

Think Before You Drive: World Model-Inspired Multimodal Grounding for Autonomous Vehicles

Haicheng Liao¹, Huanming Shen², Bonan Wang¹, Yongkang Li³,
Yihong Tang⁴, Chengyue Wang¹, Dingyi Zhuang⁵, Kehua Chen⁶,
Hai Yang⁷, Chengzhong Xu¹, Zhenning Li¹

¹University of Macau, ²UESTC, ³Purdue University, ⁴McGill University, ⁵Massachusetts Institute of Technology, ⁶University of Washington, ⁷The Hong Kong University of Science and Technology

 Equal Contribution,  Corresponding Authors

Abstract

Interpreting natural-language commands to localize target objects is critical for autonomous driving (AD). Existing visual grounding (VG) methods for autonomous vehicles (AVs) typically struggle with ambiguous, context-dependent instructions, as they lack reasoning over 3D spatial relations and anticipated scene evolution. Grounded in the principles of world models, we propose **ThinkDeeper**, a framework that reasons about future spatial states before making grounding decisions. At its core is a Spatial-Aware World Model (SA-WM) that learns to reason ahead by distilling the current scene into a command-aware latent state and rolling out a sequence of future latent states, providing forward-looking cues for disambiguation. Complementing this, a hypergraph-guided decoder then hierarchically fuses these states with the multimodal input, capturing higher-order spatial dependencies for robust localization. In addition, we present **DrivePilot**, a multi-source VG dataset in AD, featuring semantic annotations generated by a Retrieval-Augmented Generation (RAG) and Chain-of-Thought (CoT)-prompted LLM pipeline. Extensive evaluations on six benchmarks, ThinkDeeper ranks #1 on the Talk2Car leaderboard and surpasses state-of-the-art baselines on DrivePilot, MoCAD, and RefCOCO+/g benchmarks. Notably, it shows strong robustness and efficiency in challenging scenes (long-text, multi-agent, ambiguity) and retains superior performance even when trained on 50% of the data.

Date: November 25, 2025

1. Introduction

In recent years, the pursuit of fully autonomous vehicles (AVs) has captivated both industry and academia [54]. Despite remarkable technological advances, widespread public acceptance remains elusive, primarily due to concerns over the reliability of human-machine interaction and apprehensions about losing control. These challenges are magnified in complex scenes where vehicles must make split-second decisions, highlighting the need for enhanced communication between humans and machines. Visual grounding (VG) in AVs, where vehicles interpret and act on natural language commands, emerges as a pivotal innovation, empowering passengers with a direct and intuitive mode of interaction that significantly enriches the driving experience. VG in autonomous driving (AD) requires a deep understanding of both the immediate and forthcoming environment [10]. For instance, when a passenger says “merge behind the white

SUV after the crosswalk”, the AV must jointly reason about spatial relationships and textual intent: the distance to two similar SUVs at different depths, the approach of a cyclist entering the blind spot, and the state of a traffic signal that reduced to a few pixels in its field of view [17, 31]. This grounding challenge is further complicated by the inherent variability and ambiguity of natural language commands, which are often highly dependent on the driving context.

Traditional VG methods [21, 28, 65] are ill-suited for the demands of real-world AD due to two primary issues. First, they are typically designed for high-resolution, controlled datasets, and struggle with the challenging visual conditions encountered on the road, such as low light and fast-moving scenes (as detailed in Table 1). These real-world conditions often obscure critical details, causing the methods to frequently miss vital contextual cues in ambiguous or rapidly changing situations. Second, existing VG models in AD [4, 5, 35, 36, 43, 53] generally lack the 3D spatial awareness

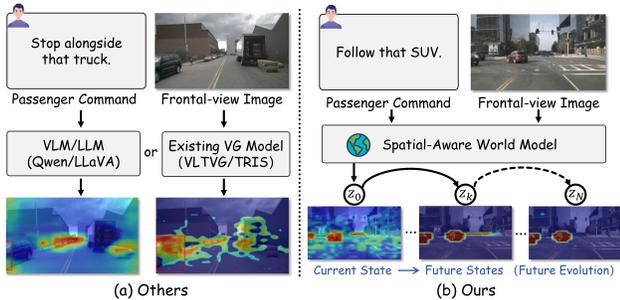


Figure 1. Comparison of visual grounding performance between existing VLMs and our model in real-world driving scenarios. (a) General-purpose VLMs (Qwen [55]/LLaVA [42]) and existing VG models [41, 60] typically fail to robustly localize targets under challenges such as motion blur, ambiguous language, and multi-agent traffic contexts. (b) In contrast, our Spatial-Aware World Model first distills the scene and command into a current latent state (z_0), and then reasons ahead by rolling out a sequence of future latent states ($z_k \rightarrow z_N$), obtaining a forward-looking perspective that enables more reliable and spatially coherent grounding.

required for complex relational reasoning. This lack of spatial intelligence limits their ability to accurately assess object distances or differentiate between immediate hazards and distant background elements. This shortcoming is particularly evident in context-dependent commands like “Avoid the cyclist ahead”, where the model must determine which cyclist is in immediate proximity and requires action, rather than simply identifying any cyclist in the scene [37].

In parallel, recent studies [46, 55] have integrated Large Language Models (LLMs) and Vision-Language Models (VLMs) to improve semantic understanding [46, 55]. While these models can mitigate ambiguity, they introduce significant drawbacks: massive data requirements, prohibitive computational costs, and high inference latency, impeding real-time deployment in AD. These limitations provoke us to ask an important question: *How can we design a visual grounding model for AD that is spatially aware, robust to ambiguity, and efficient enough for real-time operation?*

To alleviate this problem, we introduce a world model-based framework that empowers AVs to “**think deeper**” by reasoning about how the future scene is likely to evolve, effectively bridging the gap between rigid algorithmic processing and human-like contextual reasoning. Specifically, we introduce a Spatial-Aware World Model (SA-WM) to predict future states from the current scene and command, providing a forward-looking perspective for the grounding decision. As shown in Figure 1, it first constructs a command-aware current latent state that filters out irrelevant elements (e.g., roadside buildings) and then iteratively rolls out future latent states to foresee critical cues essential for decision-making. By reasoning over these prospective states before decoding, the model attains more stable decisions that improve gen-

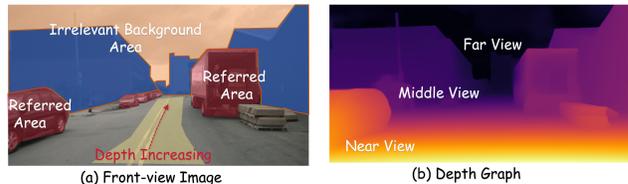


Figure 2. Illustration of depth-based spatial priors. (a) Real-world scenes have a clear 3D structure, with referred areas (nearby vehicles) at different distances than the irrelevant background (distant sky, mid-range buildings). (b) The depth map indicates depth semantics across near, middle, and far views, enabling the SA-WM to enhance its spatial awareness and filter implausible regions.

Dataset	Images	Objects	Expr. Length	Image Quality	LLM Annotations
RefCOCO [30]	19,994	50,000	3.51	High/Static	✗
RefCOCO+ [30]	19,992	49,856	3.53	High/Static	✗
RefCOCOg [30]	26,711	54,822	8.43	High/Static	✗
ReferIt [30]	19,894	96,654	3.46	High/Static	✗
Talk2Car [12]	9,217	10,519	11.01	Mid/Dynamic	✗
DrivePilot (Ours)	32,264	80,684	14.72	Mid/Dynamic	✓

Table 1. Comparison of DrivePilot with other VG datasets.

eralization and safety in AD grounding. To enhance the model’s spatial awareness, we integrate monocular depth to provide 3D localization for the vision-only pipeline. As illustrated in Figure 2, this depth signal enables the SA-WM to prioritize entities by proximity and relevance, mimicking human spatial perception. Complementing this, we introduce a cross-modal hypergraph decoder that captures higher-order relations between textual phrases and spatial regions, effectively fusing the predicted future latent states from the SA-WM to achieve robust grounding. Overall, the main contributions of this study can be summarized as follows:

- We introduce **ThinkDeeper**, the first world model-based framework for visual grounding in AD. ThinkDeeper predicts future latent states from the current scene, provides a forward-looking perspective to guide AVs to reason about the environment’s future evolution before grounding.
- We present **DrivePilot**, a multi-source dataset featuring semantic annotations generated via Retrieval-Augmented Generation (RAG) and Chain-of-Thought (CoT) prompting with LLMs. It provides a robust benchmark of dynamic, real-world driving scenes to advance VG research.
- ThinkDeeper sets SOTA on DrivePilot and MoCAD, **ranks #1** on Talk2Car C4AV Challenge leaderboard, and generalizing to excel on the RefCOCO+/g datasets. Notably, it remains robust and efficient in challenging scenes (long-text, multi-agent, ambiguity), outperforming most state-of-the-art baselines even with 50% and 75% training data.

2. Related Work

Visual Grounding in Autonomous Driving. Visual grounding has emerged as a critical component for enhancing

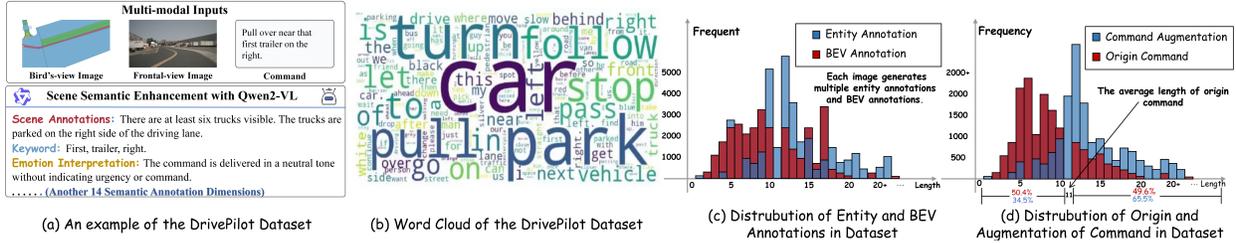


Figure 3. Overview of the proposed DrivePilot. (a) An example of the multi-source representation for a real-world scene, including RAG- and CoT-enhanced annotations (14 semantic dimensions) from Qwen2-VL. (b) Statistics of DrivePilot on the word cloud of the command distributions. (c) Distribution of entity and BEV annotations per image. (d) Length distribution of original and augmented commands.

human-machine interaction in AD. Early VG methods in this field can be categorized as one-stage or two-stage. VG methods in AD were initially divided into one-stage and two-stage methods. One-stage models [38, 60, 62] are known for their efficiency, processing images and commands in a unified architecture. However, they may falter in scenarios with densely populated or overlapping regions. In contrast, two-stage models [6, 52] first generate a set of region proposals using a pre-trained detector and then match these regions against the text. While often more precise, their performance is fundamentally capped by the quality of the initial region proposals. Recent work has employed VLMs like Qwen2-VL [55], and MiniGPT-v2 [5] for their strong cross-modal semantic reasoning. However, their prohibitive computational overhead and high inference latency are incompatible with the real-time constraints of autonomous driving systems. Consequently, this progression highlights a critical and unresolved gap: the need for developing a framework that is both semantically robust for complex driving scenes and efficient enough for real-time, on-board computation.

