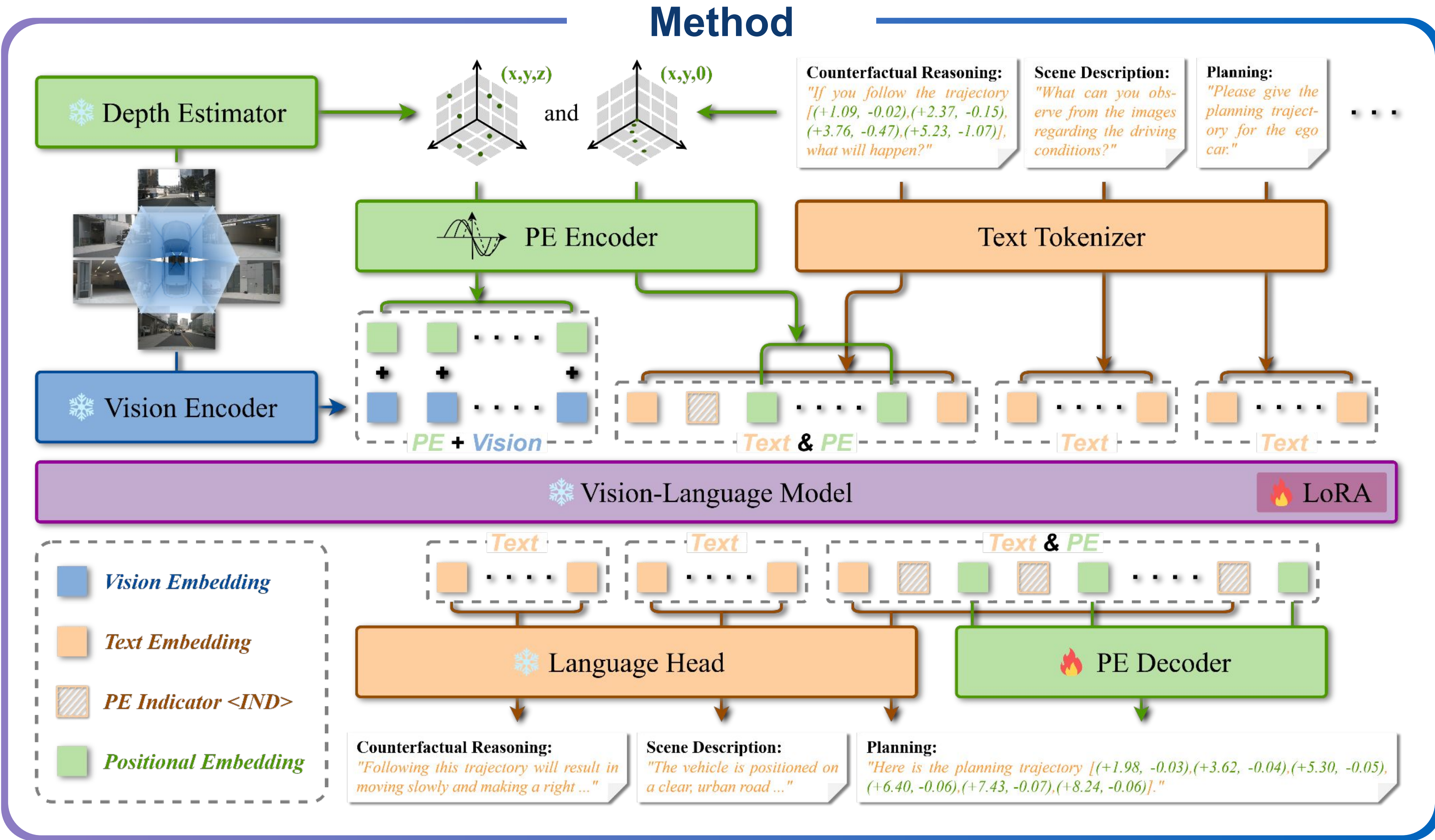
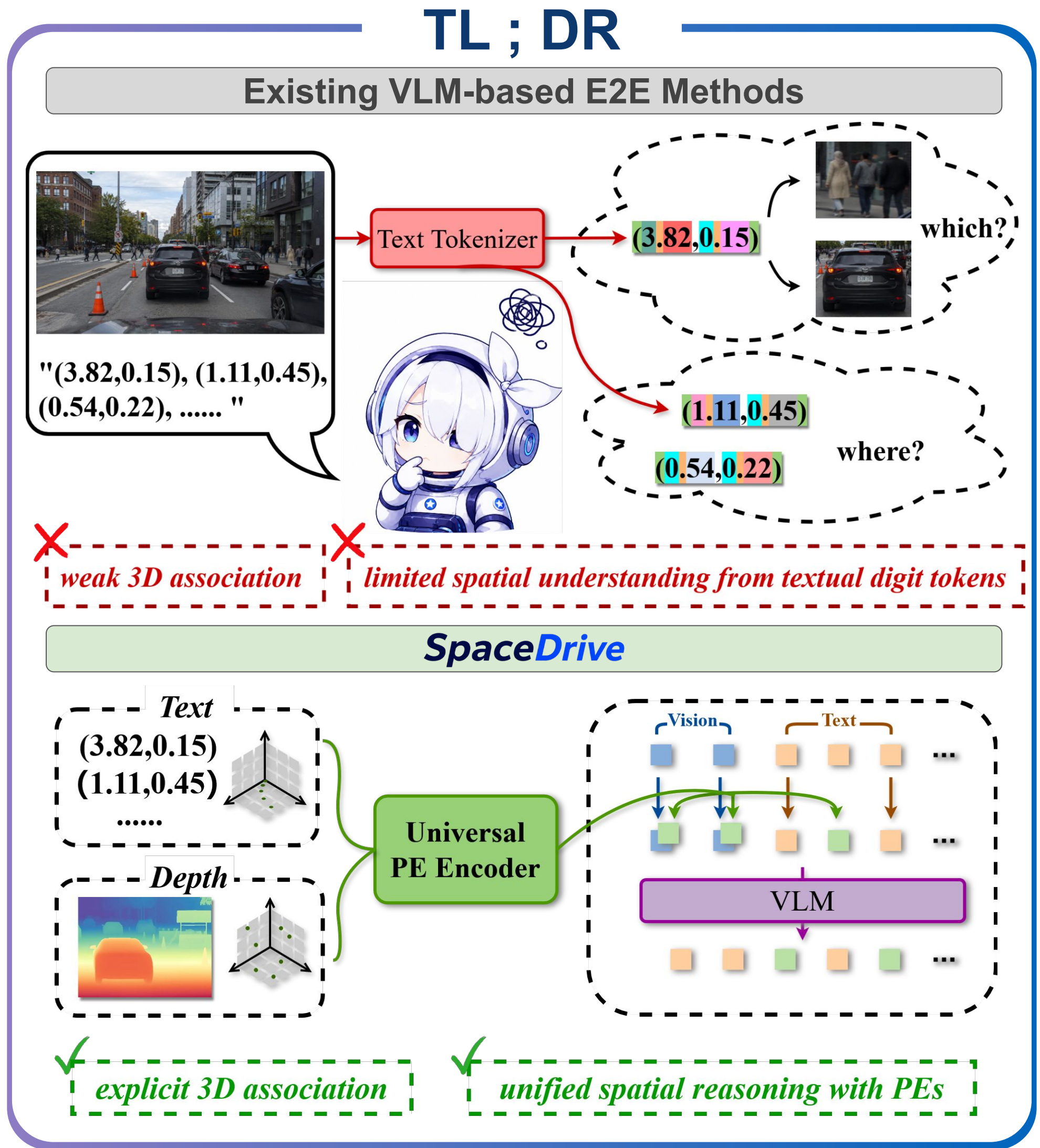


