# ELM:
# Embodied Understanding of Driving Scenes

(A 60 min Talk)

Online

August 2024

Yunsong Zhou, Shanghai Jiao Tong University / Shanghai AI Lab

# Outline

- **Introduction to End-to-end Autonomous Driving**
  - Modular vs End-to-end
  - Industry Applications
  - Research Roadmap
- **Challenges**
- **Mainstream Lineup of Work**
  - Embodied Understanding of Driving Scenarios

# Introduction

End-to-end Autonomous Driving

# Problem setting | Autonomous Driving (AD) Tasks



Bounding boxes

Waypoints

Trajectory

Perception → Prediction → Planning →

What are around?

How will they go in the future?

Where should I go?

**Challenge |** Various weathers, illuminations, and scenarios

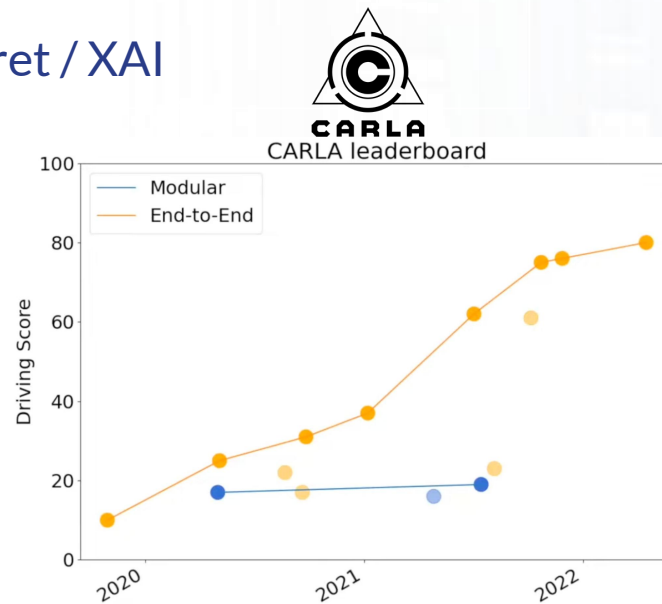OpenDriveLab

# Motivation | Why End-to-end (E2E) Autonomous Driving?

**Disadvantages**

- Closed-loop evaluation only available in simulator / onboard test

- Lack of real-world data

- Hard to interpret / XAI

*Credit to Andreas Geiger @ CVPR Workshop 2023*

**E2E vs Non-E2E**



CARLA leaderboard

https://leaderboard.carla.org/



**Performance**

**E2E Approach**

**Human Expert**

**Non E2E System**

**We are here!**

**Time/Readiness**

知乎 @曹越

*Credit to Dr. Yue Cao @ Zhihu*

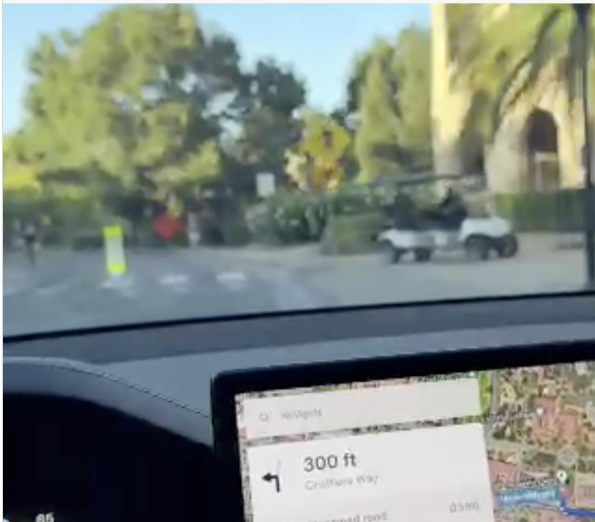OpenDriveLab

## Industry

### E2E Vehicle

v12 is reserved for when FSD is **end-to-end AI,** from images in to steering, brakes & acceleration out.

### E2E Robot



**Ashok Elluswamy** ✔
@aelluswamy

This end to end neural network approach will result in the safest, the most competent, the most comfortable, the most efficient, and overall, the best self-driving system ever produced. It's going to be very hard to beat it with anything else!

**Elon Musk** ✔ ✖ @elonmusk · Aug 26
twitter.com/i/broadcasts/1...



**Tesla Optimus** ✔ 🔳 @Tesla_Optimus · Sep 24

Optimus can now sort objects autonomously 🤖

Its neural network is trained fully **end-to-end**: video in, controls out.

*No hard-code.*
*Completely learning on its own.*
*End-to-end, video to neural network to controls.*
*Don't need map data at all, only coordinates!*
*No cellular connection needed.*

2021    2022

2023

### My Opinion

- Probably e2e as a backup module
- Massive high-quality data prevail
- Mapless is promising and feasible

OpenDriveLab

# Application｜End-to-end Autonomous Driving

**Tesla FSD Beta v12.12 rolls out to customers**

**The End-to-End Approach in Autonomous Driving**

Two types of end-to-end implementation:

**Next-Generation Automated Driving: the Power of End-to-End AI**

多模块走向端到端融合， "轻" 车更熟路

端到端感知决策大模型
小米自研 全球首次应用于量产车

- *Tesla website*
- *Mobileye  at CES 2024*
- *Wayve website*
- *DeepRoute (元戎启行) at GTC 2024*
- *Mi (小米) Automobile Technological Event 2024*

OpenDriveLab

# 研究背景与现状 - 端到端



Classroom | Wild Environment | feedback | Strong Agent

Supervisor

exploration

Weak Agent

Student

Behavioral cloning from expert trajectories | Self-evolution based on environmental feedback

## AGI（L5）
*Drive Anywhere/Anytime*

↓

泛化、智能、安全

## 80% 自动驾驶问题
已被解决（L2）

## 剩余10%
**Corner case in Daily Life**
靠规则永远写不完（L3/L4）

Bad Cases in Tesla FSD v12.3 Test drive before CVPR'24



820m  ↑↑↑  **45** km/h  60  60 ▶606m

→  加数据?
如何采?  ←
如何训?

↓

**通用的端到端具身智能体进化方法**

Open**Drive**Lab

# Roadmap | End-to-end Autonomous Driving



**CARLA** Launched
DS: 8.94

**CARLA**
DS: 24.98

**CARLA**
DS: 47.65

**CARLA**
DS: 79.95

**CARLA v2** Launched
DS: 0.01

**nuPlan** Launched
Score: 0.90

**Reinforcement Learning (RL)**
CIRL, MaRLn, GRI
**Drive in A Day**
Agent / Reward

**Policy Distillation**
WOR, Roach, TCP
**LBC**
Expert → 
Privileged Agent ⇒ Sensorimotor Agent

**Modality / Advanced Structure**
InterFuser, ThinkTwice
**Transfuser**
Transformer

**Data Generation**
Advsim, L2C
**KING**
Driving Policy $\pi_w$
Critical Scenario

ALVINN
Direction
Input

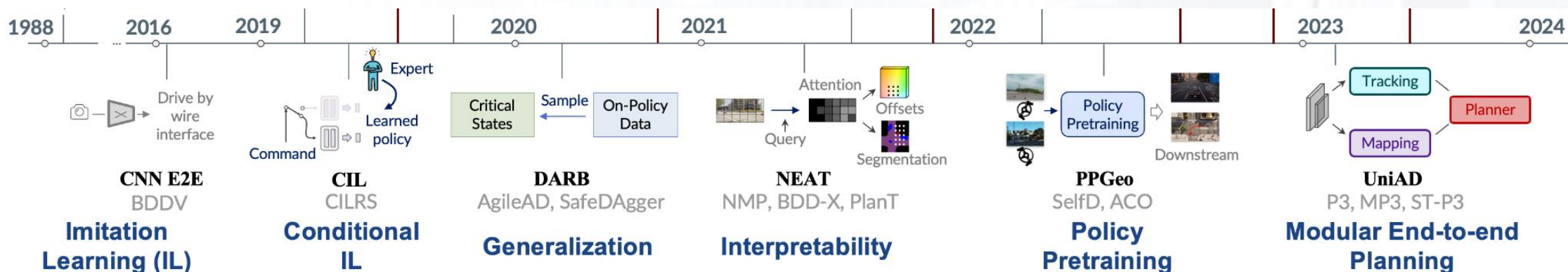1988 ... 2016 2019 2020 2021 2022 2023 2024

## Summary (1/2)

- Carla leaderboard gets much improved over the years. With new mapping / routes (Carla v2) and nuPlan benchmark, this field got so much to do.