World Model in Autonomous Driving. With the rapid progress of LLMs and modern generative models, world models have emerged as a new paradigm for model-based prediction and planning [14]. World models learn a compact latent representation of the external scenes, enabling an agent to roll out imagined future states and evaluate candidate actions before making a decision [26]. Current applications of world models in AD fall into three main categories [24]. First, for end-to-end driving [7, 20, 69], models like LAW [33] and Drive-WM [57] leverage a world model to predict future scene states, providing a robust, forward-looking context for downstream decision-making. Second, for scenario simulation [25, 49, 56], models such as DriveDreamer-2 [67] and Vista [19] utilize diffusion models or MLLMs to generate physically plausible video sequences for closed-loop simulation. Third, in representation learning [44, 68, 71], works like DriveDreamer4D [66] and FSDrive [64] use world models to couple geometry, appearance, and agent dynamics for structured 3D/4D priors. However, how world models can facilitate VG tasks for AD remains unexplored. To our knowledge, we are the first to explore the effective applica-

tion of the world model to VG tasks for autonomous vehicles, presenting preliminary studies in this emerging field.

3. DrivePilot Dataset

We introduce DrivePilot, a large-scale visual grounding dataset for AD research based on nuScenes [2]. It comprises 16,332 scenes from real-world urban environments (Singapore, Boston), capturing a wide range of driving conditions, weather, and times of day. DrivePilot is pioneering in leveraging the capabilities of Qwen2-VL [55] to generate semantic annotations. To enhance the LLM’s scene comprehension and minimize hallucinations, we employ RAG [63] to dynamically retrieve relevant external information. Furthermore, we use a zero-shot CoT prompting strategy to guide the LLM in generating step-by-step semantic annotations. The generation process is outlined in three steps:

Step-1: In-Context RAG Annotation. We first construct a knowledge base from 1,200 curated nuScenes samples, detailing agent trajectories, agent types, and road conditions. For each new scene to be annotated, we retrieve the top- k relevant scenarios from this knowledge base via cosine similarity. These retrieved examples serve as in-context cues (e.g., historical vehicle behavior in similar weather, pedestrian patterns) to guide the Qwen2-VL in generating context-aware, structured scene annotations for our dataset.

Step-2: CoT Prompting Annotation Generation. This step unfolds as a progressive dialogue, with each step directing the LLM to focus on distinct facets of the scene. As depicted in Figure 4, this CoT process involves a structured dialogue where each “thought” guides the VLM Qwen2-VL’s reasoning. The first thought focuses on understanding the overall scene, identifying key objects, and indicating their spatial relationship. The second thought analyzes command keywords and semantic intent. Subsequent thoughts prompt the Qwen2-VL to progressively consider factors such as road conditions, traffic density, notable events, and the potential behaviors of each agent. After h reasoning iterations (where h varies by sample), the insights from this chain of thought are synthesized into a cohesive semantic annotation.

Step-3: Manual Cross-check Validation. All annotations

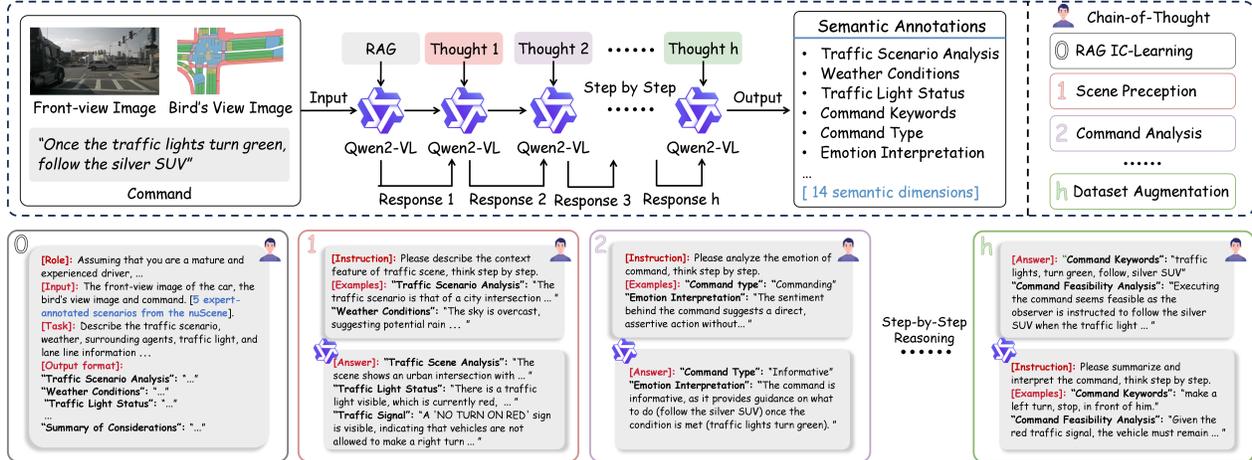


Figure 4. Illustration of the CoT prompting used in DrivePilot to generate semantic annotations for a given traffic scene. This step-by-step process involves dialogues where each “thought” guides the Qwen2-VL to understand different aspects of the scene or given command.

generated from Qwen2-VL undergo manual verification by a panel of 13 domain experts, including AV safety engineers, certified driving instructors, and postgraduate researchers. Each sample is assessed for consistency with ground-truth sensor data, adherence to traffic laws, and compliance with the European Union (EU) General Data Protection Regulation (GDPR) and the 20 AV recommendations from the EU Commission’s expert group [34]. Any discrepancies, such as misaligned object references, trigger a re-annotation process to ensure all labels meet real-world legal and operational standards. See **Appendix A** for details on these processes.

As illustrated in Figure 3, DrivePilot is a comprehensive benchmark for a range of AV tasks, including command understanding, object detection, and visual grounding. It is divided into 11,432 training samples, 2,249 validation samples, and 2,451 test samples, and provides 14 semantic dimensions for contextual richness, such as weather, traffic light status, and emotional context. Moreover, BEV images are formatted to consistently depict the target vehicle oriented to the right, with a resolution of 1200×800, while front-view images are calibrated to 1600×900. Each dataset entry is a comprehensive sample, comprising a natural language command, paired front-view and BEV images, LLM-generated scene annotations, and the precise target object location. The commands, averaging 14.72 words, are designed to challenge models with object disambiguation and interpretation of intricate textual inquiries, mirroring real-world AV navigation. By providing these diverse contexts and challenging language–vision cases, DrivePilot serves as a new standard to accelerate progress in AD grounding.

4. Methodology

The goal of this study is to create a model capable of interpreting a frontal-view image, I , and a natural language

command C . The model’s task is to pinpoint the specific area within the image I that corresponds to the destination or target object described in the command C . This requires the model to execute advanced cross-modal reasoning by seamlessly blending visual cues from I and interpretative insights from command C . In a nutshell, the model aims to determine the exact destination or object the AD is instructed to approach or identify based on the given command.

4.1. Overall Pipeline

Figure 5 shows the pipeline of ThinkDeeper, which departs from traditional ranking-based visual localization and comprises three components: (i) Multimodal Backbones, (ii) Spatial-Aware World Model, and (iii) Multimodal Decoder. Initially, the backbones extract the frontal-view image I and command C into rich vectors. Then, the SA-WM leverages these vectors to build a compact scene latent state that preserves salient visual cues while filtering background clutter, and iteratively rolls out plausible future latent states X_v to inform downstream grounding. Finally, the Multimodal Decoder fuses the predicted states with multimodal features and reasons across modalities to localize the object that best matches the natural command given by the passengers.

4.2. Multimodal Backbones

Vision Encoder. We employ a stack of Vision Transformers (ViTs) [15] to encode the image I into a visual feature map $F_v \in \mathbb{R}^{D \times H \times W}$, where D is the channel dimension and $H \times W$ matches the image resolution. In parallel, two specialist networks provide auxiliary cues: CenterNet [70] detects objects and yields an object vector $F_o \in \mathbb{R}^{N \times D}$ for the top- N objects, while ZoeDepth [1] estimates monocular depth, resulting in a depth map $F_d \in \mathbb{R}^{H \times W}$ aligned with the spatial dimensions of the visual feature map F_v .

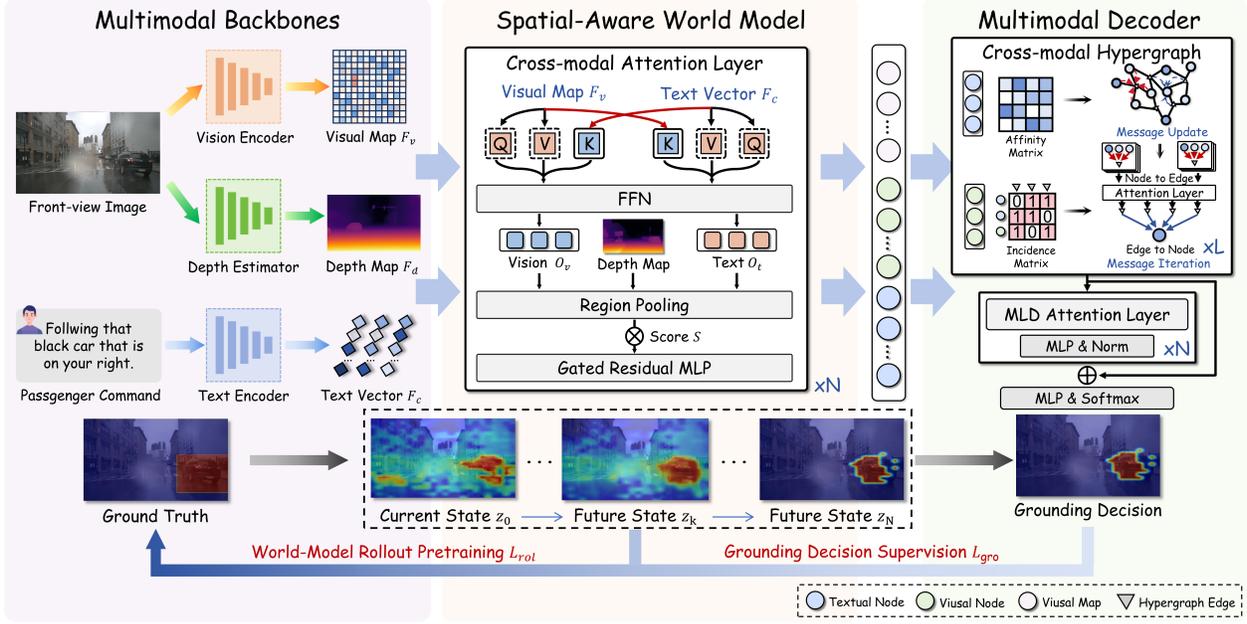


Figure 5. Overview of ThinkDeeper. Multimodal Backbones first encode the image I and command C into visual F_v , textual F_t , and depth F_d features. Next, the SA-WM takes these features and reasons ahead by rolling out a sequence of future latent states ($z_0 \rightarrow z_N$). Finally, the decoder applies a cross-modal hypergraph network to fuse these predicted states and produce the final, robust grounding decision Y .

Text Encoder. The command C is tokenized with BERT’s WordPiece tokenizer [13] and encoded by a BERT encoder, yielding text vector $F_c = \{x_1, x_2, \dots, x_L\} \in \mathbb{R}^{L \times Q}$, where L denotes the token length and Q represents the hidden size.

4.3. Spatial-Aware World Model

This component is responsible for (i) distilling the current scene into a command-aware latent state and (ii) rolling this latent forward to envisage future latent states that guide grounding. Concretely, the SA-WM operates in two phases: (1) Current State Construction and (2) Future States Rollout.