Peizheng Li<sup>\*1,2</sup>; Zhenghao Zhang<sup>\*1,4</sup>; David Holtz<sup>1</sup>; Hang Yu<sup>1,5</sup>; Yutong Yang<sup>1,6</sup>; Yuzhi Lai<sup>2</sup>; Rui Song<sup>7</sup>; Andreas Geiger<sup>2,3</sup>; Andreas Zell<sup>2</sup>

<sup>1</sup>Mercedes-Benz AG, <sup>2</sup>University of Tübingen, <sup>3</sup>Tübingen AI Center, <sup>4</sup>TU Munich, <sup>5</sup>Karlsruhe Institute of Technology, <sup>6</sup>University of Stuttgart, <sup>7</sup>UCLA



### Ablations

**Ablation 1: PE injection**

PE Injection	Avg. L2 ↓	Avg. Col. ↓	Avg. Int. ↓
baseline	2.51	4.53	6.77
+ PE on <b>Vision</b>	1.88	2.45	2.36
+ PE in <b>Text</b>	<b>1.80</b>	<b>1.88</b>	<b>4.21</b>
(+ Ego-state PE)	<b>0.32</b>	<b>0.23</b>	<b>1.27</b>

• Adding **PE on vision tokens** gives major gains. Unified PE for both **vision and text** improves further.