- RL method is prevalent in the beginning (since it's natural)

- Input modality and more advanced structure boosts the performance

OpenDriveLab

# Roadmap | End-to-end Autonomous Driving

## Summary (2/2)

- The First Neural Net based method dates back to 2016 using Imitation Learning
- Learned policy from Experts (IL), with data augmentation, could prevail in performance
- Interpretability, with explicit design in the network stands out recently
- End-to-end design comes to obsess many merits in previous attempt

Timeline:

1988 — 2016 ... — 2019 — 2020 — 2021 — 2022 — 2023 — 2024

**CNN E2E**
BDDV
**Imitation Learning (IL)**

**CIL**
CILRS
**Conditional IL**

**DARB**
AgileAD, SafeDAgger
**Generalization**

**NEAT**
NMP, BDD-X, PlanT
**Interpretability**

**PPGeo**
SelfD, ACO
**Policy Pretraining**

**UniAD**
P3, MP3, ST-P3
**Modular End-to-end Planning**

# Analogy to General Domains in  CV/NLP/Robotics

| | Domain | Method Abbreviation | | Institute / Time | Data Scale | Public? |
|---|---|---|---|---|---|---|
| **General Large Models** | NLP (LLM) | GPT-4 | | OpenAI / 2023.3 | 13T tokens | ✗ |
| | | LLaMA 2 | | Meta / 2023.7 | 2T tokens | ☑ |
| | Vision | ViT-22B | | Google / 2023.2 | 4B images | ✗ |
| | Vision Language (LLM backend) | BLIP-2 | | Salesforce / 2023.1 | 129M images-text pairs | ☑ |
| **Industrial Large Models (Application)** | Autonomous Driving  nuScenes: 4.5h | **DriveAGI (GenAD)** | | OpenDriveLab / 2023.11 | **2000 h videos (public)** | ☑ |
| | | GAIA-1 | | Wayve / 2023.6 | 4700 h videos | ✗ |
| | | World Model Demo | | Tesla / 2023.6 | Unknown (Large-scale) | ✗ |
| | Robotics (LLM backend) | PaLM-E | | Google / 2023.3 | Unknown (Large-scale) | ✗ |
| | | RT-2 | | DeepMind / 2023.7 | 1B img-text pairs / 13 robots / 17 months | ✗ |

**If taken seriously for AD: lots of compute** (at least 200 A100s) +
**massive amount of data** (at least 10k hours of diverse, high-quality data)

OpenDriveLab

# Trending in E2EAD | Driving + Language



Go straight at an intersection then turn left.
There are **construction cones** on the road.

**HAD** — human-to-vehicle driving advice, highlighting key objects.

**Explainable Driving Behavior**

Planning

**BDD-X** — one-sentence explanation of driving behavior.

*2019*

*2022*

Prediction

Full-stack

*2023*



**Action description:**    **Action justification:**

(1) The car is driving   **as** there is nothing to impede it.

# Trending in E2EAD | Driving + Language

| Dataset | Source Dataset | # Frames | Avg. captions / QA per annotated frame | Total captions / QA in Perception | Total captions / QA in Prediction | Total captions / QA in Planning | Logic among captions/QA pairs |
|---|---|---|---|---|---|---|---|
| nuScenes-QA [47] | nuScenes | 34,149 | 13.5 | 460k** | ✗ | ✗ | None |
| nuPrompt [66] | nuScenes | 34,149 | 1.0 | 35k* | ✗ | ✗ | None |
| HAD [31] | HDD | 25,549 | 1.8 | 25k | ✗ | 20k | None |
| BDD-X [30] | BDD | 26,228 | 1 | 26k | ✗ | ✗ | None |
| DRAMA [40] | DRAMA | 17,785 | 5.8 | 85k | ✗ | 17k | Chain |
| Rank2Tell [51] | Rank2Tell | 5,800 | - | - | ✗ | - | Chain |
| DriveLM-nuScenes | nuScenes | 4,871 | **91.4** | 144k* | 153k | 146k | **Graph** |
| DriveLM-CARLA | CARLA | **183,373** | 20.5 | **2.46M**** | **578k**** | **714k**** | **Graph** |

Table 1. **Comparison of DriveLM-nuScenes & -CARLA with Existing Datasets.** * indicates semi-rule-based labeling (w/ human annotators), ** indicates fully-rule-based (no human annotators), and - means publicly unavailable. DriveLM-Data significant advances annotation quantity, comprehensiveness (covering **perception, prediction and planning**), and logic (chain to **graph**).

**Rank2Tell** — reasoning for the rank of objects' importance level.

**DriveLM** — perception-prediction-planning driving description with graph VQA.

**Talk2Car** — a description of how to reach the goal point from current position.

**HAD** — human-to-vehicle driving advice, highlighting key objects.

**DRAMA** — caption about important objects and future decision.

**Explainable Driving Behavior**

**Planning**

**Prediction**

**Full-stack**

**BDD-X** — one-sentence explanation of driving behavior.

*2019*

*2022*

*2023*

**Berkeley DeepDrive**

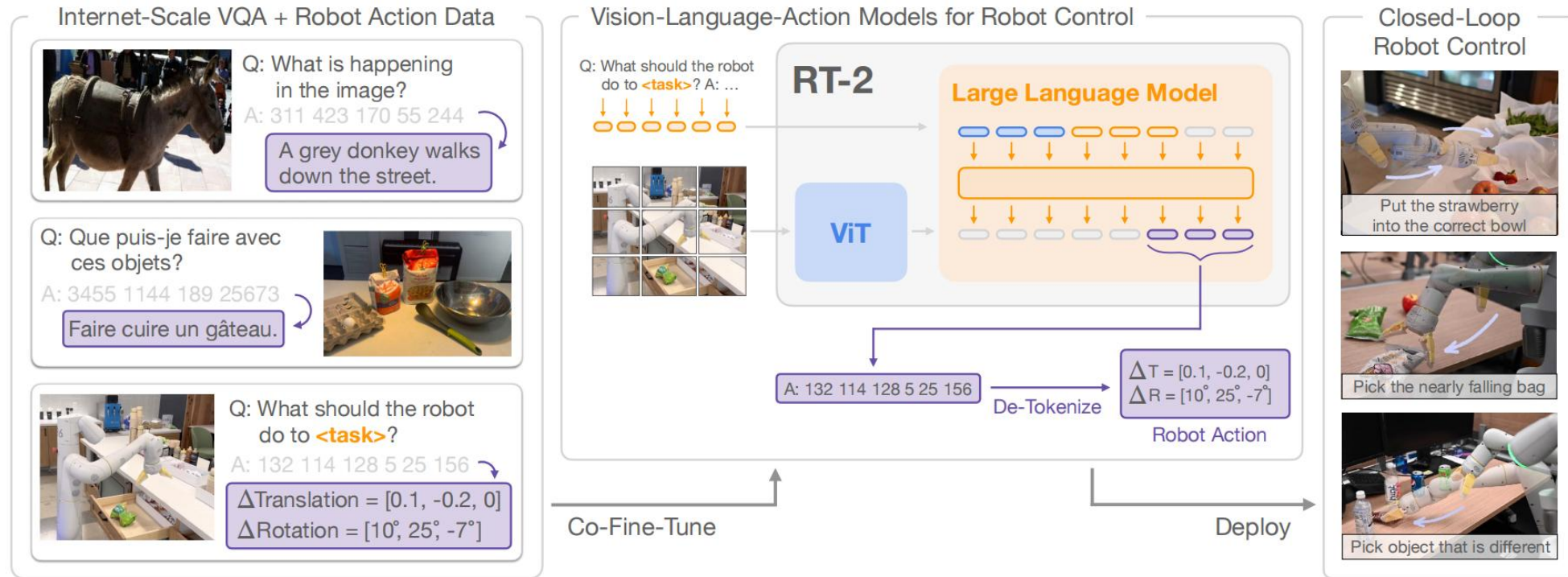The construction worker in blue dress is standing on the left side of the road. Please follow his instructions.