Current State Construction. As shown in Figure 5, the first phase constructs a compact latent state z_0 that represents the current scene. Given the visual feature map F_v , and the depth map F_d , a set of cross-modal attention layer projects them into a unified semantic space, producing vision-text vectors $O_t \in \mathbb{R}^{M \times L}$ and $O_v \in \mathbb{R}^{D \times H \times W}$. Formally,

$$\mathbf{A}_t = \phi_{\text{Softmax}} \left((F_v \mathbf{W}_q^{\text{vis}}) \otimes (F_c \mathbf{W}_k^t)^T / \sqrt{D} \right) \quad (1)$$

$$\mathbf{A}_v = \phi_{\text{Softmax}} \left((F_c \mathbf{W}_q^{\text{tex}}) \otimes (F_v \mathbf{W}_k^{\text{vis}})^T / \sqrt{D} \right) \quad (2)$$

$$O_t = \mathbf{A}_t^T \otimes (F_v \mathbf{W}_v^{\text{vis}}), \quad O_v = \mathbf{A}_v^T \otimes (F_c \mathbf{W}_v^{\text{tex}}) \quad (3)$$

Here, \mathbf{A}_v and \mathbf{A}_t are the affinity propagation of text-to-visual features and visual-to-text features, respectively. \otimes donates the matrix multiplication, while $\mathbf{W}_c^{\text{vis}}$ and $\mathbf{W}_v^{\text{tex}}$ are learnable parameters for F_v and F_c , respectively. To distill this rich representation into a compact state z_t and filter

out background clutter, we compute a fine-grained saliency score s^k for each visual patch k . Mathematically,

$$s^k = \frac{\sigma^2 \cdot \vec{\mathbf{a}}^T P(k)}{\exp \left(\left(1 - \sum_j F_v(k, j)^T O_v(k, j) \right)^2 / 2\mu \right)} \quad (4)$$

Here, $P(k)$ is a depth-derived prior from F_d that biases attention toward physically plausible regions, μ and σ are learnable parameters, and j is the j -th channel of dimension D . The saliency map at layer l is denoted $s^{(l)} = \{s^k\}_{k=1}^n$, and all layer-wise maps are collected as score $S = \{s^{(1)}, \dots, s^{(L)}\}$. Regions with low text–visual affinity or inconsistent depth geometry are assigned low scores and suppressed. Subsequently, we apply Region Pooling [22] over the candidate regions on score S to obtain an aggregated map \tilde{S} , which is broadcast and used to gate the object vector F_o , generating the current latent state: $z_0 = \phi_{\text{MLPs}} \left(F_o \odot \tilde{S} + F_o \right)$, where \odot is the Hadamard product and ϕ_{MLPs} is the linear projection. The resulting latent highlights command-relevant structure (objects, geometry, intent cues), while suppressing background clutter.

Future States Rollout. Grounding in dynamic scenes requires not only an accurate understanding of the current scene but also forward-looking reasoning to anticipate future developments. Accordingly, the second phase predicts a sequence of future latent states that think visually about how the scene is likely to evolve under the com-

mand. Starting from the latent state z_0 , a gated residual MLP realizes the recurrent transition f_θ in latent space: $z_{k+1} = f_\theta(\{z_k\}_{k=1}^{N-1}, O_t)$, step-by-step producing the future latent states $Z_v = \{z_1, z_2, \dots, z_N\}$. Here, O_t provides language conditioning and encodes geometric constraints. The predictions are in latent space rather than pixels, capturing prospective saliency, geometry-aware attention, and intent cues most informative for the final grounding process. The Multimodal Decoder then consumes Z_v to localize the target with spatial consistency. See **Appendix B** for details.

4.4. Multimodal Decoder

The decoder aims to fuse textual intent with the SA-WM’s future latent states X_v , reason over higher-order relations between phrases and spatial regions, and output a final grounding decision. To this end, we present a cross-modal hypergraph network to model the interplay between visual and textual modalities, guided by the future latent states, enabling multimodal fusion and localization refinement.

Hypergraph Construction. We first define the hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where the node set $\mathcal{V} = Z_v \cup X_t$ comprises N visual nodes $Z_v = \{z_1, \dots, z_N\}$ and L textual nodes $X_t = \{x_1, \dots, x_L\}$. \mathcal{E} is the hyperedge set. To couple modalities, we compute an affinity matrix A_{ij} for every vision–text pair (z_i, x_j) in a shared embedding space. We then construct hyperedges based on this matrix. For each visual node z_i , we form a corresponding hyperedge $\mathcal{E}_j \in \mathcal{E}$ (where j indexes the hyperedge) by selecting the top- k textual nodes $\{x_k\}$ with the highest affinity A_{ik} . Formally, the affinity and the resulting hyperedge feature e_j (the average of its constituent text nodes), which can be formally represented as follows:

$$A_{ij} = \left(\bar{\mathbf{a}}^T [\mathbf{W}_v z_i \| \mathbf{W}_t x_j] \right), \quad e_j = \frac{1}{|\mathcal{E}_j|} \sum_{k \in \mathcal{E}_j} x_k \quad (5)$$

where $\bar{\mathbf{a}}$ is a single-layer feed-forward projection, and \mathbf{W}_v , \mathbf{W}_t denoted the transformation matrices, respectively.

Hyperedge Weighting and Aggregation. The hypergraph topology is represented by an incidence matrix $\mathbf{H} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{E}|}$. Each element \mathbf{h}_{ij} is an attention coefficient quantifying the importance of hyperedge e_j to a given node x_i :

$$\mathbf{h}_{ij} = \frac{\exp(\text{LeakyReLU}(\bar{\mathbf{a}}^T [\mathbf{W} \mathbf{x}_i \| \mathbf{W} \mathbf{e}_j]))}{\sum_{k \in \mathcal{N}(i)} \exp(\text{LeakyReLU}(\bar{\mathbf{a}}^T [\mathbf{W} \mathbf{x}_i \| \mathbf{W} \mathbf{e}_k]))} \quad (6)$$

where $\mathcal{N}(i)$ denotes the neighborhood set of the i -th node, **Message Update and Iteration.** Inter-node communication proceeds via hypergraph convolution. Mathematically,

$$\mathbf{X}^{(l+1)} = \phi \left(\mathbf{D}_v^{-1/2} \mathbf{H} \mathbf{W}_e \mathbf{D}_e^{-1} \mathbf{H}^T \mathbf{D}_v^{-1/2} \mathbf{X}^{(l)} \Theta_g^l \right) \quad (7)$$

where \mathbf{D}_e and \mathbf{D}_v are the diagonal matrices of edge and vertex degrees. Θ_g^l is a learnable parameter at layer l , \mathbf{W}_e is the diagonal hyperedge weight matrix, and ϕ is the activation function. After processing through the hypergraph network,

the output features $\tilde{\mathbf{X}}$ are split into visual node features $\{\tilde{z}_v^i\}_{i=1}^N$ and textual node features $\{\tilde{x}_t^j\}_{j=1}^L$ for decoding.

Feature Decoding and Grounding. The updated visual and textual node features are processed by a multi-layer dynamic (MLD) attention to yield the probability distribution $P(\mathbf{Y} | \tilde{\mathbf{X}})$ over the visual nodes for the final grounding \mathbf{Y} .

$$P(\mathbf{Y} | \tilde{\mathbf{X}}) = \phi_{\text{MLD}}[\underbrace{\tilde{z}_v^1, \tilde{z}_v^2, \dots, \tilde{z}_v^N}_{\text{visual nodes}}, \underbrace{\tilde{x}_t^1, \tilde{x}_t^2, \dots, \tilde{x}_t^L}_{\text{textual nodes}}] \quad (8)$$

where the \tilde{z}_v^i and \tilde{x}_t^j are the output node feature from $\tilde{\mathbf{X}}$, respectively. See **Appendix C** for details on this decoder.

4.5. Training

We adopt a two-stage training pipeline: (1) World-Model Rollout Pretraining L_{rol} , which supervises the latent dynamics to predict future scene evolution; and (2) Grounding Decision Supervision L_{gro} , which trains the grounding module for target localization. See **Appendix D** for more details.

5. Experiments

5.1. Experimental Setup

Data Segmentation. We evaluate ThinkDeeper on three real-world datasets: Talk2Car, DrivePilot, and MoCAD [36], collectively forming the Full test set. To probe robustness, we further segment DrivePilot and MoCAD into two specialized test sets: Long-text and Corner-case. The Long-text set includes commands exceeding 23 words, which we expanded up to 50 words using Qwen2-VL to test complex linguistic handling. Correspondingly, the Corner-case set comprises three subsets: (1) 165 visual constraint scenarios (e.g., occlusions, low visibility), (2) 175 multi-agent interaction scenarios, and (3) 185 ambiguous command scenarios, designed to evaluate robustness in challenging real-world conditions. Additionally, we benchmark performance on the standard RefCOCO+/g datasets using the segmentation methodology from the classic VG model VLTVG [60].

Evaluation Metric. In accordance with the C4AV challenge, we utilize the $IoU_{0.5}$ score as the evaluation metric.

Implementation Details. ThinkDeeper is trained on $4 \times$ NVIDIA A100 GPUs, with a two-stage training process: 15 epochs in the first stage and 40 epochs in the second. We use a batch size of 32 and the AdamW optimizer with an initial learning rate of 10^{-4} . See **Appendix E** for more details.

5.2. Comparison to State-of-the-arts (SOTA)

Overall Performance. Table 2 reports the performance of ThinkDeeper against SOTA baselines. ThinkDeeper consistently outperforms all SOTA baselines across all datasets. For Talk2Car and DrivePilot, our model achieves improvements of 7.9% and 2.7% in $IoU_{0.5}$, respectively, over the best-performing baselines UNINEXT and CAVG. For the MoCAD dataset, it reduces errors by at least 3.8%, indicating the superior generalization of world model-based design.

Model	Backbone	Talk2Car	MoCAD		DrivePilot		Corner-case Test sets			Long-text
			test	val	test	val	Visual Const.	Multi-agent	Ambiguous	val
AttnGrounder [48]	ResNet-50	61.32	62.34	64.35	62.31	64.57	62.74	64.82	64.31	57.25
CMSVG [50]	EfficientNet	68.61	67.66	68.47	68.87	69.93	69.39	66.77	67.83	62.21
TransVG [11]	ResNet-101	65.83	68.14	70.85	66.52	68.42	68.12	66.34	69.25	65.45
CMRT [45]	ResNet-152	69.11	69.42	68.83	69.54	70.37	67.12	66.20	62.23	64.25
MDERT [29]	ResNet-101	70.52	66.74	70.23	71.35	72.15	68.35	65.37	68.38	62.72
VL-BERT [9]	ResNet-101	70.03	71.42	70.54	71.47	72.36	70.29	70.14	69.84	66.70
RSD-LXMERT [3]	ResNet-101	72.64	72.35	71.46	73.37	74.52	70.22	71.87	63.44	65.80
VLTVG [60]	ResNet-101	63.33	67.14	68.26	65.37	68.49	68.51	66.22	70.24	68.80
Grounding DINO [43]	ViT	68.15	67.92	68.48	69.50	70.10	66.17	65.85	67.24	63.15
UNINEXT [58]	ViT	70.87	70.62	71.34	71.35	73.47	69.26	68.78	71.29	65.32
CAVG [36]	ViT	74.62	72.44	73.25	75.52	76.48	68.39	67.36	69.45	64.36
MiniGPT-v2 [5]	Llama-2	61.15	60.89	61.72	62.85	63.27	59.14	58.33	55.41	56.78
LLaVA-NeXT (13B) [42]	LLM	42.31	43.45	44.02	43.98	44.15	41.22	40.71	37.84	39.03
Qwen2.5-VL-7B [55]	VLM	47.31	48.20	49.10	50.06	50.84	45.12	46.37	47.05	41.92
Qwen2.5-VL-72B [55]	VLM	56.17	57.10	57.85	58.92	59.74	53.43	54.25	55.17	49.83
Qwen3-VL-8B [59]	VLM	56.19	57.25	58.16	59.05	59.85	53.55	54.49	55.25	50.13
ThinkDeeper (50%)	ViT+Hyper.	68.84	70.05	73.43	70.25	72.16	66.26	68.37	72.51	71.76
ThinkDeeper	ViT+Hyper.	76.64	75.20	75.76	77.27	79.52	72.23	73.59	74.36	73.18

Table 2. Performance comparison of ThinkDeeper (marked in purple) and SOTA baselines. **Bold** values are the best performance.

