioning with a higher correlation to human judgment than other automatic language metrics like Cider [57] or Meteor [33]. The GPT Score is based on the DriveLM implementation with two smaller changes. Since we could not directly compare to their numbers, because of a different evaluation set those changes should not have any impact on the conclusions drawn. The first change is using GPT-4 (gpt-4o-2024-08-06) instead of GPT-3.5. In addition, we add "Just rate the similarity of the content not the sentence structure. If the content is completely different rate with 0." to the prompt as we found this to be more accurate.

**Action Dreaming.**

For the Action Dreaming evaluation, we use Success Rate as the metric. Each category has its own definition of success, which we detail in the following:

- **Slow down**: We calculate the target speed for each waypoint of the predicted speed waypoints. Those target speeds represent the target speeds for future timesteps. We do linear regression to get the slope.

| | Method | Ability (%) ↑ | | | | | |
|---|---|---|---|---|---|---|---|
| | | Merging | Overtaking | Emergency Brake | Give Way | Traffic Sign | **Mean** |
| w/ dist. | TCP [62] | 16.18 | 20.00 | 20.00 | 10.00 | 6.99 | 14.63 |
| | TCP-ctrl | 10.29 | 4.44 | 10.00 | 10.00 | 6.45 | 8.23 |
| | TCP-traj | 8.89 | 24.29 | 51.67 | 40.00 | 46.28 | 28.51 |
| | ThinkTwice [29] | 27.38 | 18.42 | 35.82 | 50.00 | 54.23 | 37.17 |
| | DriveAdapter [28] | 28.82 | 26.38 | 48.76 | 50.00 | 56.43 | 42.08 |
| w/o dist. | AD-MLP [65] | 0.00 | 0.00 | 0.00 | 0.00 | 4.35 | 0.87 |
| | UniAD-Tiny [24] | 8.89 | 9.33 | 20.00 | 20.00 | 15.43 | 14.73 |
| | UniAD-Base [24] | 14.10 | 17.78 | 21.67 | 10.00 | 14.21 | 15.55 |
| | VAD [31] | 8.11 | 24.44 | 18.64 | 20.00 | 19.15 | 18.07 |
| | TCP-traj w/o distillation | 17.14 | 6.67 | 40.00 | 50.00 | 28.72 | 28.51 |
| | SimLingo-BASE (LB2.0 model) | **60.00** | **60.00** | 78.33 | 50.00 | 77.89 | 65.25 |
| | SimLingo | 54.01±2.63 | 57.04±3.40 | **88.33±3.34** | 53.33±5.77 | 82.45±4.73 | **67.03±2.12** |

Table 8. **Multi-Ability Results of Bench2Drive.** We outperform the existing methods in all abilities and can improve in average by 25 percentage points.

| | DS ↑ | Stat ↓ |
|---|---|---|
| WPs | 3.21 | 0.68 |
| +Path | 4.49 | 0.0 |

(a) **Output.**

| | DS ↑ |
|---|---|
| Clip ViT | 6.87 |
| w/o pretr. | 0.45 |
| Resnet-34 | 2.71 |

(b) **Vision encoder**.

| | DS ↑ |
|---|---|
| 1300 | 3.93 |
| 1800 | 4.49 |
| 2100 | 6.87 |
| 2400 | 6.35 |

(c) **Early stopping**.

Table 9. Ablations of different parts of SimLingo-BASE, showcasing the superiority of the disentangled output representation and the large impact of the correct threshold for early stopping. The score of the default configuration is highlighted in gray. All numbers are official Leaderboard 2.0 scores.

Success is defined as the slope being smaller than $-0.05 * v$, with $v$ being the current ego speed.

- **Speed up**: Same calculation as for *Slow down* but we use $slope > 0.05 * v$.
- **Target Speed**: Since we do not know if the target speed can be reached in the prediction horizon of the waypoints we compare the predictions with the ground truth actions instead of directly comparing to the target speed. We use two rules defining success: First, if the predicted target speed inferred from the last two waypoints is in a 20% range of the instructed target speed. Second, if the predicted target speed inferred from the last two waypoints is in the 20% range of the speed of the last two waypoints of the ground truth speed waypoints. This can be different from the instructed target speed due to limitations in the acceleration rates of the vehicle.
- **Lane Change**: We compare the final waypoint of the predicted path waypoint with the ground truth dreamer path and the ground truth expert path waypoints. We define the lane change as successful when the predicted final location is closer to the dreamer's final location than the expert final location.