**LINGO-1** — commentary for explaining driving behaviours.

**For now, language into driving is marginal (trivial).**
**Serves only as a "commentator". We haven't verified (or seen) the effectiveness.**

OpenDriveLab

# Insight | VLM in Robotics / Embodied AI



- How vision-language models trained on Internet-scale data can be incorporated directly into **end-to-end robotic control**

- Goal: to **boost generalization** and enable emergent semantic reasoning

- Robotic tasks naturally fits into language at dissecting tasks step by step using language (prompt).

- Is it the <u>right way</u> to open the language tool box as does in Robotics for Autonomous Driving?

**Key ingredient(s): huge amount of data (not public) + language prompt to dissect tasks**

OpenDriveLab

Grand Challenge 2024

# CVPR 2023 AD Challenge - Recap

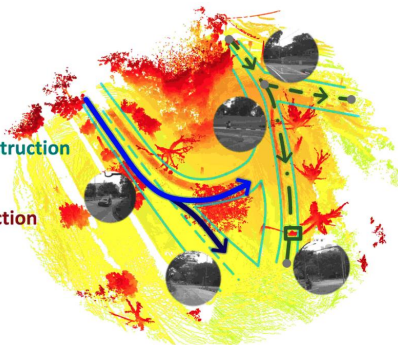https://opendrivelab.com/AD23Challenge.html

## Overview

### Tracks

**Track 1**
OpenLane Topology

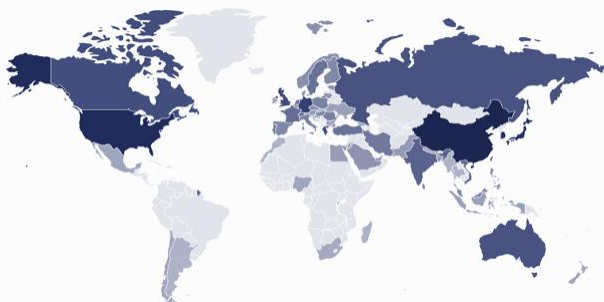**Track 2**
Online HD Map Construction

**Track 3**
3D Occupancy Prediction

**Track 4**
nuPlan Planning

### Internationalization

### Submissions

| | | | | | |
|---|---|---|---|---|---|
| 19 | 32 | 34 | 36 | 72 | 110 |
| AV2 2023 | Waymo 2023 | Apollo 2022 | Carla 2020 | Huawei 2022 | Ours |

### Teams

| | | |
|---|---|---|
| ~20 | ~40 | 277 |
| BDD / Berkeley 2022 | AV2 2022 | Ours |

> **2.3k**
Submissions

> **15**
Contries / Regions

> **68k**
Website / Social Media Views

## Competing Teams

## Attention from All Sides

**NVIDIA**
HOME  AI  DATA CENTER  DRIVING  GAMING  PRO GRAPHICS  ROBOTICS  HEALTHCARE  STARTUPS  AI

**NVIDIA Research ... Award at CVPR**
New work introduces sta...
June 15, 2023 by CALISA COLE

#CVPR2023 #IndustrialAI #automateddriving #research #competition #nuplan

**BOSCH**
Bosch Center for Artificial Intelligence (BCAI)
54,190 人關注
＋ 關注
瀏覽完整頁面

旷视科技
已認證賬號
你經常看 自动驾驶 相关内容
＋ 关注

近日，为期三个月的 CVPR 2023 自动驾驶国际挑战赛比赛结果揭晓。旷视研究院在OpenLane Topology 赛道中击败 30 余支国内外队伍，夺得冠军。

**42dot**  Company  Technology  Mobility  AEV  Careers  Blog
Tech
**Team 42dot Wins 2nd Place in the Autonomous Driving Challenge at CVPR 2023**

**International Teams ... Advance AV Resear... 2023 nuPlan Challe...**
September 06, 2023 | NUPLAN

**Challenges of Mass Production Autonomous Driving in China**
And the Recent Progress from Xpeng Motors in 2023

Patrick Langechuan Liu · Follow
...rds Data Science · 7 min read · Jun 19, 2023

ころんびあ
@colum2131

実はTuringでCVPR Challenge 2023参加していましたー！
自分はTrack 1のOpenLane Topologyという走行レーンや交通信号機・標識の検出や、それら間にトポロジー(意味的な接続性)があるか分類するタスクにチームで参加 ...

什么样的方案，夺得了CVPR自动驾驶挑战赛冠军？

雷峰网  机器之心  2023-06-19 19:21  发表于北京

你經常看 自动驾驶 相关内容

自动驾驶圈正被"去高精地图"风暴席卷。

今年6月中旬，一年一度的AI顶会 CVPR 2023 在加拿大举行，同期还有场自动驾驶国际挑战赛，这场挑战赛吸引了来自15个国家的270多支队伍参赛，有2300多件作品投稿，赛况激烈。

**36Kr**  旷视刷榜了什么自动驾驶比赛？
旷视研究院参加的这个比赛，是CVPR 2023专门面向自动驾驶感知决策系统设立的挑战赛。

**OpenDriveLab**

# CVPR 2024 Autonomous Grand Challenge

**120,000 USD Bonus!**

https://opendrivelab.com/challenge2024/

## Timeline

**2023.12**
Contest tracks open for comments

**2024.02**
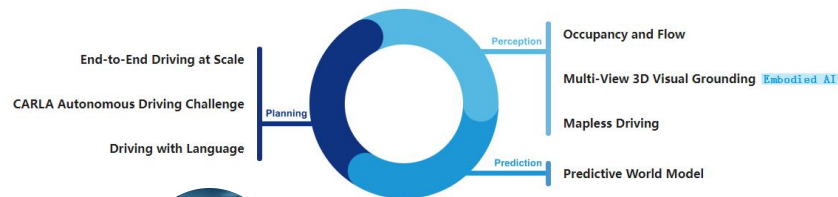Track Announcement

**2024.03**
Challenge Starts
Dataset / Evaluation Metrics
Public

**2024.06**
Workshop @ Seattle
Competition ends
Announcement of winning teams

## Seven Brand-new Tracks

End-to-End Driving at Scale — Perception — Occupancy and Flow
CARLA Autonomous Driving Challenge — Multi-View 3D Visual Grounding — Embodied AI
Driving with Language — Planning — Mapless Driving
Prediction — Predictive World Model

End-to-End Driving at Scale

Predictive World Model

Occupancy and Flow

Multi-View 3D Visual Grounding

CARLA AD Challenge

Driving with Language

Mapless Driving

## Cooperation from All Sides

OpenGVLab

OpenRobotLab

清华大学 Tsinghua University

Hugging Face

Motional

NVIDIA

CARLA

美团 Meituan Autonomous Delivery

EBERHARD KARLS UNIVERSITÄT TÜBINGEN

WAYVE

机器之心 Synced

OpenDriveLab

# Embodied Understanding of Driving Scenarios

Yunsong Zhou[1,2*]    Linyan Huang[1*]    Qingwen Bu[1,2*]    Jia Zeng[1]    Tianyu Li[1,3]

Hang Qiu[4]    Hongzi Zhu[2†]    Minyi Guo[2]    Yu Qiao[1]    Hongyang Li[1†]