Dataset	Method	Venue	val/val-g	test A/val-u	test B/test-u
RefCOCO	TransVG [11]	ICCV	81.02	82.72	78.35
	VILLA [18]	NeurIPS	81.65	87.40	74.48
	VLTVG [60]	CVPR	83.21	86.78	78.45
	SeqTR [4]	ECCV	78.22	81.47	73.80
	TransCP [53]	TPAMI	<u>84.25</u>	87.38	<u>79.78</u>
	ThinkDeeper ^s	-	75.72	79.49	77.22
	ThinkDeeper	-	85.74 ±0.4	87.78 ±0.5	80.64 ±0.2
RefCOCO+	FAQA [61]	ICCV	56.81	60.23	49.6
	VLTVG [60]	CVPR	72.36	77.21	64.8
	MPCCT [4]	PR	73.28	78.96	63.59
	TransCP [53]	TPAMI	73.07	78.05	63.35
	VILLA [18]	NeurIPS	<u>76.05</u>	<u>81.62</u>	<u>65.70</u>
	ThinkDeeper ^s	-	70.58	73.07	59.53
	ThinkDeeper	-	77.72 ±0.5	82.10 ±0.3	66.71 ±0.4
RefCOCO-g	NMTree [40]	CVPR	64.62	65.87	66.44
	VLTVG [60]	CVPR	<u>72.53</u>	74.90	73.88
	RvG-Tree [27]	TPAMI	-	66.95	66.51
	TransCP [53]	TPAMI	73.07	78.05	63.35
	ReSC-Large [62]	ECCV	63.12	67.30	67.20
	VILLA [18]	NeurIPS	-	<u>75.90</u>	<u>75.93</u>
	ThinkDeeper ^s	-	65.42	69.97	63.72
ThinkDeeper	-	72.73 ±0.2	80.90 ±0.5	77.72 ±0.4	

Table 3. Quantitative comparison of ThinkDeeper and SOTA baselines using $IoU_{0.5}$ metric. **Bold** and underlined denote the best and second-best scores. ThinkDeeper^s is the model trained on 75% of the training data. Results are averaged over ≥ 3 random seeds.

Model Robustness. ThinkDeeper excels in challenging scenarios, surpassing the next best baseline (UNINEXT) by 6.4% on the Corner-case set and 9.7% on the Long-text set. Similarly, on the Long-text set, it achieves 73.18 $IoU_{0.5}$ (+4.38 over VLTVG), highlighting its robustness in parsing extended commands, resolving linguistic ambiguity, and grounding in multi-agent scenes. Conversely, we observe that recent LLMs (MiniGPT-v2, LLaVA-NeXT, Qwen2-VL) perform poorly on these grounding tasks, with scores lagging 15-25 points behind SOTA. This is expected, as these general-

Components	Ablation Methods				
	A	B	C	D	E
Vision Encoder	✓	✓	✓	✓	✓
SA-WM (-Future)	✓	✗	✓	✓	✓
SA-WM	✓	✓	✗	✓	✓
Cross-modal Decoder	✓	✓	✓	✗	✓
$IoU_{0.5}$	72.33	68.27	62.70	70.42	77.27

Table 4. Ablation studies for core component in DrivePilot.

Method	Backbone	Param.	Inference Time	$IoU_{0.5}$
VLTVG [60]	ResNet-101	152.18M	55ms	69.72
UNITER [8]	ResNet-101	112.00M	58ms	73.42
RSDLayerAttn [3]	ResNet-101	112.28M	54ms	74.12
Grounding DINO [43]	ViT	210.32M	88ms	68.15
CAVG [36]	ViT	172.78M	69ms	74.50
ThinkDeeper	ViT	135.81M	39ms	76.64

Table 5. Efficiency comparison of ThinkDeeper and baselines on Talk2Car, benchmarked using an NVIDIA A40 (40GB) GPU.

purpose models lack the inductive biases for high-precision localization. Their inability to exploit 3D spatial context and depth cues to resolve ambiguity underscores the need for our depth-aware and spatial-aware world-model design tailored to real-world autonomous driving development.

Model Generalizability. We further evaluate ThinkDeeper on RefCOCO+/g benchmarks. As reported in Table 3, it achieves impressive performance across all nine evaluation splits, surpassing the SOTA baselines by at least 2.9%, 3.0%, and 3.5%, respectively. These results illustrate that although ThinkDeeper is developed for noisy, safety-critical on-road scenes, prioritizing lightweight real-time operation, it generalizes strongly to high-quality, densely annotated benchmarks, showcasing robust cross-domain applicability.

Data Efficiency. When trained on only 50% or 75% of the dataset, ThinkDeeper surprisingly beat most baselines

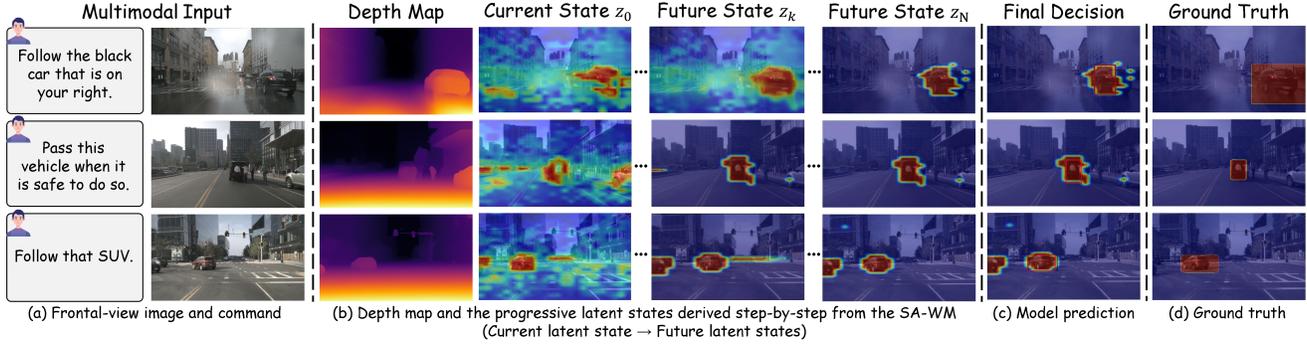


Figure 6. Qualitative results of depth map, current/future latent states, and model performance on the DrivePilot dataset.

trained on full datasets across all test sets. This shows its efficiency and scalability, reducing training data requirements while maintaining high performance in corner-case scenes.

5.3. Comparison of Model Complex and Efficiency

Table 5 reports the inference efficiency of ThinkDeeper across 2,048 randomly selected scenes from Talk2Car, evaluated on an NVIDIA A40 GPU (70 TOPS). Despite utilizing fewer parameters than CAVG and Grounding DINO, it achieves the highest $IoU_{0.5}$ (78.64), while maintaining a competitive average inference time of 39ms per sample, outperforming most SOTA baselines in both accuracy and efficiency. This performance meets the computational requirements for Level 3 AD systems (20-30 TOPS), enabling seamless deployment on various on-board Neural Processing Units (NPU) and Data Processing Units (DPU), such as Tesla FSD (245 TOPS) and NVIDIA Thor-U (500 TOPS).

5.4. Ablation Study

Table 4 showcases the ablation study for each core component in ThinkDeeper. Our full model (Method E) achieves the highest score (77.27 $IoU_{0.5}$) and serves as the reference. Specifically, Method A removes depth-derived prior from the Vision Encoder, causing a 6.3% $IoU_{0.5}$ drop, underscoring the significance of depth information in refining spatial awareness. Method B omits the Future States Rollout, restricting the model to only the current latent state z_0 . This incurs a severe 11.7% performance drop, validating that static reasoning is insufficient and forward-looking inference is essential for resolving spatial ambiguity. Method C removes the entire SA-WM module, revealing a catastrophic 14.57 $IoU_{0.5}$ collapse. This confirms that our world model’s ability to distill a command-aware current latent state (z_0) and project future latent states ($\{z_k\}_{k=1}^N$) is a fundamental intermediate cue for robust grounding decisions. Finally, Method D replaces our hypergraph with a standard GCN, which also suffers notable performance degradation. This highlights the hypergraph’s superior ability to encode higher-order visual-textual and multi-agent relations compared with

pairwise graph updates. Finally, replacing our hypergraph with a standard GCN in Method D suffers a notable 6.85 $IoU_{0.5}$ degradation ($\downarrow 8.9\%$), highlighting the hypergraph’s superior ability to encode higher-order visual-textual relations over pairwise graph updates.

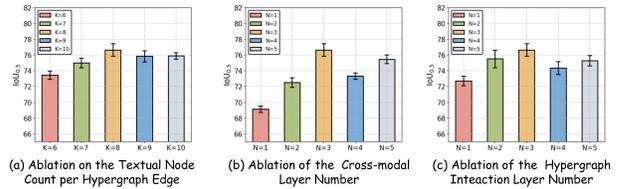


Figure 7. Ablation results of ThinkDeeper’s hyperparameter.

5.5. Hyperparameter Sensitivity Analysis

Figure 7 reports our sensitivity analysis on the Talk2Car. We find that performance steadily improves as we increase the number of cross-modal attention layers, enabling richer interaction between the textual command and the SA-WM’s latent states. The gains saturate at $N = 3$, achieving our peak score of 76.64. Similarly, performance peaks when each visual node’s hyperedge connects to $K = 6$ textual nodes, indicating this provides an optimal balance of contextual information without introducing noise from irrelevant phrases. These results suggest that a moderately deep (3-layer) and broad (6-node) decoder is ideal for reasoning over our proposed world model’s envisioned future states, while overly complex architectures yield diminishing returns.

5.6. Visualization and Interpretability Analysis

Figure 6 reveals how our model “thinks deeper” by visualizing the latent states generated by the SA-WM in challenging scenes. These states serve as intermediate evidence for the final grounding decision. The SA-WM first distills the scene into a command-aware current latent state z_0 , which encodes relevant objects, geometry, and saliency. It then rolls out future latent states $\{z_k\}_{k=1}^N$ that anticipate prospective changes

in saliency and interaction patterns. This distilled, forward-looking representation provides the Decoder with a set of high-quality, filtered cues, enabling it to resolve ambiguity. For example, in the third row of Figure 6(b) (command: “Follow that SUV.”), the model’s initial attention is diffuse. However, as the SA-WM reasons about the “follow” and “SUV” concepts over its future states, it successfully isolates the target vehicle from background clutter. This internal filtering process directly enables the decoder to lock onto the correct SUV. These results support our quantitative findings: by treating current and imagined future latents as structured intermediate evidence, ThinkDeeper achieves robust grounding in low-light, occluded, crowded, and ambiguous scenes common in on-road applications. We have included additional quantitative results of ThinkDeeper in **Appendix F**.

6. Conclusion

This paper introduces ThinkDeeper, the first world model-based framework for visual grounding in AD. By leveraging imagined future states as intermediate reasoning steps, ThinkDeeper effectively bridges the visual-linguistic gap, enabling more accurate and efficient grounding in real-world scenes. Importantly, we present DrivePilot, a multi-source dataset with extensive LLM-generated annotations, providing a challenging benchmark for VG research. ThinkDeeper achieves SOTA performance, ranking #1 on the C4AV leaderboard for Talk2Car, DrivePilot, and MoCAD, while also showcasing impressive performance on the RefCOCO+/g datasets, highlighting the potential of world models for robust cross-modal reasoning in fully autonomous driving.