- **Objects (Collisions)**: This describes the task of driving towards or crashing into specific objects. We first look at the path. If the path of the expert trajectory and the ground truth dreamer trajectory is different (Average Displacement Error $ADE > 1.0$) we count it as success if the predicted path is closer to the ground truth dreamer path than to the expert path ($ADE_{pred2expert} > ADE_{pred2dreamer}$). If the dreamer path is nearly identical to the expert path ($ADE < 1.0$) the instruction is about correct speed predictions (e.g., if the instruction is "drive towards a dynamic object" it is important to get the speed right and not just the path. The success is then defined as $ADE_{pred2dreamer} < 1.0$ and the average predicted speed is within 30% of the ground truth dreamer speed.

## C. Quantitative Results

### C.1. Ablations on CARLA Leaderboard 2.0

We provide additional results on SimLingo-BASE. *Output representation.* Tab. 9a compares the DS on the Leaderboard for the different output representations. As the goal of the additional path prediction is improved lateral control, we also report the collisions with static layout as this is mainly caused due to bad steering. With the disentangled representation, we can reduce the layout collision from 0.68 to 0 showing the strength of additional path predictions.

*Vision-Language and CLIP pretraining.* We ablate the pretraining of the vision encoder and train the same model from scratch. Tab. 9b 'w/o pretr.' shows that the pretraining stage is essential for good driving performance (longer training can further improve the performance but is unlikely to reach the performance of the pretrained model). Additionally, we show a comparison to the widely used Resnet-34 pretrained on ImageNet. The decreased performance (2.71 vs.

|         | Bench2Drive | |
|---------|:-----------:|:-------------:|
|         | DS ↑ | SR(%) ↑ |
| w/o CoT | 84.41±1.76 | 64.84±2.42 |
| with CoT | **85.07±0.95** | **67.27±2.11** |

Table 10. **Inference Mode**. BEnch2Drive results for using Sim-Lingo with chain-of-thought (CoT) during inference and without. We see a small improvement when using CoT.

6.87 DS) indicates the importance of the larger ViT and the internet-scale image-language pretraining.

*Early stopping.* As described in the metric section the DS on long routes is not optimal and favors models that do not complete the full route. We ablate the thresholds for the early stopping as it is not trivial to calculate the perfect trade-off as the routes and density of scenarios are secret (however a rough function of the expected DS can be calculated, which we used to get a rough range). Tab. 9c shows the Leaderboard DS for a given traveled distance in meters. This hyperparameter has a big impact on the final score.

*Leaderboard variance.* We submitted our base model SimLingo-BASE with an early stopping threshold of 2100 and 2400 three times to the leaderboard to get an estimate of the evaluation variance. For the 2100 model, we obtain the following scores: 6.9, 5.5, and 5.3 resulting in a mean DS of 5.9 with a standard deviation of 0.87. The base model with a threshold of 2400 obtained 6.4, 6.3, and 4.8 resulting in a mean of 5.83 with a standard deviation of 0.90.

### C.2. Bench2Drive Multi-Ability Results

Tab. 8 shows the Bench2Drive Multi-Ability metrics. Consistent with the findings in the main paper SimLingo outperforms existing works. Especially in the abilities *Merging, Overtaking, Emergency Brake* and *Traffic Sign* we get a large boost in performance. *Give way* is still challenging.

### C.3. Chain-of-Thought Inference Mode

We ablate the Chain-of-Thought (CoT) inference mode of SimLingo in Tab. 10. When using the Commentary task as CoT we see a small improvement so we decided to use this as the default mode. However, since without CoT the performance does not drop much, using the model without CoT is a feasible option especially when inference speed is important.