[1] OpenDriveLab, Shanghai AI Lab    [2] Shanghai Jiao Tong University

[3] Fudan University    [4] University of California, Riverside

*Credits:*
*https://github.com/OpenDriveLab/ELM,*
*https://arxiv.org/abs/2403.04593*

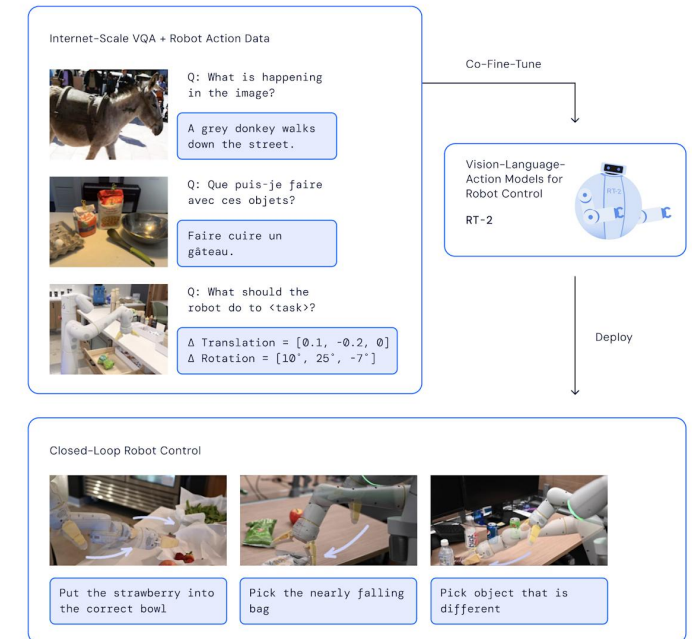Introduction

# ELM - Introduction

**Embodied Understanding**

> " This makes transfer of internet scale knowledge to robots more direct, and may provide a more scalable class of approaches in the future. "

**Interacting** with environments & **reasoning** via common sense



4x speed

"push the blue block to the kethcup"

4x speed

"push the blue block to the tabasco"

Google DeepMind

RT-2

RT-X

PaLM-E 562B

BAIR
BERKELEY ARTIFICIAL INTELLIGENCE RESEARCH

TOYOTA RESEARCH INSTITUTE

*Credits:*
*Google DeepMind, RT-2, RT-X, PaLM-E*

**Internet-Scale VQA + Robot Action Data**

Q: What is happening in the image?
A grey donkey walks down the street.

Q: Que puis-je faire avec ces objets?
Faire cuire un gâteau.

Q: What should the robot do to <task>?
Δ Translation = [0.1, -0.2, 0]
Δ Rotation = [10°, 25°, -7°]

Co-Fine-Tune

Vision-Language-Action Models for Robot Control
RT-2

Deploy

**Closed-Loop Robot Control**

Put the strawberry into the correct bowl

Pick the nearly falling bag

Pick object that is different

OpenDriveLab

# ELM - Introduction

**Embodied Understanding**

**Vision Language Models**

**Interacting** with environments & **reasoning** via common sense

Only focus on 2D domain, *i.e.,* **description**

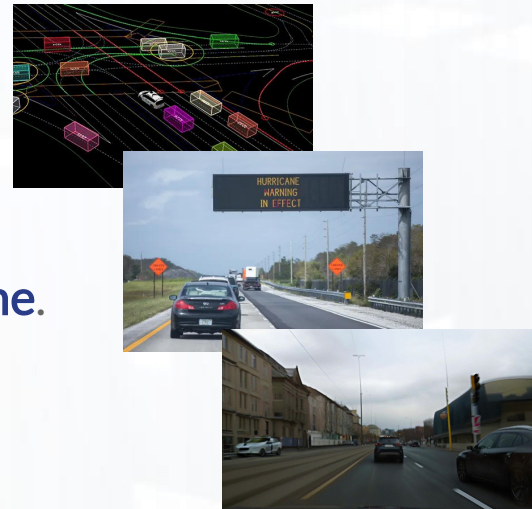*Credits: DriveLM; Open-sourced Data Ecosystem in Autonomous Driving*

# ELM - Introduction

**Embodied Understanding**

**Vision Language Models**

**Embodied Language Model**

**Interacting** with environments & **reasoning** via common sense

Only focus on 2D domain, *i.e.,* **description**

Expanding vanilla VLMs to driving scenes

- Task: embodied understanding of driving scenarios.

- Capabilities: description, **localization**, **memorization**, **forecasting**.

- Model: **ELM** with long-horizon **space** and **time**.

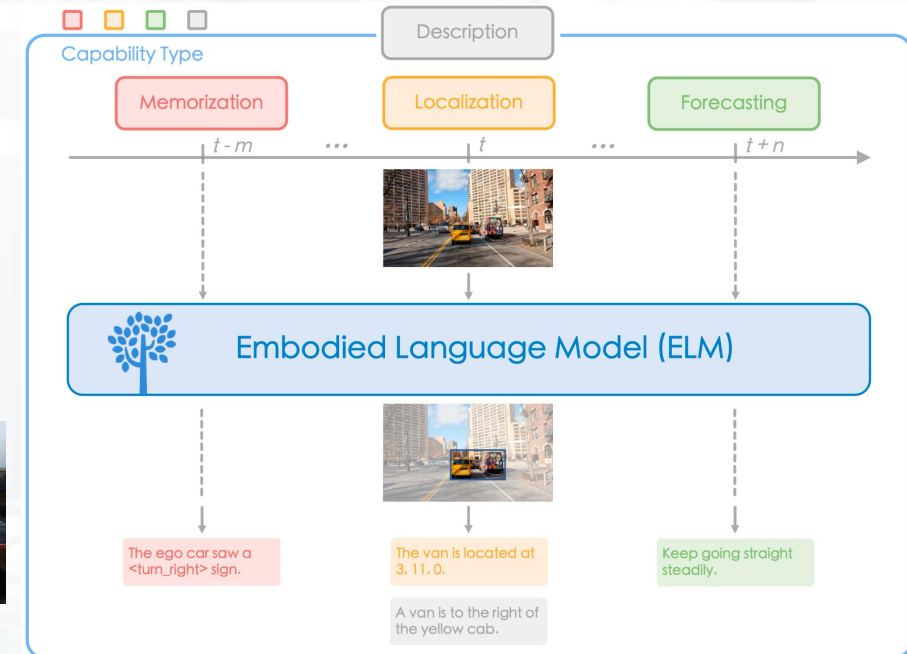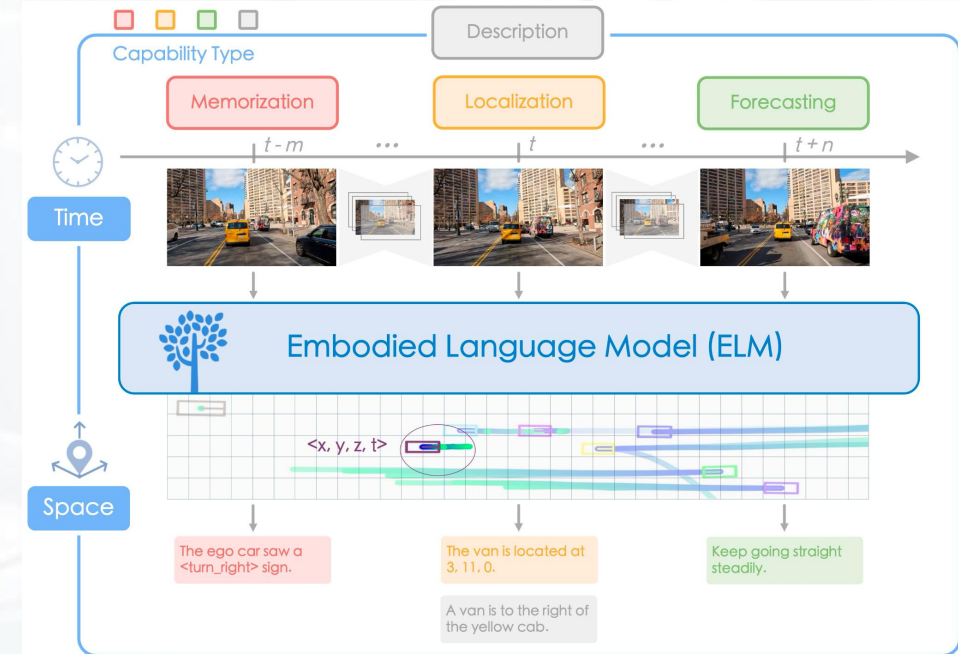- Benchmark: A spectrum of tasks in an embodiment setting.