Acknowledgements

This work was supported by the Science and Technology Development Fund of Macau [0122/2024/RIB2, 0215/2024/AGJ, 001/2024/SKL], the Research Services and Knowledge Transfer Office, University of Macau [SRG2023-00037-IOTSC, MYRG-GRG2024-00284-IOTSC], National Natural Science Foundation of China [Grants 52572354], the Shenzhen-Hong Kong-Macau Science and Technology Program Category C [SGDX20230821095159012], the State Key Lab of Intelligent Transportation System [2024-B001], and the Jiangsu Provincial Science and Technology Program [BZ2024055].

References

- [1] Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023. 4
- [2] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multi-modal dataset for autonomous driving. In *Proceedings of the IEEE/CVF CVPR*, pages 11621–11631, 2020. 3, 5
- [3] Hou Pong Chan, Mingxi Guo, and Cheng-Zhong Xu. Grounding commands for autonomous vehicles via layer fusion with region-specific dynamic layer attention. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 12464–12470. IEEE, 2022. 7
- [4] Chongqing Chen, Dezhi Han, and Chin-Chen Chang. Mpcct: Multimodal vision-language learning paradigm with context-based compact transformer. *Pattern Recognition*, 147:110084, 2024. 1, 7
- [5] Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechu Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023. 1, 3, 7
- [6] Long Chen, Wenbo Ma, Jun Xiao, Hanwang Zhang, and Shih-Fu Chang. Ref-nms: Breaking proposal bottlenecks in two-stage referring expression grounding. In *Proceedings of the AAAI conference on artificial intelligence*, pages 1036–1044, 2021. 3
- [7] Yuntao Chen, Yuqi Wang, and Zhaoxiang Zhang. Drivinggpt: Unifying driving world modeling and planning with multimodal autoregressive transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 26890–26900, 2025. 3
- [8] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *European conference on computer vision*, pages 104–120, 2020. 7
- [9] Hang Dai, Shujie Luo, Yong Ding, and Ling Shao. Commands for autonomous vehicles by progressively stacking visual-linguistic representations. In *Computer Vision–ECCV Workshops*, pages 27–32, 2020. 7
- [10] Ming Dai, Lingfeng Yang, Yihao Xu, Zhenhua Feng, and Wankou Yang. Simvg: A simple framework for visual grounding with decoupled multi-modal fusion. *Advances in neural information processing systems*, 37:121670–121698, 2024. 1
- [11] Jiajun Deng, Zhengyuan Yang, Tianlang Chen, Wengang Zhou, and Houqiang Li. Transvg: End-to-end visual grounding with transformers. In *Proceedings of the IEEE/CVF ICCV*, pages 1769–1779, 2021. 7
- [12] Thierry Deruyttere, Simon Vandenhende, Dusan Grujicic, Luc Van Gool, and Marie-Francine Moens. Talk2car: Taking control of your self-driving car. *arXiv preprint arXiv:1909.10838*, 2019. 2, 5
- [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 5
- [14] Jingtao Ding, Yunke Zhang, Yu Shang, Yuheng Zhang, Zefang Zong, Jie Feng, Yuan Yuan, Hongyuan Su, Nian Li, Nicholas Sukiennik, et al. Understanding world or predicting future? a comprehensive survey of world models. *ACM Computing Surveys*, 58(3):1–38, 2025. 3
- [15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Syl

- vain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 4
- [16] Jessica López Espejel, El Hassane Ettifouri, Mahaman Sanoussi Yahaya Alassan, El Mehdi Chouham, and Walid Dahhane. Gpt-3.5, gpt-4, or bard? evaluating llms reasoning ability in zero-shot setting and performance boosting through prompts. *Natural Language Processing Journal*, 5:100032, 2023. 1
- [17] Anja K Faulhaber, Anke Dittmer, Felix Blind, Maximilian A Wächter, Silja Timm, Leon R Sütfeld, Achim Stephan, Gordon Pipa, and Peter König. Human decisions in moral dilemmas are largely described by utilitarianism: Virtual car driving study provides guidelines for autonomous driving vehicles. *Science and engineering ethics*, 25:399–418, 2019. 1
- [18] Zhe Gan, Yen-Chun Chen, Linjie Li, Chen Zhu, Yu Cheng, and Jingjing Liu. Large-scale adversarial training for vision-and-language representation learning. *Nips*, pages 6616–6628, 2020. 7
- [19] Shenyan Gao, Jiazhi Yang, Li Chen, Kashyap Chitta, Yihang Qiu, Andreas Geiger, Jun Zhang, and Hongyang Li. Vista: A generalizable driving world model with high fidelity and versatile controllability. *Advances in Neural Information Processing Systems*, 37:91560–91596, 2024. 3
- [20] Zeyu Gao, Yao Mu, Chen Chen, Jingliang Duan, Ping Luo, Yanfeng Lu, and Shengbo Eben Li. Enhance sample efficiency and robustness of end-to-end urban autonomous driving via semantic masked world model. *IEEE Transactions on Intelligent Transportation Systems*, 25(10):13067–13079, 2024. 3
- [21] Liang Geng, Jianqin Yin, Gang Chen, and Qingxuan Jia. Pseudo-ev: Enhancing 3d visual grounding with pseudo embodied viewpoint. *IEEE Transactions on Circuits and Systems for Video Technology*, 2025. 1
- [22] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015. 5
- [23] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010. 6
- [24] Yanchen Guan, Haicheng Liao, Zhenning Li, Jia Hu, Runze Yuan, Yunjian Li, Guohui Zhang, and Chengzhong Xu. World models for autonomous driving: An initial survey. *IEEE Transactions on Intelligent Vehicles*, 2024. 3
- [25] Yanchen Guan, Haicheng Liao, Chengyue Wang, Xingcheng Liu, Jiaxun Zhang, and Zhenning Li. World model-based end-to-end scene generation for accident anticipation in autonomous driving. *Communications Engineering*, 4(1):144, 2025. 3
- [26] David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. *Advances in neural information processing systems*, 31, 2018. 3
- [27] Richang Hong, Daqing Liu, Xiaoyu Mo, Xiangnan He, and Hanwang Zhang. Learning to compose and reason with language tree structures for visual grounding. *IEEE TPAMI*, pages 684–696, 2019. 7
- [28] Haojun Jiang, Yuanze Lin, Dongchen Han, Shiji Song, and Gao Huang. Pseudo-q: Generating pseudo language queries for visual grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15513–15523, 2022. 1
- [29] Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *Proceedings of the IEEE/CVF ICCV*, pages 1780–1790, 2021. 7
- [30] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 787–798, 2014. 2
- [31] Sebastian Krügel and Matthias Uhl. Autonomous vehicles and moral judgments under risk. *Transportation research part A: policy and practice*, 155:1–10, 2022. 1
- [32] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022. 6
- [33] Yingyan Li, Lue Fan, Jiawei He, Yuqi Wang, Yuntao Chen, Zhaoxiang Zhang, and Tieniu Tan. Enhancing end-to-end autonomous driving with latent world model. *arXiv preprint arXiv:2406.08481*, 2024. 3
- [34] Zhenning Li et al. Steering the future: Redefining intelligent transportation systems with foundation models. *CHAIN*, 1(1): 46–53, 2024. 4, 3
- [35] Haicheng Liao, Yongkang Li, Chengyue Wang, Yanchen Guan, Kahou Tam, Chunlin Tian, Li Li, Chengzhong Xu, and Zhenning Li. When, where, and what? a benchmark for accident anticipation and localization with large language models. In *ACM International Conference on Multimedia (ACM MM), Oral Presentation*, pages 8–17, 2024. 1
- [36] Haicheng Liao, Huanming Shen, Zhenning Li, Chengyue Wang, Guofa Li, Yiming Bie, and Chengzhong Xu. Gpt-4 enhanced multimodal grounding for autonomous driving: Leveraging cross-modal attention with large language models. *Communications in Transportation Research*, 4:100116, 2024. 1, 6, 7
- [37] Haicheng Liao, Hanlin Kong, Bonan Wang, Chengyue Wang, Wang Ye, Zhengbing He, Chengzhong Xu, and Zhenning Li. Cot-drive: Efficient motion forecasting for autonomous driving with llms and chain-of-thought prompting. *IEEE Transactions on Artificial Intelligence*, 2025. 2, 1
- [38] Yue Liao, Si Liu, Guanbin Li, Fei Wang, Yanjie Chen, Chen Qian, and Bo Li. A real-time cross-modality correlation filtering method for referring expression comprehension. In *Proceedings of the IEEE/CVF CVPR*, pages 10880–10889, 2020. 3
- [39] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 5
- [40] Daqing Liu, Hanwang Zhang, Feng Wu, and Zheng-Jun Zha. Learning to assemble neural module tree networks for visual

- grounding. In *Proceedings of the IEEE/CVF ICCV*, pages 4673–4682, 2019. 7
- [41] Fang Liu, Yuhao Liu, Yuqiu Kong, Ke Xu, Lihe Zhang, Bao-cui Yin, Gerhard Hancke, and Rynson Lau. Referring image segmentation using text supervision. In *Proceedings of the IEEE/CVF ICCV*, pages 22124–22134, 2023. 2
- [42] 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. 2, 7
- [43] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 1, 7
- [44] Yifan Lu, Xuanchi Ren, Jiawei Yang, Tianchang Shen, Zhangjie Wu, Jun Gao, Yue Wang, Siheng Chen, Mike Chen, Sanja Fidler, et al. Infinitcube: Unbounded and controllable dynamic 3d driving scene generation with world-guided video models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 27272–27283, 2025. 3
- [45] Shujie Luo, Hang Dai, Ling Shao, and Yong Ding. C4av: learning cross-modal representations from transformers. In *Computer Vision–ECCV 2020*, pages 33–38, 2020. 7
- [46] Yunsheng Ma, Wenqian Ye, Can Cui, Haiming Zhang, Shuo Xing, Fucui Ke, Jinhong Wang, Chenglin Miao, Jintai Chen, Hamid Rezaatofighi, et al. Position: Prospective of autonomous driving—multimodal LLMs world models embodied intelligence AI alignment and mamba. In *Proceedings of the Winter Conference on Applications of Computer Vision (WACV)*, pages 1010–1026, 2025. 2
- [47] Mayug Maniparambil, Chris Vorster, Derek Molloy, Noel Murphy, Kevin McGuinness, and Noel E O’Connor. Enhancing clip with gpt-4: Harnessing visual descriptions as prompts. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 262–271, 2023. 1
- [48] Vivek Mittal. Attngrounder: Talking to cars with attention. In *Computer Vision–ECCV Workshops*, pages 62–73, 2020. 7
- [49] Chaojun Ni, Guosheng Zhao, Xiaofeng Wang, Zheng Zhu, Wenkang Qin, Guan Huang, Chen Liu, Yuyin Chen, Yida Wang, Xueyang Zhang, et al. Recondreamer: Crafting world models for driving scene reconstruction via online restoration. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 1559–1569, 2025. 3
- [50] Nivedita Rufus, Unni Krishnan R Nair, K Madhava Krishna, and Vineet Gandhi. Cosine meets softmax: A tough-to-beat baseline for visual grounding. In *Computer Vision–ECCV Workshops*, pages 39–50, 2020. 7
- [51] Seyed Sadegh Mohseni Salehi, Deniz Erdogmus, and Ali Gholipour. Tversky loss function for image segmentation using 3d fully convolutional deep networks. In *International workshop on machine learning in medical imaging*, pages 379–387. Springer, 2017. 5
- [52] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*, 2019. 3
- [53] Wei Tang, Liang Li, Xuejing Liu, Lu Jin, Jinhui Tang, and Zechao Li. Context disentangling and prototype inheriting for robust visual grounding. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(5):3213–3229, 2024. 1, 7
- [54] Chengyue Wang, Haicheng Liao, Zhenning Li, and Chengzhong Xu. Wake: Towards robust and physically feasible trajectory prediction for autonomous vehicles with wavelet and kinematics synergy. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025. 1
- [55] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 2, 3, 7
- [56] Xiaofeng Wang, Zheng Zhu, Guan Huang, Xinze Chen, Jiayang Zhu, and Jiwen Lu. Drivedreamer: Towards real-world-drive world models for autonomous driving. In *European conference on computer vision*, pages 55–72. Springer, 2024. 3
- [57] 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14749–14759, 2024. 3
- [58] Bin Yan, Yi Jiang, Jiannan Wu, Dong Wang, Ping Luo, Zehuan Yuan, and Huchuan Lu. Universal instance perception as object discovery and retrieval. In *Proceedings of the IEEE/CVF CVPR*, pages 15325–15336, 2023. 7
- [59] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025. 7
- [60] Li Yang, Yan Xu, Chunfeng Yuan, Wei Liu, Bing Li, and Weiming Hu. Improving visual grounding with visual-linguistic verification and iterative reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9499–9508, 2022. 2, 3, 6, 7
- [61] Zhengyuan Yang, Boqing Gong, Liwei Wang, Wenbing Huang, Dong Yu, and Jiebo Luo. A fast and accurate one-stage approach to visual grounding. In *Proceedings of the IEEE/CVF ICCV*, pages 4683–4693, 2019. 7
- [62] Zhengyuan Yang, Tianlang Chen, Liwei Wang, and Jiebo Luo. Improving one-stage visual grounding by recursive subquery construction. In *Computer Vision–ECCV 2020*, pages 387–404, 2020. 3, 7
- [63] Jianhao Yuan, Shuyang Sun, Daniel Omeiza, Bo Zhao, Paul Newman, Lars Kunze, and Matthew Gadd. Rag-driver: Generalisable driving explanations with retrieval-augmented in-context learning in multi-modal large language model. *arXiv preprint arXiv:2402.10828*, 2024. 3
- [64] Shuang Zeng, Xinyuan Chang, Mengwei Xie, Xinran Liu, Yifan Bai, Zheng Pan, Mu Xu, and Xing Wei. Futuresightdrive: Thinking visually with spatio-temporal cot for autonomous driving. *arXiv preprint arXiv:2505.17685*, 2025. 3
- [65] Yang Zhan, Yuan Yuan, and Zhitong Xiong. Mono3dvg: 3d visual grounding in monocular images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 6988–6996, 2024. 1