### D. Qualitative Results

We provide more qualitative results of different navigational commands in Fig. 6. Red dots are the predicted path waypoints and green are the predicted speed waypoints. Each row is captured from a closed-loop run while changing the navigation command. The second row shows how the model can successfully differentiate between different situations and adapt its behavior given a certain command. The vehicle starts in the right lane. When giving the instruction "Turn left..." the model initiates a lane change to the left lane. After finishing this lane change it stays on this lane even though the command is still "Turn left...". So the model learns to reason about the meaning of the different lanes, that it is forbidden to go on an oncoming lane, and that "Turn left" not always mean to do a 90-degree turn. In the third row, we prompt the model with a misleading instruction: "Turn right" on an intersection without a lane going to the right. The third and fourth image shows that the model is slightly confused but when unrolling closed-loop the car still goes straight and stays on the road. In the fourth row, we show out-of-distribution commands. The model is still able to choose a valid path when using a command that does not make sense like "I really like my dog". When using the command "Why is there a tree on the right side?" the model still picks the left turn, indicating that it does not just overfit on single words like *left* or *right* but also takes the context into account. The last image with the command "Do a U-Turn" shows that the model is not capable of following concepts it has never seen during training.

Fig. 7, Fig. 8, Fig. 9, Fig. 10 and Fig. 11 shows qualitative results for the different Dreamer modes. In most cases, the model correctly follows the instructions even if it clearly goes against the expert driving behavior (e.g., accelerating at a red traffic light). We also include some of the failure cases of the model in the red boxes where the model ignores the instructions.
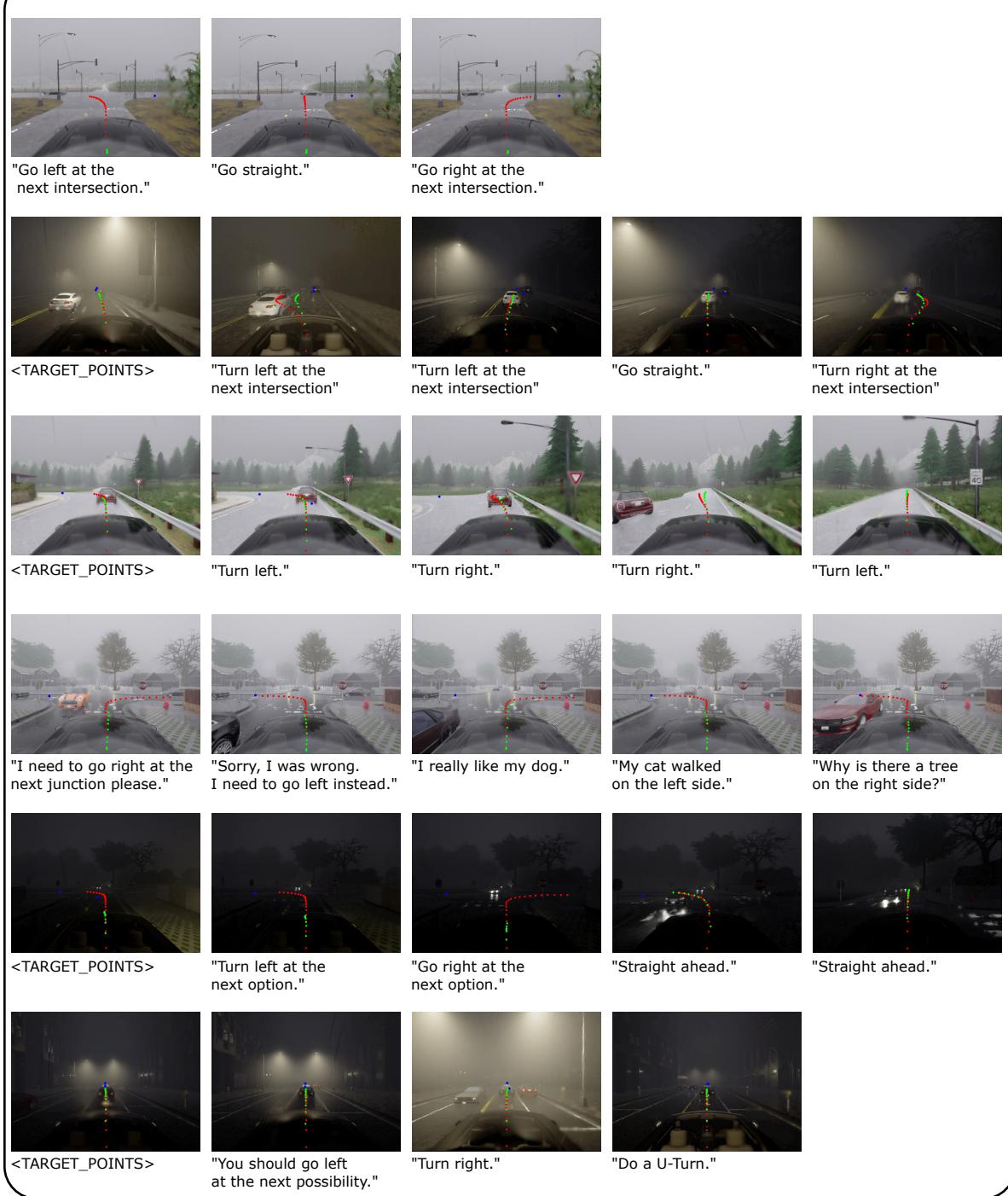
Figure 6. **Navigational Commands.** We show closed-loop results for different in-distribution and out-of-distribution commands. Red: path waypoints, green: speed waypoints.
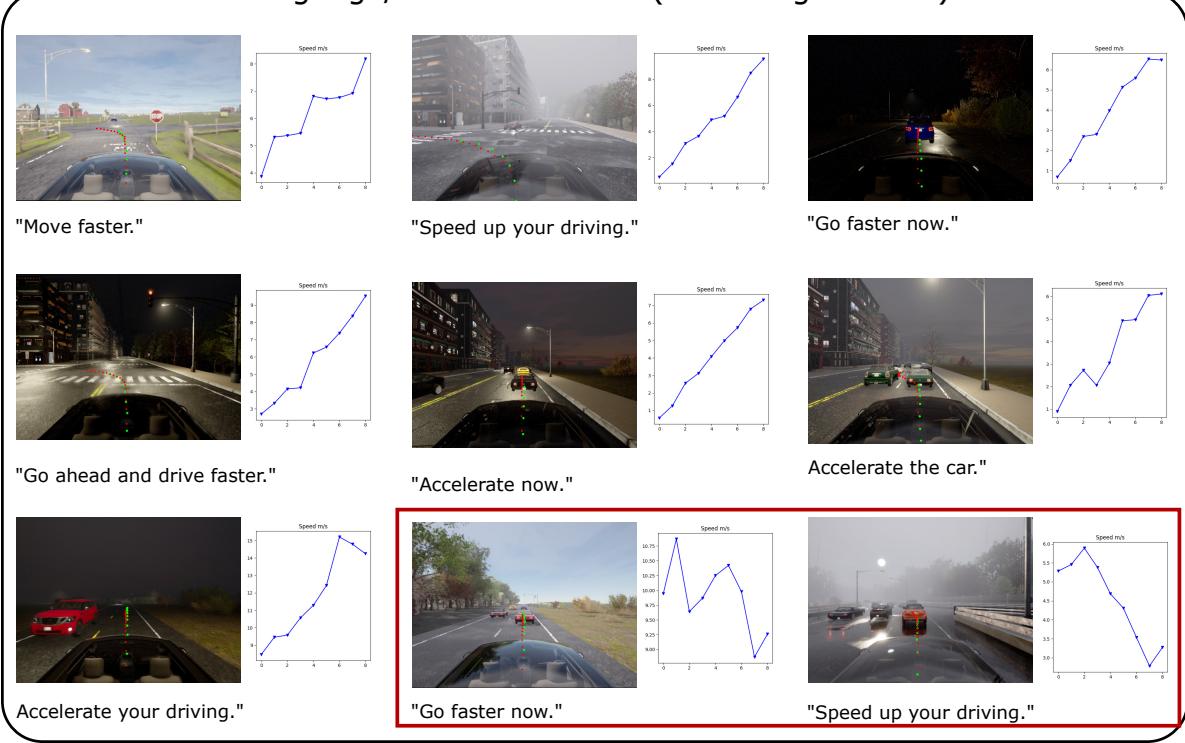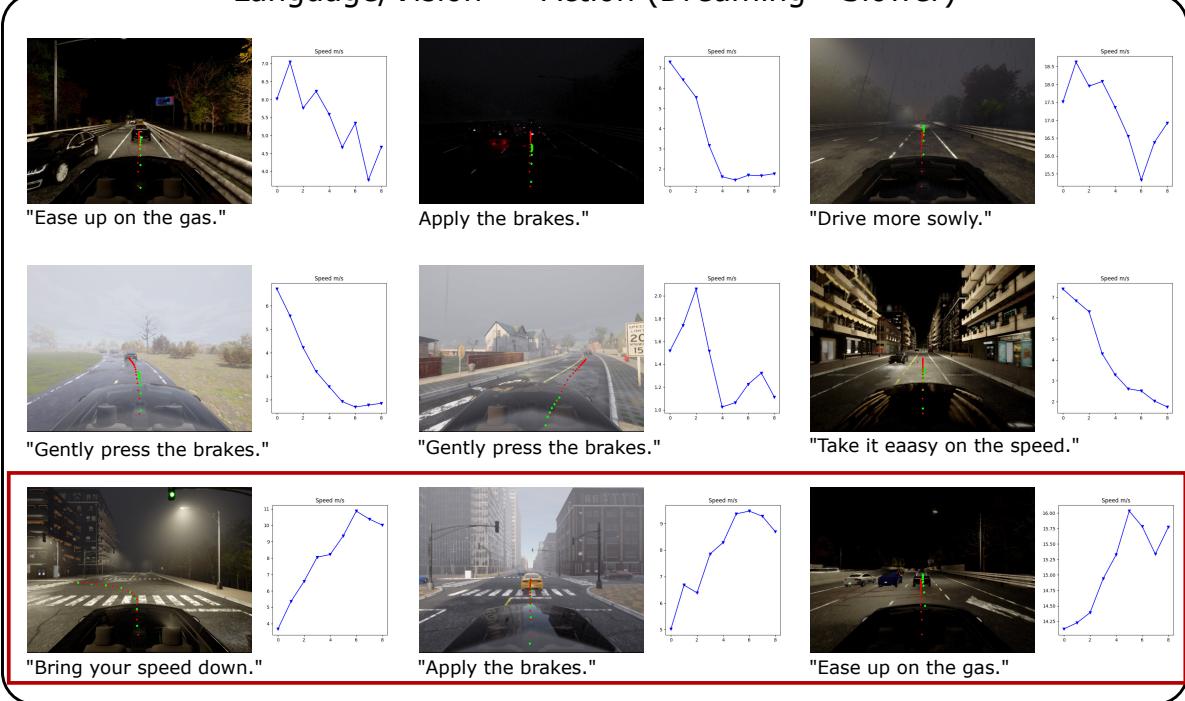
Figure 7. **Dreamer mode - Faster.** We show in blue the predicted speed curve (inferred from the speed waypoints). Inside the red border are examples where the model does not follow the command correctly. Red: path waypoints, green: speed waypoints.