# ELM - Introduction

| Embodied Understanding | Vision Language Models | Embodied Language Model |
|---|---|---|
| **Interacting** with environments & **reasoning** via common sense | Only focus on 2D domain, *i.e.,* **description** | Expanding Vanilla VLMs to Driving Scenes |

- Task: embodied understanding of driving scenarios.

- Capabilities: description, **localization**, **memorization**, **forecasting**.

- Model: **ELM** with long-horizon **space** and **time**.

- Benchmark: A spectrum of tasks in an embodiment setting.



OpenDriveLab

# Embodied Understanding

## Surrounding Narration
Q: Please give me a visual description of the object to the front of the ego car.
A: A yellow taxi is driving on a street passing by trees.
65.2 | 73.2

## Traffic Sign Inquiry
Q: Has the ego vehicle seen any traffic sign before?
A: Yes. The car has seen 1 <go_straight> before.
63.1 | 76.5

## Action & Decision
Q: What are safe actions to take for the ego vehicle?
A: Slow down to Keep a safe distance.
68.7 | 74.4

### Capability Type
- Memorization
- Localization
- Forecasting

$t-m$ ... $t$ ... $t+n$

Time
Space

## Embodied Language Model (ELM)

$<x, y, z, t>$

The ego car saw a <turn_right> sign.
The van is located at 3, 11, 0.
Keep going straight steadily.
A van is to the right of the yellow cab.

## Egocentric Narration
Q: Please give a caption for this action.
A: He turns the steering wheel to cross the intersection ahead.
16.9 | 26.5

## Moment Recap
Q: What happened 30 seconds before in the history?
A: He drived pass the motorcycle.
13.1 | 22.6

## Event Query
Q: What happened between "starts the vehicle" and "passes the black car"?
A: He drives along the right side of the road.
27.3 | 38.0

## Tracking
Q: What is the location of <c2> 2 seconds before in the history?
A: Car, 2, 11, -1.
3.0 | 14.0

## Box Detection
Q: Please provide the coordinates and category of <c1>.
A: Car, -3, 29, -1.
5.1 | 51.6

## Box Prediction
Q: What is the location of <c0> in next 3 seconds.
A: Pedestrian, 2, 13, 0.
3.6 | 15.1

## Activity Prediction
Q: What will happen in the next 3 seconds in the future?
A: He will drive through the junction.
11.5 | 18.1

Legend:
- BLIP2-flant5 (baseline)
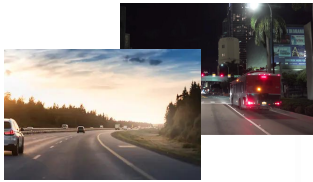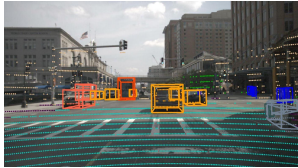- ELM on nuScenes
- ELM on Ego4D

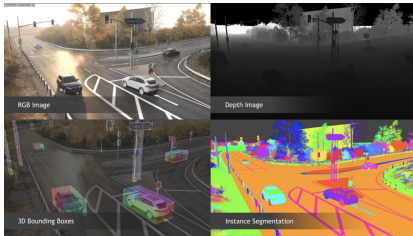We expand a wide spectrum of **new tasks** to fully leverage large language models in an embodiment setting
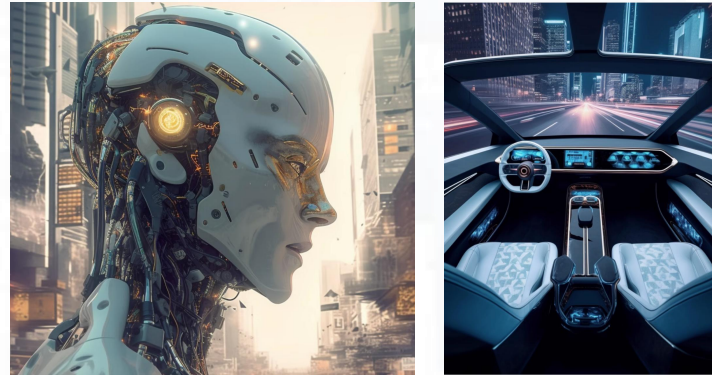
OpenDriveLab

# ELM - The Big Picture

## Data-centric Pipeline

## Pre-training DriveCore

## Benchmark

### Data Collection

Motional
ARGO AI
lyft

YouTube

### Data Generation

RGB Image | Depth Image
3D Bounding Boxes | Instance Segmentation

### Autonomous Driving (nuScenes)

### Robotics (Ego4D)

Introducing
Ego4D

## Universal Foundation Model
for autonomous driving

How to formulate?
What's the objective goal?
**Space** and **time**?

*Partial photo by courtesy of online resources.*

OpenDriveLab

# Embodied Understanding

## Location Labeling

Prompt — 1. GPT-4 → Meta Template — 2. Sampling → Inspector — 3. Select → Template

Point Cloud (x, y, z)
Pixel <u, v>

Data Batch — 4. Sampling & Filling → Generated Data

Q: What are the 3D spatial coordinates of the pixel at <Pixel Location>

A: Location: x, y, z, Category: Car

## Description Labeling

nuScenes    Waymo
Youtube    Ego4D

**Example Image**

**Refined Caption**

The driving scene in the image **features a busy city street** with···

The **traffic light is green**··· ✓

The ego-vehicle **should proceed with caution and**···

Manual Revision ↑ 8.

World Data — 1. Crawl → Raw Data → Data Batch — 4. Feed → Vision-Language model

2. Sampling → 3. Select

Blacklist — Discard ← Quality Check

5. Generate → Caption Batch — 6. Sampling → Generated Caption — 7. ✓

**Feedback** The traffic light is **red**··· ✕

Rerun ← Quality Check → Discard

## Data Status

| Method | Pre-train Data | # | Time | City | Anno |
|---|---|---|---|---|---|
| LLaVA [48] | COCO [47] | 150K | - | - | Des |
| VideoChat [44] | Self-Collected | 18K | - | - | Des |
| Vid-ChatGPT [56] | ActivityNet-200 [6] | 100K | - | - | Des |
| nuScences-QA [61] | nuScenes [7] | 460K | 5.5 | 2 | Des |
| DriveGPT4 [82] | BDD-X [38] | 28K | 77 | 4 | Des |
| LLM-driver [11] | Self-Collected | 160K | - | - | Des |
| **ELM (Ours)** | nuScenes [7] | 7.4M | 5.5 | 2 | Des, Loc |
|  | Waymo [70] | 450K | 6.4 | 3 | Des |
|  | YouTube | 1.1M | 1474 | 700 | Des |
|  | Ego4D [25] | 300K | 1638 | 74 | Des |

(a) Global Distribution

(b) in USA    (c) in China

Discarded

## Annotation - Quality check



### Location Labeling

Prompt → 1. GPT-4 → Meta Template → 2. Sampling → Inspector → 3. Select → Template → Generated Data

Point Cloud (x, y, z)

Pixel <u, v>

Data Batch → 4. Sampling & Filling → Generated Data

**Generated Data**
Q: What are the 3D spatial coordinates of the pixel at <Pixel Location>
A: Location: x ,y, z, Category: Car

### Description Labeling

nuScenes, Waymo, Youtube, Ego4D

**Example Image**

**Refined Caption**
The driving scene in the image **features a busy city street** with···
The **traffic light is green**··· ✓
The ego-vehicle **should proceed with caution and**···

Manual Revision ↑ 8.