- [66] Guosheng Zhao, Chaojun Ni, Xiaofeng Wang, Zheng Zhu, Xueyang Zhang, Yida Wang, Guan Huang, Xinze Chen, Boyuan Wang, Youyi Zhang, et al. Drivedreamer4d: World models are effective data machines for 4d driving scene representation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 12015–12026, 2025. [3](#)
- [67] Guosheng Zhao, Xiaofeng Wang, Zheng Zhu, Xinze Chen, Guan Huang, Xiaoyi Bao, and Xingang Wang. Drivedreamer-2: Llm-enhanced world models for diverse driving video generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 10412–10420, 2025. [3](#)
- [68] Wenzhao Zheng, Weiliang Chen, Yuanhui Huang, Borui Zhang, Yueqi Duan, and Jiwen Lu. Occworld: Learning a 3d occupancy world model for autonomous driving. In *European conference on computer vision*, pages 55–72. Springer, 2024. [3](#)
- [69] Yupeng Zheng, Pengxuan Yang, Zebin Xing, Qichao Zhang, Yuhang Zheng, Yinfeng Gao, Pengfei Li, Teng Zhang, Zhongpu Xia, Peng Jia, et al. World4drive: End-to-end autonomous driving via intention-aware physical latent world model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 28632–28642, 2025. [3](#)
- [70] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019. [4](#)
- [71] Sicheng Zuo, Wenzhao Zheng, Yuanhui Huang, Jie Zhou, and Jiwen Lu. Gaussianworld: Gaussian world model for streaming 3d occupancy prediction. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 6772–6781, 2025. [3](#)

Appendix

A. DrivePilot Dataset

A.1. Step-1: In-Context RAG Annotation

To enhance LLM reasoning with real-world driving knowledge, we implement a two-tier Retrieval-Augmented Generation (RAG) framework. This process grounds AV visual grounding annotations in empirical driving data, ensuring contextual accuracy and reduced hallucination rates. We curate a comprehensive, multimodal knowledge base, including 1,200 expert-annotated scenarios from the nuScene Dataset, covering agent trajectories, road topology, and traffic rule compliance. For each query pair (input image and command), we execute a three-phase retrieval process:

Feature Encoding. A pre-trained vision backbone (Fast R-CNN) extracts dense visual embeddings from the input image, capturing spatial relationships and object semantics. Simultaneously, a BERT-based language model encodes textual features from the given command, ensuring a rich semantic representation for cross-modal alignment.

Cross-Modal Retrieval. The extracted embeddings are used to retrieve the top- k ($k = 5$) most relevant scenes from the knowledge base via cosine similarity. As shown in Figure A3, the retrieved samples include both human-annotated metadata (weather conditions, agent behavior) and raw sensor data from the traffic reports and nuScene dataset, enhancing scene understanding and context recall.

Knowledge Infusion. We template the retrieved scenarios into a structured prompt, instructing the LLM to ground its reasoning in both common driving behaviors, such as yielding to pedestrians in rain, and traffic regulations like interpreting ambiguous signals. This retrieval-augmented prompting ensures that the model’s decisions are contextually grounded, enhancing both situational awareness and the trustworthiness of its reasoning.

A.2. Step-2: CoT Annotation Generation

LLMs like Qwen and LLaVA excel in natural language understanding but are not inherently trained for AD or VG tasks. Prior studies [16, 37, 47] have shown that structured prompt engineering can significantly improve LLMs’ zero-shot visual description performance. To leverage this potential, we develop a progressive Chain-of-Thought (CoT) prompting strategy for generating context-aware semantic annotations, enabling dataset augmentation and semantic enrichment without fine-tuning. Specifically, we explore the utility of Qwen in augmenting and refining existing multimodal data (frontal and BEV images, paired with natural language commands) using few-shot or zero-shot prompting techniques. This enables dataset expansion and improved semantic annotation generation without costly and time-consuming fine-tuning. This CoT process includes:

<p>Question/Prompt: Assuming that you are a mature and experienced driver, please talk about what information you got about driving based on the front view and top-down view of the car given. Please describe the observable traffic scenario, e.g., "The parking lot is busy with trucks."</p> <p>Qwen2-VL Answer: The road is wet and likely slippery due to rain, with multiple vehicles on the road, some with headlights on, indicating low visibility conditions. It is raining, as evidenced by the wet road surface and raindrops on the camera lens.</p> <p>Question/Prompt: Here is an explanation of sentiment analysis: 1.Urgent: Urgent commands are associated with potential traffic hazards and demand immediate action to ensure safety. These include situations where there might be pedestrians in or near the vehicle’s path, circumstances that necessitate sudden braking, or scenarios that require evasive maneuvers to avoid collisions or other dangers. 2.Commanding: Direct and non-negotiable instructions. These are authoritative but less intense than urgent instructions. 3.Informative: Neutral instructions that convey useful information without a strong emotional tone. Informative commands are primarily characterized by their objective nature, focusing on providing factual details about a particular subject or object. These commands tend to be descriptive, where the content predominantly revolves around elucidating specific attributes or features of an entity without conveying a pressing need for action. Please Interpret the sentiment behind the command using the keywords.</p> <p>Qwen2-VL Answer: Commanding. The command suggests a desire for efficiency and convenience, particularly due to bad weather, but does not convey an immediate safety concern.</p>	 <p>Frontal-view Image</p>  <p>Bird's-view Image</p> <p>"After this car on our left there should be an entrance to a parking lot. Take it. I want to be as close as possible to my work with this bad weather."</p> <p>Command</p>
---	--

Figure A1. An example of the Qwen2-VL leveraging lens-less cueing technology to interpret driving scenarios and generate driving maneuvers, eliminating the need for specialized fine-tuning.

Scene Semantic Enhancement. Beyond object-level annotations, we introduce Scene Semantic Enhancement to enrich high-level contextual information. As shown in Table A1 and Figure A1, we employ CoT prompting, guiding Qwen to generate 14 categories of contextual metadata, including:

- Scene Descriptions: Summarizing overall environmental context, including road surface conditions, weather, lighting, and potential hazards in the traffic scenes.
- Emotion Interpretation: Analyzing pedestrian and driver intent to predict potential interaction risks.
- Road Condition Summaries: Detailing surface quality, lane markings, and obstacles affecting navigation.
- Traffic Signal & Sign Interpretations: Detailing surface quality, lane markings, and obstacles.

To diversify training data, we randomly replace 30% of the original command texts with this CoT-generated augmentation text during training, ensuring diverse linguistic patterns and robust generalization. Furthermore, keyword-based augmentation is introduced to append relevant context cues to the original command prompts, aiding semantic disambiguation. For example, semantic keywords such as “low visibility” and “intersection” are added as auxiliary hints to commands in scenarios with obstructed vision or Multi-agent. Notably, this step-by-step process involves dialogues where each “thought” guides Qwen2-VL to understand different aspects of the scene or command. The first thought fo-

Image Annotations

Visual input

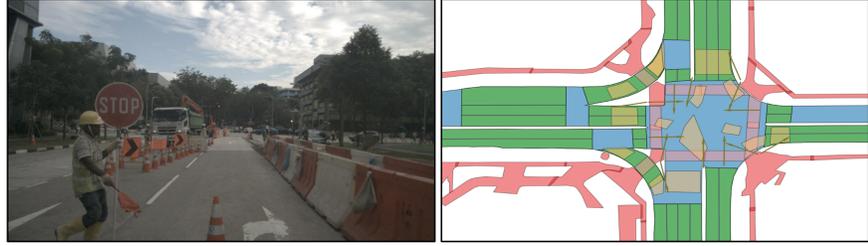


Figure A2. An example of how a large language model interprets visual information: the left image represents the front view, while the right image corresponds to the BEV image.

Prompt

You are an AI visual assistant tasked with providing semantic enhancement for visual grounding tasks. Given an image containing multiple bounding boxes and its corresponding BEV (Bird's-Eye View) map, your role is to annotate each region by supplementing and describing the objects within it. In the BEV map, green areas represent drivable regions, while red areas indicate non-drivable regions. You should concisely summarize the behavior of objects in each region and their interactions with the surrounding environment, then describe the specific content of the provided bounding boxes, including object types, actions, positions, and relative locations between objects. Below are the coordinates for each entity object in the forward view and BEV map, along with the corresponding annotation. Example:

- **Entities Annotation 1:** (label 1: [x1,y1,x2,y2]) Description of the features in this scene region.
- **BEV Annotation 1:** (label 1: [x1,y1,x2,y2]) Description of the spatial positioning in this region.

Entities Bounding box

Person:[0.116,0.45,0.314,0.998], Truck[0.322,0.403,0.488,0.609], Sign: [0.202,0.338,0.329,0.553], Car [0.685,0.541,0.756,0.630], etc.

BEV coordinates

Entities BEV coordinates, Drivable area coordinates, Non-drivable area coordinates.

LLM Annotations

Entities Annotation 1: (Person: [0.116,0.45,0.314,0.998]) A worker in safety gear is holding a stop sign, managing traffic at a road construction site.

Entities Annotation 2: (Truck: [0.202,0.338,0.329,0.553]) A truck is involved in construction activities or waiting to proceed, on the left side of the road ahead.

Entities Annotation 3: (Stop Sign: [0.202,0.338,0.329,0.553]) A warning sign being held by a worker crossing the road.

Entities Annotation 4: (Car: [0.685,0.541,0.756,0.630]) A black sedan approaching head-on in the right lane.

...