Figure 8. **Dreamer mode - Slower.** We show in blue the predicted speed curve (inferred from the speed waypoints). Inside the red border are examples where the model does not follow commands correctly. Red: path waypoints, green: speed waypoints.
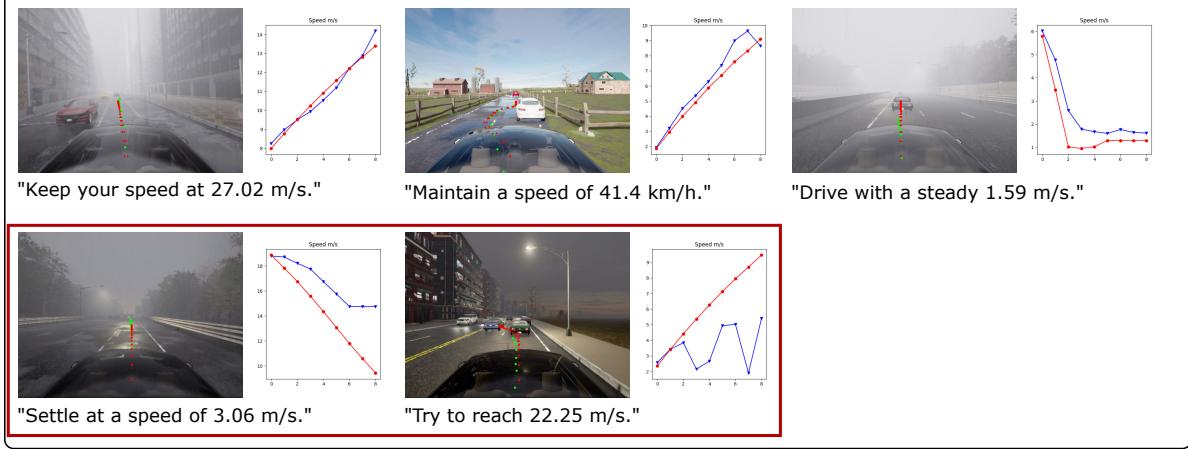
9