World Data → 1. Crawl → Raw Data
2. Sampling → 3. Select (Quality Check)
Discard → Blacklist

Raw Data → Data Batch → 4. Feed → Vision-Language model → 5. Generate → Caption Batch → 6. Sampling → 7. Generated Caption → Discard (Quality Check)

**Feedback**
The traffic light is **red**··· ✗ → Rerun

## Examples of Description Labels



**Raw:** The scene is a busy city street with a red traffic light, a McDonald's restaurant, and a crosswalk. The traffic elements include traffic lights, cars, and pedestrians. The ego-vehicle is **traveling in the left lane**, and the traffic light is red···

**Refined:** The scene is a busy city street with a McDonald's restaurant, and a crosswalk. The traffic elements include traffic lights, cars, and pedestrians. The ego-vehicle is traveling in the **right** lane, and the traffic light is red···

**Q: How should the ego-vehicle drive at this point?**

A: The ego-vehicle should follow the traffic light's instructions and **wait for the light to turn green before proceeding.**

**Raw:** The driving scene in the image features a car driving under a bridge, with a truck carrying a load of wood in the background. The car is **traveling in the left lane**, and it appears to be moving in the same direction as the truck.

**Refined:** The driving scene in the image features a car driving under a bridge, with a truck carrying a load of wood in the background. The car is traveling in the **middle** lane, and it appears to be moving in the same direction as the truck.

**Q: Are there any potential dangers in the scenario?**

A: The ego-vehicle should continue **driving through the intersection**, following the traffic light's instructions. It is important to **be cautious and attentive to the pedestrians to ensure safety for all involved.**
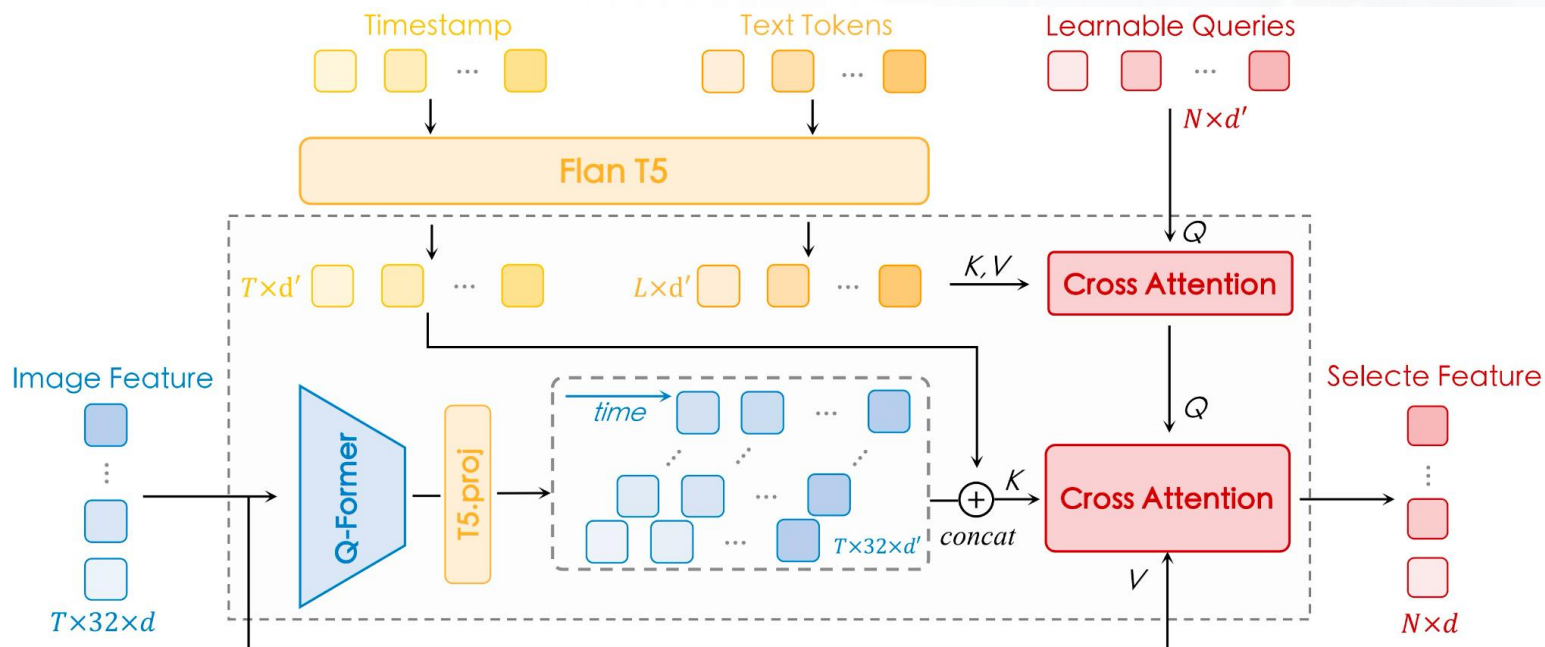
OpenDriveLab

# Embodied Understanding

- Utilize text prompts as guidance to select tokens wisely.
- Alignment among text, timestamps, and videos.



Examples of Selection

Experiments

# Embodied Understanding  Benchmark

- Performing tasks for embodied understanding .

| Tasks | Fine-tune Dataset | Capability | | | | Statistics | | |
|---|---|---|---|---|---|---|---|---|
| | | Description | Localization | Memorization | Forecasting | S($m$) / R($m$) | T($s$) / F | # |
| Surrounding Narration | nuScenes [7] | ✓ | | | | 30 / 5 | 0.5 / 1 | 142K |
| Traffic Sign Inquiry | | ✓ | | ✓ | | 30 / 1 | 3.5 / 7 | 20K |
| Action & Decision | | ✓ | | | ✓ | 30 / 5 | 3.5 / 7 | 301K |
| Box Detection | | | ✓ | | | 50 / 1 | 0.5 / 1 | 232K |
| Tracking | | | ✓ | ✓ | | 50 / 1 | 3.5 / 7 | 131K |
| Box Prediction | | | ✓ | | ✓ | 50 / 1 | 3.5 / 7 | 133K |
| Egocentric Narration | Ego4D [25] | ✓ | | | | 20 / 3 | 3 / 1 | 357K |
| Moment Recap | | ✓ | | ✓ | | 20 / 3 | 60 / 20 | 70K |
| Event Query | | ✓ | | ✓ | | 20 / 3 | 60 / 20 | 70K |
| Activity Prediction | | ✓ | | | ✓ | 20 / 3 | 60 / 20 | 69K |



**Surrounding Narration**
Q: What is the **status** of the cars that are **to the front of** the ego car?
BLIP2-T5: Three cars are parked.
Ours: Many cars are parked and one is moving.
GT: Many cars are parked and many are moving.

**Traffic Sign Inquiry**  $t-n$  $t$
Q: Has the ego vehicle seen any **traffic sign** before?
BLIP2-T5: There is no traffic sign in the scene.
Ours: The ego vehicle has seen 1 go_straight, and 1 turn_right before.
GT: The ego vehicle has seen 1 go_straight, and 1 turn_right before.

**Box Detection**
Q: Provide the 3D location of the **2D point <c, 856, 481>** in the image.
BLIP2-T5: Location: [2.20, 8.85, -0.03], Car
Ours: Location: [3, 12, 0], Car
GT: Location: [3.38, 13.56, 0.003], Car

**Action & Decision**
Q: What **actions** taken by the ego vehicle can lead to a collision with **<c, 817.5, 668.3>**
BLIP2-T5: Sharp left turn.
Ours: Accelerate and go straight.
GT: Accelerate and go straight.