**Ablation 3: PE normalization**

$\alpha_{PE}$	Learnable	Avg. L2 ↓	Avg. Col. ↓	Avg. Int. ↓
1.0	✗	2.34	3.63	8.46
0.1	✗	2.43	3.79	9.42
0.02	✗	2.22	2.71	10.17
1.0	✓	1.82	2.04	4.62
<b>0.1</b>	✓	<b>1.80</b>	<b>1.88</b>	<b>4.21</b>
0.02	✓	1.86	2.03	5.42

• **Learnable PE** scaling is crucial for training stability.

**Ablation 5: Depth estimators**

$f_{dep.}$	Avg. L2 ↓	Avg. Col. ↓	Avg. Int. ↓
DepthAnythingV2	1.76	1.95	3.96
UniDepthV2	1.80	1.88	4.21

• Effectiveness of SpaceDrive is **independent** from specific pre-trained depth models.

**Ablation 2: PE encoder & decoder**

Encoder ( $\phi$ )	Decoder ( $\psi$ )	Avg. L2 ↓	Avg. Col. ↓	Avg. Int. ↓
Sine-Cosine	Coordinate-wise	<b>1.80</b>	<b>1.88</b>	<b>4.21</b>
MLP	Coordinate-wise	1.96	3.17	6.76
RoPE	Coordinate-wise	1.93	3.71	11.40
Sine-Cosine	Sine-Cosine	1.87	2.62	9.20
Sine-Cosine	Task-specific	1.93	2.41	5.58

• Sine-Cosine inherently incorporates **relative position modeling**.  
• Inverting Sine-Cosine encoding is **worse than** learnable decoder.

**Ablation 4: Adaptability across VLMs**

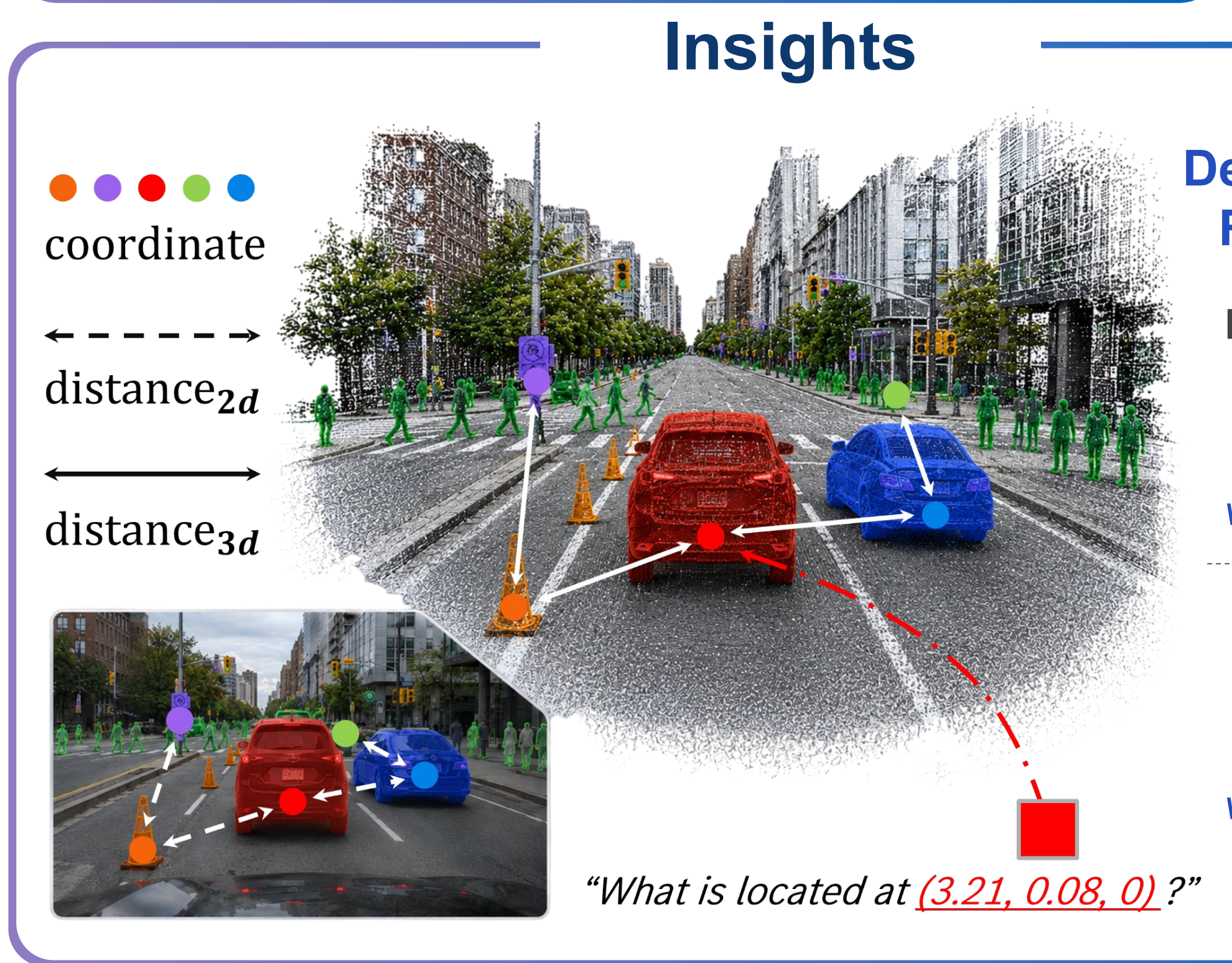
Method	Backbone	Avg. L2 ↓	Avg. Col. ↓	Avg. Int. ↓
SpaceDrive	LLaVA	1.82	2.44	4.08
SpaceDrive	Qwen-VL	1.80	1.88	4.21
SpaceDrive+	LLaVA	0.31	0.23	1.42
SpaceDrive+	Qwen-VL	0.32	0.23	1.27