Figure 9. **Dreamer mode - Target speed.** We show in red line plot the ground truth speed curve and in blue the predicted one (inferred from the speed waypoints). Inside the red border are examples where the model does not follow commands correctly. Red: path waypoints, green: speed waypoints.
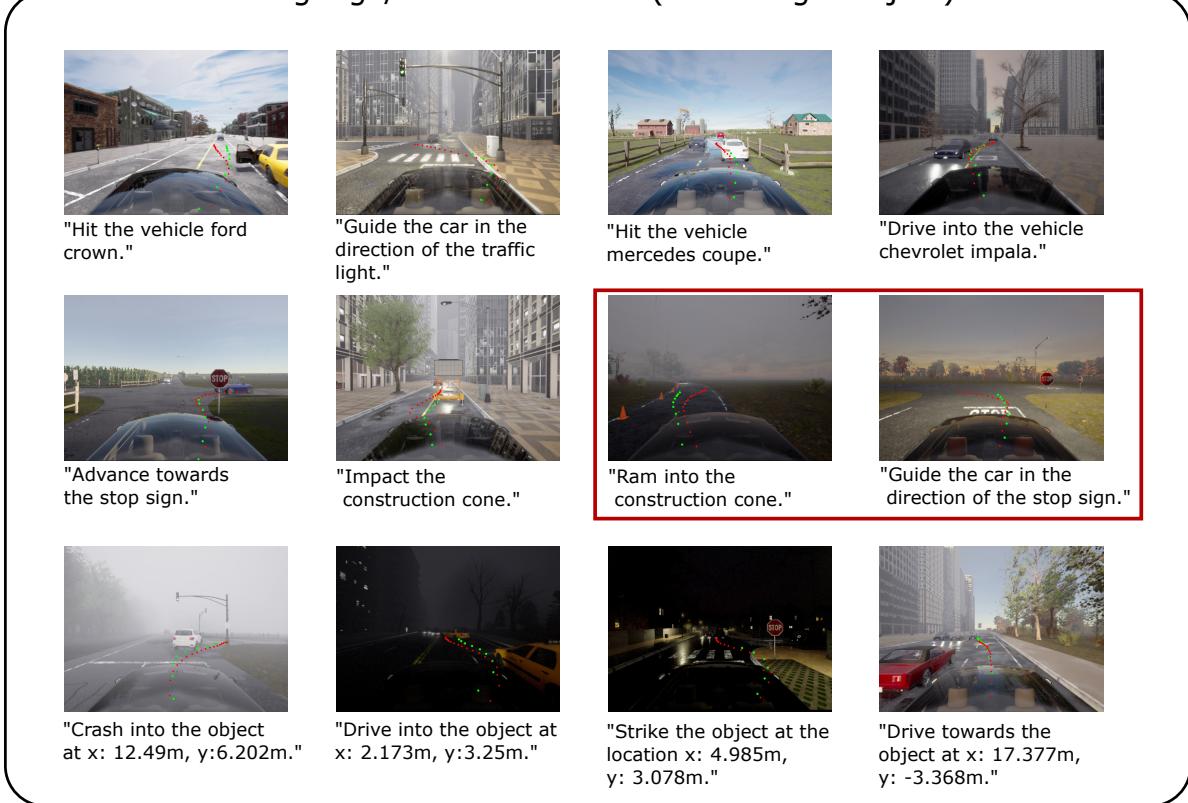


Figure 10. **Dreamer mode - Object (Collision).** Inside the red border are examples where the model does not follow dreamer mode command correctly. Red: path waypoints, green: speed waypoints.

Figure 11. **Dreamer mode - lane changes.** Inside the red border are examples where the model does not follow dreamer mode commands correctly. Red: path waypoints, green: speed waypoints.