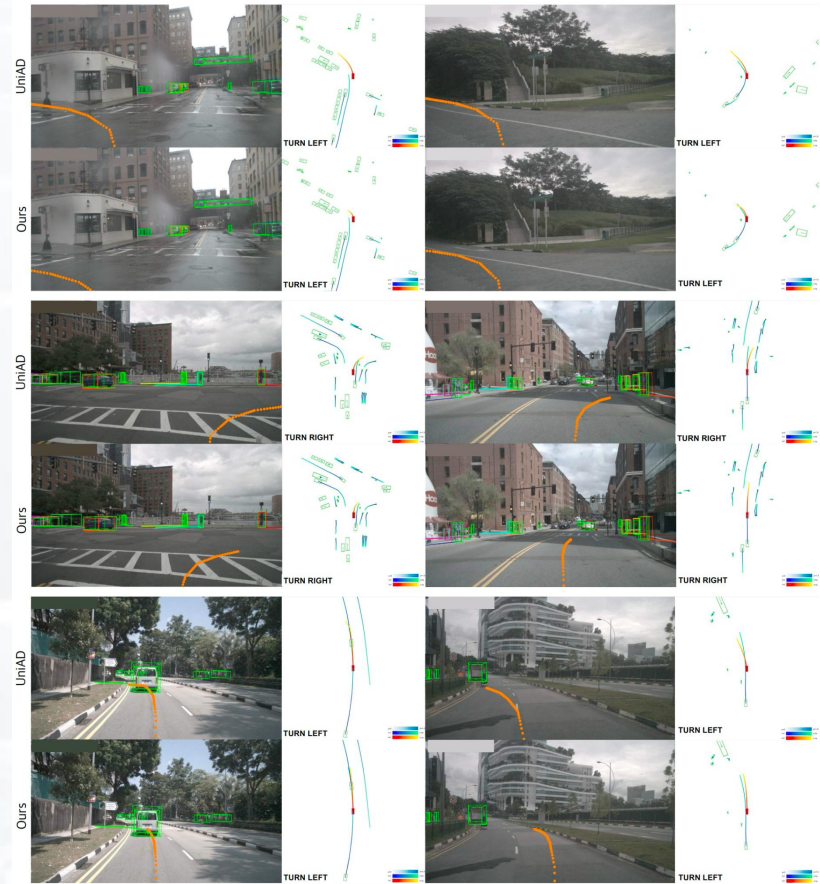**Tracking**  $t-2s$  $t$
Q: What are the 3D spatial coordinates of the pixel at **<c, 1001.1, 544.6>** **2.0 seconds ago?**
BLIP2-T5: Location: [-2.31, 19.62, -1.21], Barrier
Ours: Location: [-1, 15, 0], Pedestrian
GT: [-0.35, 15.32, -0.01], Pedestrian

**Box Prediction**  $t$  $t+3s$
Q: Determine the 3D location in the scene of the 2D pixel at **<c, 865.3, 505.9> 3 seconds later?**
BLIP2-T5: Location: [1.18, 29.96, 1.56], Car
Ours: Location: [0, 38, 0], Pedestrian
GT: Location: [-0.12, 38.45, 0.05], Pedestrian

**Memory Recap**  $t-12.1s$  $t$
Q: "What **happened 12.1 seconds before?**
BLIP2-T5: C stirs the meat in the frying pan
Ours: C moves the meat from the frypan to the plate
GT: C moves the meat from the frypan to the plate

## End-to-end Planning

**Benchmark**



Q: What is the object **to the front of** the ego car?

BLIP2-T5: One bicycle is parked.
Ours: One bicycle is without a rider.
GT: One bicycle is without a rider.

*t − n*      *t*

Q: Has the ego vehicle seen any **traffic sign** before?

BLIP2-T5: Yes. The ego vehicle has seen 1 turn_right before.
Ours: Yes. The ego vehicle has seen 1 turn_left before.
GT: Yes. The ego vehicle has seen 1 turn_left before.

Queried Object

Q: Find the 3D position in the scene of the 2D pixel at **<c, 737.3, 472.7>.**

BLIP2-T5: Location: [-0.2, 25.2, 1.3], truck
Ours: Location: [-2, 32, 1], bus
GT: Location: [-2.16, 34.06, 0.77], bus

Q: What is the future state of **<c3,826.7,521.7>?**

BLIP2-T5: Still.
Ours: Moving forward.
GT: Moving forward.

*t − 1.5s*      *t*

Queried Object

Q: Find the 3D world position of the 2D point **<c, 626.9, 558.3>** in the image. 1.5 seconds ago.

BLIP2-T5: Location: [-0.87, 9.2, -0.3], Car
Ours: Location: [-1, 11, 0], Car
GT: Location: [-1.03, 11.17, -0.52], Car

*t*      *t + 1s*

Queried object

Q: Calculate the 3D world location of the pixel with 2D coordinates **<c, 734.5 , 673.2>.** 1.0 seconds later.

BLIP2-T5: Location: [-0.2, 4.15, -1.3], Car
Ours: Location: [0, 8, 0], Car
GT: [-0.05, 8.00, -0.75], Car

*t − 37.8s*      *t*

Q: "What **happened 37.8 seconds before?**

BLIP2-T5: C pens the fridge
Ours: C puts the water in the fridge
GT: C puts the water in the fridge

## Benchmark



**Event Query**

'C opens the water bottle'    'C closes the water bottle'

Q What happened between 'C opens the water bottle' and 'C closes water bottle'?

BLIP2-T5: C pours water from the water bottle
Ours: C drinks water
GT: C drinks water

'C picks up the drill'    'C puts the power drill on the wood'

Q: What happened between 'C picks up the drill' and 'C puts the power drill on the wood'?
BLIP2-T5: C picks the wood
Ours: C drills the wood
GT: C drills the wood

**Egocentric Narration**

Q: **Give a caption.**
BLIP2-T5: C moves the wires
Ours: C cuts black wire
GT: C cuts black wire

Q: **Give a caption.**
BLIP2-T5: C drops the sieve in the sink
Ours: C pours the potatoes in the plate
GT: C pours the potatoes into the sieve.

**Event Forecasting**

t    t + 0.7s

Q: What will happen in the next 0.7 seconds in the future?
BLIP2-T5: C drops the knife on the cutting board
Ours: C puts the slices of cabbage in the sieve
GT: Cputs the slices of cabbage in the sieve with her left hand

t    t + 1.2s

Q: What will happen in the next 1.2 seconds in the future?
BLIP2-T5: C drops the knife on the cutting board
Ours: C cuts the onion with the knife
GT: C cuts the onion with the knife

## Experiments

### OOD Evaluation

| Methods | Tracking | | Box Detection | | Box Prediction | | Traffic Sign Inquiry | | | Surrounding Narration | | | Action & Decision | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Pr@1 | Pr@2 | Pr@1 | Pr@2 | Pr@1 | Pr@2 | C | R | B | C | R | B | C | R | B |
| BLIP2-opt [27] | 0.1 | 0.1 | 0.1 | 0.2 | 0.2 | 0.5 | 23.0 | 26.9 | 20.5 | 8.1 | 19.7 | 21.2 | 8.4 | 11.5 | 11.1 |
| BLIP2-flant5 [27] | 3.0 | 6.0 | 5.1 | 10.5 | 3.6 | 6.3 | 63.1 | 39.4 | 31.4 | 65.2 | 64.9 | 27.9 | 68.7 | 71.4 | 43.1 |
| LLaMA-Ada. [17] | 6.1 | 10.5 | 8.3 | 14.9 | 7.5 | 12.5 | 68.3 | 66.6 | 61.6 | 67.0 | 77.5 | **60.1** | 72.3 | 76.8 | **64.7** |
| LLaVA [32] | 5.5 | 9.3 | 28.5 | 31.2 | 6.1 | 10.2 | 51.1 | 58.5 | 50.8 | 64.9 | 64.6 | 41.2 | 64.4 | 62.4 | 57.9 |
| Otter [26] | 10.0 | 17.2 | 41.8 | 46.9 | 8.9 | 15.8 | 62.8 | 41.1 | 32.4 | 60.0 | 64.2 | 13.3 | 69.2 | 73.2 | 53.0 |
| VideoChat [28] | 0.4 | 0.9 | 0.1 | 0.3 | 0.1 | 0.2 | 25.3 | 21.9 | 11.7 | 21.7 | 29.2 | 12.2 | 29.6 | 33.2 | 13.1 |
| Vid-ChatGPT [36] | 0.1 | 0.6 | 0.1 | 1.0 | 0.3 | 1.2 | 49.6 | 57.1 | 48.6 | 61.0 | 69.6 | 37.2 | 53.6 | 58.5 | 43.5 |
| **ELM (Ours)** | **14.0** | **23.3** | **51.6** | **56.9** | **15.1** | **24.4** | **76.5** | **71.2** | **63.9** | **73.2** | **78.7** | 29.8 | **74.4** | **83.3** | 41.2 |