• The unified spatial representation works on **many** VLM backbones.

**Ablation 6: Robustness to Depth Errors**

Setting	Depth Noise	Avg. L2 ↓	Avg. Col. ↓	Avg. Int. ↓
Training & Inference	w. -5% Global Shift	1.86	1.80	5.22
Training & Inference	w. 5% Global Shift	1.86	1.93	4.41
Training & Inference	w. 5% Random Noise	1.83	2.16	4.44
Inference only	w. -2.5% Global Shift	1.80	1.89	4.22
Inference only	w. 2.5% Global Shift	1.80	1.86	4.24
Inference only	w. 2.5% Random Noise	1.80	1.89	4.17

• PE serves as hints for **spatial association in training**, not an extra modality input.



### Benchmarks

**A. Open-loop planning on nuScenes**

Method	Avg. L2 ↓	Avg. Col. ↓	Avg. Int. ↓
UniAD	0.46	0.37	1.59
UAD	0.30	0.27	1.37
Drive-WM	0.80	0.26	-
EMMA	0.32	-	-
RDA-Driver	0.80	0.32	-
DriveVLM	0.40	0.27	-
ORION	0.34	0.37	-
OmniDrive-Q	1.98	3.79	4.59
OmniDrive-Q++	0.33	0.30	3.00
<b>SpaceDrive</b>	<b>1.80</b>	<b>1.88</b>	<b>4.21</b>
<b>SpaceDrive+</b>	<b>0.32</b>	<b>0.23</b>	<b>1.27</b>

**SpaceDrive** achieves the **best** open-loop performance among VLM-based methods on unScenes.

**B. Closed-loop planning on Bench2Drive**

Method	Driving Score ↑	Success Rate ↑
ReAL-AD	41.17	11.36
Dual-AEB	45.23	10.00
VDRive	66.25	50.51
StuckSolver	70.89	50.01
DriveMoE	74.22	48.64
ETA	74.33	48.33
VLR-Drive	75.01	50.00
ORION	77.74	54.62
SimLingo	85.07	67.27
<b>SpaceDrive+</b>	<b>78.02</b>	<b>55.11</b>

**SpaceDrive** achieves **2nd-best** in VLM-based planners on Bench2Drive w/o specific closed-loop optimizations.