BEV Annotation 4: (Car: [0.685,0.541,0.756,0.630]) A sedan positioned at the intersection behind a truck.

Table A1. Illustration of LLM-driven scene understanding. Given paired front-view images and BEV maps, the LLM is prompted with region definitions, bounding boxes, and BEV coordinates to produce rich semantic annotations. The prompt instructs the model to summarize scene context, describe object behaviors, and provide detailed entity-level grounding (types, actions, positions, and relations). By integrating 2D visual data with 3D spatial coordinates (BEV), the model generates fine-grained semantic descriptions. This demonstrates the LLM's ability to reason about object behaviors and interactions within a dynamic traffic environment.

cuses on understanding the scene and identifying key objects and their dynamics. The second thought analyzes command keywords and emotions. After h iterations of reasoning and updating (with iterations varying per sample for the actual situation), insights from each thought are synthesized into a cohesive semantic scene annotation.

A.3. Step-3: Manual Cross-check Validation

To ensure accuracy, reliability, and compliance, all LLM-generated annotations undergo rigorous manual cross-checking by 13 domain experts, including AV safety engineers, certified driving instructors, and researchers. The validation process follows a multi-stage review pipeline:

Consistency Verification. Annotations are checked against ground-truth LiDAR, radar, and camera sensor data from nuScenes to validate spatial accuracy. For example, object positions are checked against 3D bounding box coordinates to verify the accuracy of object locations, spatial relationships, and motion dynamics, while temporal consistency checks are performed for dynamic objects, ensuring that annotations align with actual vehicle or pedestrian movement trajectories across frames. Moreover, objects identified in BEV annotations are manually compared to their real-world positioning to eliminate false positives, resolve spatial ambiguities, and ensure correct depth perception.

Compliance with Traffic and Safety Regulations. Each annotation undergoes validation against standardized traffic regulations, including compliance with EU right-of-way regulations, lane discipline and right-of-way rules, traffic light and stop sign recognition, speed limits, braking distance considerations, and pedestrian and cyclist priority laws. Mislabeling of objects or failure to recognize critical traffic elements, like crosswalk violations and illegal lane changes, triggers reannotation to meet legal standards.

Regulatory Compliance and Ethical Considerations. Each annotation is audited for adherence to traffic laws (EU right-of-way regulations) and ethical guidelines. Legal experts ensure compliance with the General Data Protection Regulation (GDPR) (anonymizing license plates/faces via pixelation) and the EU Commission’s 20 AV recommendations [34], such as avoiding high-risk phrasing, like “accelerate through a yellow light”, ensuring that safety-critical decisions made by AV models reflect ethically responsible driving behavior. In addition, safety engineers further validate scene dynamics against ISO 26262 functional safety standards, flagging ambiguous scenarios without road conditions for reannotation.

Discrepancy Handling and Iterative Reannotation. Any detected misalignments, incomplete labels, or logical inconsistencies (e.g., an object classified as static but marked with movement vectors) are flagged for correction. Then, a tiered reannotation process is employed, where flagged samples are re-evaluated by at least two independent reviewers to minimize human bias and maintain annotation consistency.

In cases of annotation ambiguity, a consensus-based verification is conducted, where multiple reviewers deliberate and adjust labels. Finally, the Qwen2-VL regenerates annotations using revised prompts incorporating expert feedback, ensuring alignment with driving norms.

Final Validation and Benchmarking. Post-correction, a final review pass ensures that all annotations adhere to quality benchmarks and align with industry-grade AV datasets. Moreover, generates statistical reports through user surveys to analyze annotation consistency rates, discrepancy resolution efficiency, and regulatory compliance metrics to inform continuous improvement of the dataset validation pipeline.

B. Spatial-Aware World Model

The calculation of depth-derived prior $P(k)$ is derived from the depth graph F_d , which encodes the depth information of the input visual data. To ensure a consistent depth representation, F_d is normalized to the range $[0, 1]$ using an Exponential Decay Function. This function assigns higher values to closer objects while attenuating the influence of distant regions, ensuring depth-aware feature refinement. The transformation can be formally expressed as follows:

$$F_D^{\text{nor}}(x) = \exp(-\alpha \cdot F_D(x)) \quad (9)$$

where α is a decay rate hyperparameter that regulates depth sensitivity, preserving finer details for nearby objects while suppressing distant regions. Pixels corresponding to objects at infinite depth are set to zero, excluding them from further processing to improve computational efficiency.

Next, the normalized depth map F_D^{nor} is passed through a Multi-Layer Perceptron (MLP), which employs a piecewise activation function. This allows the model to adaptively emphasize depth regions based on their visual importance, refining spatial awareness in visual grounding. Formally,

$$P(x) = \phi_{\text{MLP}}(F_D^{\text{nor}}(x)) \quad (10)$$

where ϕ_{MLP} represents the MLP transformation, mapping depth-normalized features for downstream tasks.

C. Multi-Layer Dynamic Attention Mechanism

Traditional approaches to traffic scene analysis primarily focus on localized and top-level feature extraction, often overlooking the intricate interplay between scene semantics and textual cues. This limitation hinders a model’s ability to comprehend contextual relationships in complex driving environments. Our study highlights the importance of integrating semantic and contextual information in traffic scene analysis beyond simple detection through fusion vectors. To address this, we propose an MLD attention mechanism embedded within the Multimodal Decoder to capture dynamic textual-semantic interactions. Our approach enables precise

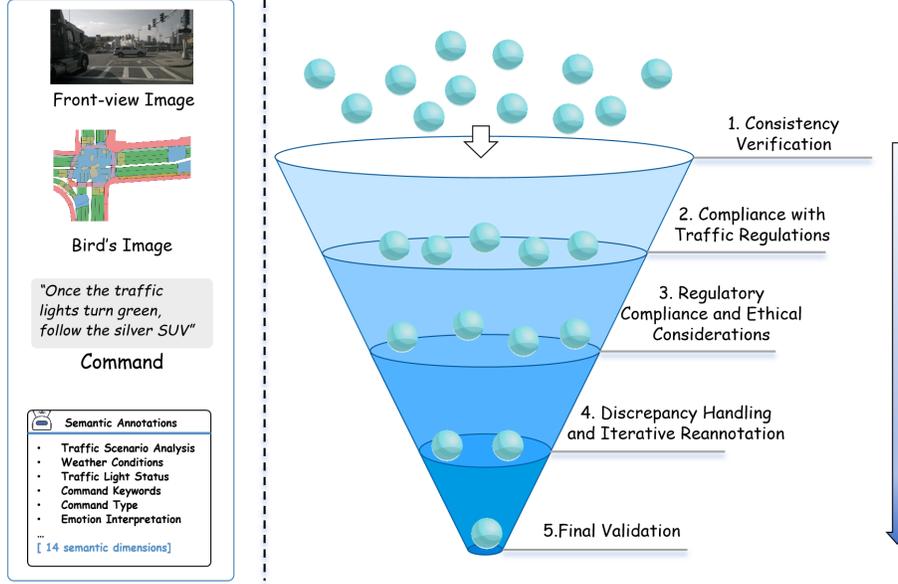


Figure A3. Example workflow for manual cross-checking process. The audit team evaluates LLM-generated annotations based on multiple input modalities, following a rigorous five-step review process: (1) Consistency Verification, ensuring alignment with ground-truth sensor data; (2) Compliance with Traffic and Safety Regulations, assessing adherence to legal driving standards; (3) Regulatory Compliance and Ethical Considerations, validating data privacy and ethical guidelines; (4) Discrepancy Handling and Iterative Reannotation, resolving annotation inconsistencies through expert review; and (5) Final Validation and Benchmarking, confirming dataset integrity before inclusion. Each LLM-generated annotation undergoes a rigorous multi-stage verification process to ensure accuracy, reliability, and compliance with industry standards, guaranteeing high-quality dataset integrity for real-world autonomous driving applications.

identification of command-relevant regions by integrating multi-modal reasoning, significantly improving model performance in complex and challenging traffic scenarios.

Cross-Modal Feature Transformation. Given a visual vector $O_v \in \mathbb{R}^{N \times C}$ and corresponding textual feature vectors $O_t \in \mathbb{R}^{C \times L}$, we first apply linear transformations to extract queries (Q), keys (K), and values (V) for each modality. This process can be formally defined as follows:

$$Q_v = W_v^Q O_v, K_v = W_v^K O_v, V_v = W_v^V O_v \quad (11)$$

$$Q_t = W_t^Q O_t, K_t = W_t^K O_t, V_t = W_t^V O_t \quad (12)$$

Then, we can compute the attention matrices using the scaled dot-product attention mechanism. Formally,

$$\begin{cases} A_{tt} = \phi_{\text{Softmax}} \left(\frac{Q_t (K_t)^\top}{\sqrt{d_k}} \right), A_{tv} = \phi_{\text{Softmax}} \left(\frac{Q_t (K_v)^\top}{\sqrt{d_k}} \right) \\ A_{vv} = \phi_{\text{Softmax}} \left(\frac{Q_v (K_v)^\top}{\sqrt{d_k}} \right), A_{vt} = \phi_{\text{Softmax}} \left(\frac{Q_t (K_t)^\top}{\sqrt{d_k}} \right) \end{cases} \quad (13)$$

Using these attention scores, we refine the visual and textual feature representations as follows:

$$O'_v = A_{vv} V_v + A_{vt} V_t + O_v \quad (14)$$

$$O'_t = A_{tt} V_t + A_{tv} V_v + O_t \quad (15)$$

These feature updates enable cross-modal interactions, allowing the model to dynamically integrate visual and textual contexts, improving the alignment of multimodal data.

Multi-Layer Feature Fusion. To progressively refine the representations across multiple layers, the i -th layer outputs are denoted as O_v^i and O_t^i . The final feature representation is obtained through hierarchical fusion. Mathematically,

$$F_{out} = \phi_{\text{MLP}} \left(\phi_{\text{GELU}} \left(\underbrace{O_v^0 | O_v^1 | \dots | O_v^{n-1}}_n \right) \right) \quad (16)$$

where F_{out} denotes the final output of MLD layers and ϕ_{GELU} represents the Gaussian Error Linear Unit (GELU) activation function, defined as follows:

$$\phi_{\text{GELU}}(x) = x\phi(x) = \frac{x}{\sigma\sqrt{2\pi}} \int_0^x e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt \quad (17)$$

Where $\phi(x) = P(X < x)$ for X following a normal distribution $N(0, 1)$. Moreover, μ is the mean, σ is the standard deviation, and e is the base of the natural logarithm.

D. Training Loss

To train our model, we calculate loss by the prediction set $\{s, y\}$ and ground-truth set $\{\hat{s}, \hat{y}\}$, with Z_v being the output of SA-WM. To fully utilize the information from each visual block, we require the model to simultaneously calculate and fit the IoU values between each visual block and the ground-truth region, which is denoted as y .