(a) **nuScenes.** We outperform the best previous methods on most metrics across the six tasks on nuScenes which validates the generality of our model.

| Method | ADE↓ | FDE↓ | Time(s)↓ |
|---|---|---|---|
| Command Mean | 7.98 | 11.41 | - |
| UniAD-single [34] | 4.16 | 9.31 | 0.56 |
| Flamingo [3] | 2.78 | 5.31 | 1.47 |
| **ELM** | **2.28** | **4.27** | 1.61 |

| Methods | Moment Recap | | | Event Query | | | Egocentric Narration | | | Activity Prediction | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | C | R | B | C | R | B | C | R | B | C | R | B |
| BLIP2-opt [27] | 1.2 | 8.9 | 6.8 | 7.8 | 28.4 | 14.7 | 5.2 | 19.8 | 10.7 | 2.7 | 18.7 | 9.6 |
| BLIP2-flant5 [27] | 13.1 | 31.9 | 12.5 | 27.3 | 33.0 | 16.6 | 16.9 | 33.5 | 15.4 | 11.5 | 31.2 | 11.3 |
| LLaMA-Ada. [17] | 11.2 | 30.2 | 12.3 | 37.5 | **47.2** | 28.1 | 18.4 | 34.2 | 15.3 | 13.1 | 31.2 | 12.8 |
| LLaVA [32] | 9.6 | 28.3 | 12.1 | **39.8** | 44.6 | **29.9** | 6.5 | 28.2 | 11.6 | 8.4 | 28.0 | 13.0 |
| Otter [26] | 11.4 | 29.6 | 10.5 | 27.1 | 38.3 | 19.1 | 14.1 | 31.4 | 13.9 | 11.1 | 29.4 | 10.3 |
| VideoChat [28] | 13.2 | 32.5 | 13.8 | 34.5 | 42.2 | 26.4 | 20.7 | 35.0 | **17.6** | 12.1 | 32.4 | 14.1 |
| Vid-ChatGPT [36] | 10.0 | 31.1 | 13.3 | 27.9 | 36.5 | 20.9 | 10.2 | 21.7 | 10.4 | 9.4 | 30.5 | 12.6 |
| **ELM (Ours)** | **22.6** | **36.7** | **19.4** | 38.0 | 43.1 | 27.6 | **26.5** | **37.7** | 16.9 | **18.1** | **34.1** | **17.0** |

(b) **Ego4D.** We extended the model to Ego4D dataset and verified the generality of our token bank module on four tasks.

| Methods | # param |
|---|---|
| BLIP2-opt | 2.7B |
| BLIP2-flant5 | 2.7B |
| LLaMA-Ada. | 7B |
| LLaVA | 7B |
| Otter | 7B |
| VideoChat | 7B |
| Vid-ChatGPT | 7B |
| **ELM (Ours)** | 2.7B |

(c) **Adopted LLM params.**
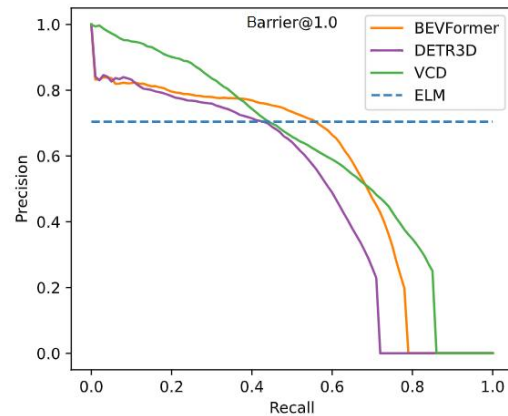
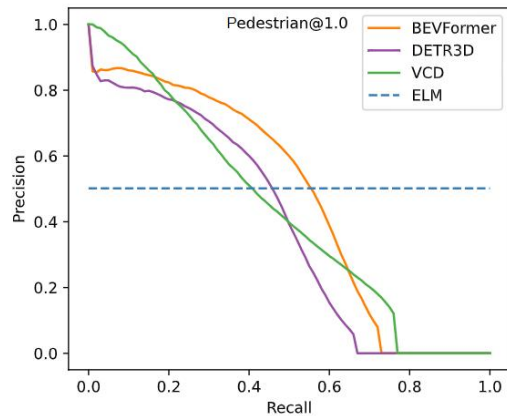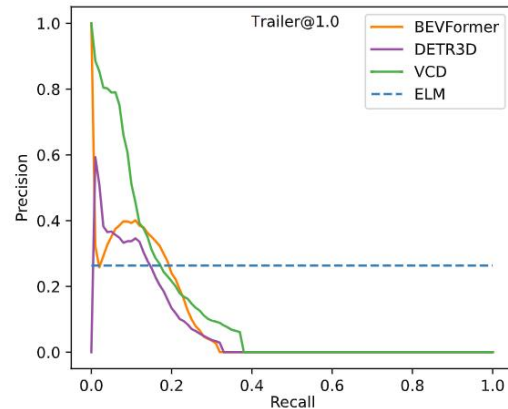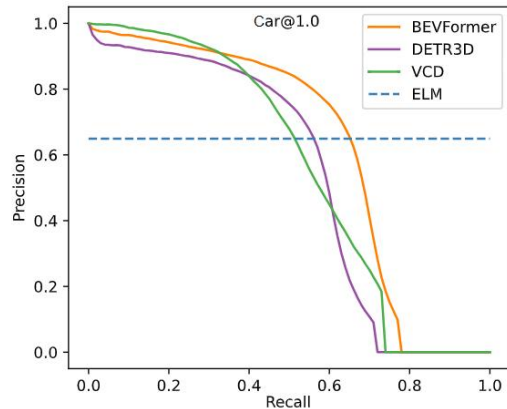Q: What is the unusual about the driving scene?

A: The scene shows a busy city street with a **construction site** on the side of the road. There is a traffic light visible in the background. The unusual element in this scene is **the presence of a construction site, which may cause traffic disruptions or delays.**
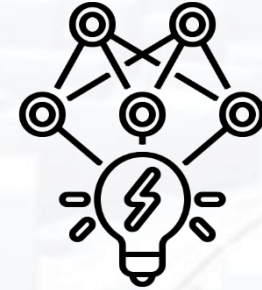
## Experiments

# Embodied Understanding



## Driving-specific Inputs

ELM cannot handle common setting such as LiDAR or multi-view images as input, limiting its information source.

## Closed-loop Planning

ELM is evaluated under an open-loop scheme, while closed-loop planning is necessary to see if it can handle corner cases.

## Efficiency Constraints

Inheriting the drawbacks of LLMs, ELM suffers from long inference time, which may impact practical implementation.

OpenDriveLab

# One-page Takeaway

- End-to-end Autonomous Driving
  - Challenge: **Generalization & Explainability**
  - Recent trend: use vision language model to **embed "world knowledge"** to solve the challenges.

- ELM: Embodied Understanding of Driving Scenarios
  - Revive driving scene understanding by delving into **embodied** settings, along with capacities, tasks, and rubrics.
  - Expand vanilla VLMs to process long horizon **space** and **time** (open-world data & module design).

END