Stage-1: World-Model Rollout Pretraining. In the first stage, traditional loss functions such as dice loss and cross-entropy often struggle with imbalance issues in ground truth

representation due to the expansive observation areas in real traffic scenarios, numerous reference objects and samples, and a low proportion of positive objects. To address this, we draw inspiration from the field of medical imaging detection of minor lesions and integrate an adapted Tversky loss \mathcal{L}_{tve} [51] with Focal loss \mathcal{L}_{foc} [39]. Specifically, we initialize the value of P to 1 during the calculation of S in the first stage to accelerate the model’s visual-text mapping learning. This strategy leverages prior scores S and the downsampled ground-truth mask vector F_{GTmask} , interpolated via bilinear sampling, to compute the World-Model Rollout Pretraining loss. Given elements p_i and g_i corresponding to the i -th position of S and \hat{y} , the loss L_{rol} is formulated as follows:

$$\begin{aligned} L_{rol} &= \sum_{i=1}^N [\lambda_{tve} \mathcal{L}_{tve}(p_i, g_i) + \lambda_{foc} \mathcal{L}_{foc}(p_i, g_i)] \\ &= \lambda_{tve} - \frac{\lambda_{tve} |p_1 \cap g_1|}{|p_1 \cap g_1| + \alpha |p_1 \cap g_0| + \beta |p_0 \cap g_1|} \\ &\quad - \lambda_{foc} (\alpha_t (1 - p_t)^\gamma \log(p_t)) \end{aligned} \quad (18)$$

where the λ_{tve} and λ_{foc} are both the hyperparameters. Moreover, p_1 represents the set of scores in S predicted as positive examples, while p_0 denotes negative examples. Correspondingly, $|p_1 \cap g_1|$ represents true positives (TP), while $|p_1 \cap g_0|$ represents false negatives (FN). By adjusting the weighting parameters α and β , the model dynamically modulates its focus between false positives (FP) and false negatives (FN), improving class imbalance handling performance.

For focal loss \mathcal{L}_{foc} , the probability p_t is defined as:

$$p_t = \begin{cases} p, & \text{if class } t \text{ is positive} \\ 1 - p, & \text{otherwise} \end{cases} \quad (19)$$

where p is the model’s estimated probability of the class being present. γ is a focusing parameter that reduces the weight of well-classified examples, prioritizing hard, misclassified cases. α_t is a class-specific weighting factor to further adjust the impact of positive and negative samples. This loss enhances model flexibility, ensuring robust adaptation to imbalanced classes while improving visual-text alignment.

Stage-2: Grounding Decision Supervision. In the second stage, to prevent the model from forgetting previously acquired patch-level task knowledge, we employ a multi-task learning strategy to optimize both Binary Cross-Entropy (BCE) loss \mathcal{L}_{bce} and L1 loss \mathcal{L}_{L1} . We calculate loss L_{gro} by the prediction set $\{s, y\}$ and ground-truth set $\{\hat{s}, \hat{y}\}$, with the future latent states S is the SA-WM’s output that represents the mask annotations of referred object. Formally,

$$L_{gro} = \lambda_{bce} \mathcal{L}_{bce}(y, \hat{y}) + \lambda_{L1} \mathcal{L}_{L1}(S, \hat{s}) \quad (20)$$

where λ_{bce} and λ_{L1} are the hyperparameters. Overall, this diversity loss term incentivizes the model to further capture cross-modal interactions and progressively produce target predictions consistent with the commander’s intent.

E. Experiments Setups

E.1. Benchmarks

To evaluate our model’s effectiveness, we conduct experiments on the dataset zoo: Talk2Car, DrivePilot, MoCAD, and RefCOCO, RefCOCO+, and RefCOCog. These datasets provide diverse and complex real-world scenarios for benchmarking visual grounding in the field of autonomous driving.

Talk2Car. The Talk2Car dataset [12], an extension of the NuScenes dataset [2], consists of 11,959 natural language commands across 9,217 images captured in urban landscapes of Singapore and Boston. This dataset includes a variety of conditions, such as different times of day and weather scenarios, offering a challenging and diverse benchmark. The commands, averaging 11 words, contain complex instructions (e.g., “Parallel park behind the black car on our right”), requiring precise semantic understanding and scene reasoning. It enhances the NuScenes with bounding box annotations across 850 videos, with 55.94% of commands originating from Boston and 44.06% from Singapore. A detailed linguistic analysis reveals an average of 11.01 words per command, comprising 2.32 nouns, 2.29 verbs, and 0.62 adjectives, highlighting the linguistic diversity and complexity of the dataset. Each video is associated with an average of 14.07 commands, enriching the contextual learning process.

DrivePilot. We introduce DrivePilot, the first dataset to leverage Qwen’s linguistic capabilities for detailed semantic scene annotation using regularized prompts. This dataset categorizes urban scenes across 14 dimensions, including weather conditions, emotional context, and agent interactions. Each dataset entry comprises a natural language command, paired front-view and BEV images, scene annotations generated by Qwen2-VL, and precise target object locations. The dataset is designed to challenge models with object disambiguation and complex query interpretation, closely reflecting real-world AV navigation challenges.

MoCAD. The dataset originates from the first Level 4 autonomous bus deployed in Macau and has been continuously tested since 2020. The dataset spans 300+ hours of real-world driving, including data sets from a 5-kilometer campus route, a more extensive 25-kilometer city and urban road collection, and various open traffic situations observed under varying weather, time, and traffic density conditions. It comprises over 13,000 scene images and nearly 40,000 scene objects, with an average command length of 12.5 words, providing a rich dataset for visual grounding research. A distinctive aspect of MoCAD is its Macau-based driving environment, where right-hand driving contrasts with regions that enforce left-hand driving. This contrast introduces unprecedented challenges for VG in AVs, particularly in terms of context adaptation and cross-domain generalization.

RefCOCO. RefCOCO is built from the ReferItGame, a two-player interface for collecting referring expressions. It

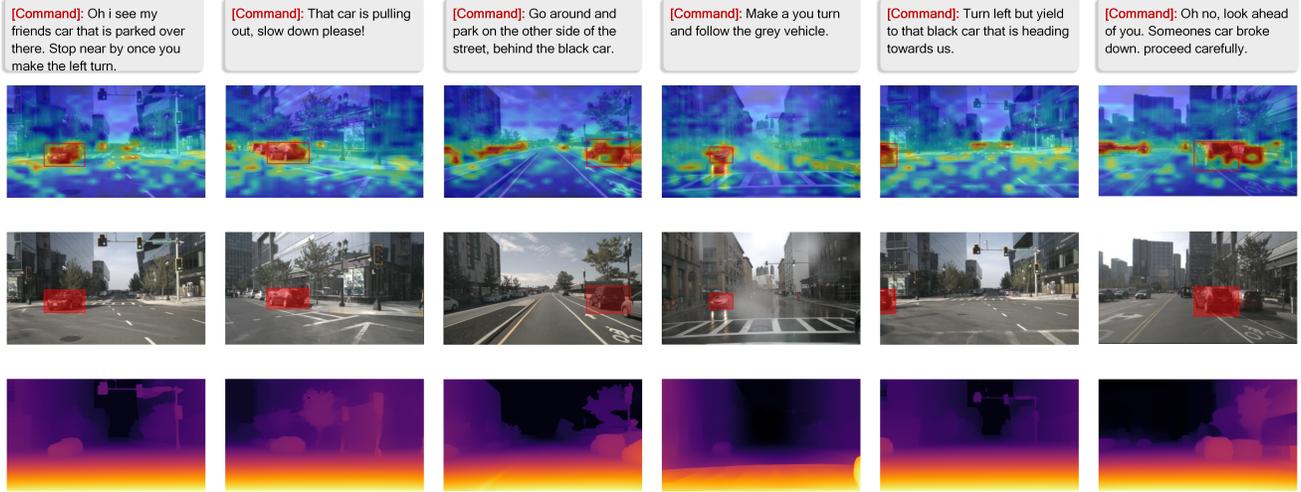


Figure A4. Visualization of ThinkDeeper’s multimodal grounding. The first line corresponds to the query commands for each scene, while the second row shows the future latent states Z_v generated by the SA-WM alongside the corresponding prediction. The third row highlights the ground truth regions in red mask areas, while the fourth row presents the depth map, providing spatial context for scene understanding.

contains 142,209 expressions for 50,000 objects in 19,994 images. The dataset is split into train/val/test sets, with the test set further divided into Test A (images with multiple people) and Test B (images with only objects).

RefCOCO+. Similar to RefCOCO, this dataset was also collected through the ReferItGame but with a restriction that prohibits the use of location-based descriptions. This constraint encourages the use of appearance-based expressions, making the task more challenging. RefCOCO+ comprises 141,564 referring expressions for 49,856 objects in 19,992 images. The dataset shares the same split structure as RefCOCO, including the Test A and Test B subdivisions.

RefCOCog. Collected in a non-interactive setting, RefCOCog features longer and more complex referring expressions, averaging 8.4 words per expression compared to 3.5 words in RefCOCO. It includes 95,010 expressions for 49,822 objects across 25,799 images. It is split into training, validation, and test sets, focusing on more descriptive and detailed language, which poses additional challenges for language comprehension and visual grounding models.

E.2. Implementation Details

Overall Configuration. Input images are resized to 384×384 pixels, and text expressions are truncated to a maximum of 50 tokens. During training, we apply a text augmentation strategy: with 30% probability, original descriptions are replaced by LLM-augmented versions, followed by keyword appending after a [SEP] token. We use the AdamW optimizer with a batch size of 32 and adopt a learning rate warmup over the first 10% of training steps. All components except the vision-language extractors (ViT and BERT) are trained with an initial learning rate of 10^{-4} . ViT and BERT

are initialized using BLIP [32] pre-trained weights, while other modules employ Xavier initialization [23].

Text Encoder. We use a pre-trained BERT model for text embedding extraction, configured with 16 hidden layers and an embedding vocabulary of 30,524 tokens. We also enforce a maximum sentence length of 50 tokens, with a layer normalization epsilon of $1e-12$, and a hidden size calibrated to $d = 768$ for linguistic input processing.

Vision Encoder. We use a Vision Transformer-Base (ViT-B) as the vision encoder with a 4:1 MLP to embed dimension ratio and 12-head multi-head attention, extracting a 24×24 visual token stream, and 3 cross-modal attention layers.

Spatial-Aware World Model. We use a three-layer cross-modal attention layer ($N = 3$) to compute prior scores, where each layer’s output is projected through a linear head. Each cross-modal attention block uses a hidden size of $D = 768$ for both text and visual inputs, with 8 attention heads and a dropout rate of 0.1. The learnable parameters in the discriminative module are initialized with $\mu = 1.0$ and $\sigma = 1.0$. Training proceeds in two stages: we first optimize the model with the world-model rollout loss L_{rol} for 15 epochs, followed by grounding-focused training with L_{gro} for an additional 55 training epochs (i.e., 70 epochs total).

Multimodal Decoder. Our architecture employs a cross-modal hypergraph where each visual node connects to its top $L = 8$ text nodes to form hyperedges, selected according to an affinity matrix computed by a 1536-dim MLP. Hypergraph attention uses 4 heads with LeakyReLU (negative slope 0.2), a 768-dim hidden layer, and 0.2 dropout. The multi-layer dynamic attention stack consists of 6 attention blocks (each followed by a linear layer), with 12-head multi-head attention, 768-dim hidden size, and 0.2 dropout.

E.3. Corner-case and Long-text Test Sets

To assess model robustness under real-world challenging conditions, we curate four specialized test subsets from the DrivePilot and MoCAD datasets. These include restricted visibility, multi-agent interactions, ambiguous prompts, and long, complex commands. In the multi-agent set, we focus on scenes containing more than 16 targets. The visual constraints set consists of scenarios with impaired visibility caused by nighttime conditions, fog, rain, camera obstructions, or low-resolution images. To evaluate the model’s ability to handle linguistic ambiguity, we identify unclear or potentially ambiguous commands, categorizing them as the ambiguous set. Additionally, recognizing that longer commands often introduce irrelevant details or increased complexity, we select commands exceeding 23 words and categorize them into the long-text test set, designed to assess the model’s capacity to process intricate instructions.

F. Quantitative Results

As shown in Figure A4, we present more quantitative results generated by ThinkDeeper in the DrivePilot dataset. These qualitative results showcase the superior performance of our proposed world model-inspired framework in integrating multimodal, real-world commands by jointly leveraging language, spatial, and visual cues, leading to improved multimodal grounding accuracy. These examples highlight our model’s ability to achieve robust localization under challenging conditions like high multi-agent and ambiguous scenes.