Open AriveLab



端到端自动驾驶:前沿与挑战 End-to-end Autonomous Driving: Challenges and Frontiers

Li Chen OpenDriveLab at Shanghai Al Lab June 10, 2024

Outline

• 课程目标

- 掌握端到端自动驾驶系统的特点与基础方法
- 了解端到端自动驾驶算法的典型问题与挑战
- 了解端到端自动驾驶前沿算法





Outline

- 端到端自动驾驶系统概述 / Introduction
 - 背景与动机 / Motivation
 - **发展历程** / Roadmap
- 基础方法 / Method
- 验证与评测 / Benchmarking
- 关键研究内容与挑战 / Frontiers and Challenges
- 未来可能与讨论 / Future Trends and Discussions
- Q&A







端到端自动驾驶系统概述 / Introduction

Autonomous Driving (AD) Tasks



Challenge | Various weathers, illuminations, and scenarios







- Independent teams for module development
 - Dataset friendly
 - Quantitatively evaluation for intermediate tasks
 - Great interpretability
- Parallel onboard deployment

Error accumulation, information loss. **Results, instead of features** are traversed across modules.

- Openset problems
- Long-tail problems

Open 🔁 rive Lab

Motivation | Why End-to-end Autonomous Driving?



End-to-end autonomous driving system - A suite of fully differentiable programs that:

- take raw sensor data as input
- produce a plan and/or low-level control actions as output



Motivation | Why End-to-end Autonomous Driving?

Advantages

- + Simplicity in combining all modules into a single model that can be joint trained
- + Preventing cascading errors in modular design
- + Directly optimized toward the ultimate task, planning / trajectory prediction
- + Computational efficiency (all shared backbone), production-level friendly





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

Open PriveLab

Trending | End-to-end Autonomous Driving

No hard-code



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

E2E Robot

Industry



Tesla Optimus 🤣 🖬 @Tesla_Optimus · Sep 24 Optimus can now sort objects autonomously 🧝



Open sriveLab

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

Ashok Elluswamy

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!

Selon Musk ⊘ ⊠ @elonmusk - Aug 26 twitter.com/i/broadcasts/1...



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.

My Opinion

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



Trending | End-to-end Autonomous Driving

Loss function

World state

Neural network

And many others ...

NVIDIA

Perception

Loss function



Synthetic / Retrieval

World state

fully know

artists

Sensor

data

sim sensors

physics based

Driving Input, 10⁸ dimensions



Q GNSS

complex physics, (AI) traffic models

Loss function Loss function

Planning

Control

AV stack

Prediction

Basic Sat-nav Map

Vehicle State

+ other sensing modalities where required, e.g. RADAR



Decoded human-interpretable intermediate representations



Semantics, geometry, motion prediction.



Industry

comma.ai

- Openpilot is an open source driver assistance system.
- Openpilot performs the functions of Automated Lane Centering (ALC) and Adaptive Cruise Control (ACC) for 250+ supported car makes and models.



https://arxiv.org/abs/2206.08176





Open AriveLab

Trending | End-to-end Autonomous Driving Academia / Hybrid EBERHARD KARLS UCLA 9 UNIVERSITAT TÜBINGEN MetaDrive, ACO, CAT, etc Transfuser, KING, Misconceptions, DA-RB, etc. **USC** XAS LBC, WoR, LAV, etc. GPT-Driver, Agent-Driver, etc The University of Texas at Aust Uber UNIVERSITY OF NMP, P3, MP3, UniSim (sort of) GenAD, etc Berkeley TORONTO Шaabi BOSTON 清莱大学 LbW, SelfD, CaT, AnyD, etc. SparseDrive, ADAPT, etc UNIVERSITY **ETH** zürich Roach, etc **Open PriveLab** THE HONG KONG UNIVERSITY OF SCIENCE AND TECHNOLOGY MMFN, Carl-lead, PMP, etc Work from our Team introduced later GRI Horizon Robotics VAD, VAD v2

Open riveLab

Roadmap | End-to-end Autonomous Driving



Summary (1/2)

- CARLA leaderboard gets much improved over the years. With new mapping / scenarios (leaderboard 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



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





Public Opinions on Our Survey

Paper https://arxiv.org/pdf/2306.16927.pdf

Repo (paper collection)
<u>https://github.com/OpenDriveLab/End-to-end-Autonomous-Driving</u>

Alex Kendall 🤣 @alexgkendall

This is a fantastic, comprehensive and forward-looking survey of academic literature about end-to-end machine learning for autonomous driving. It is a very timely publication as the field is exploding with interest right now.

I'm aligned with the paper's conclusions on open algorithmic challenges. There's loads of insight around opportunities like world modelling, language, foundation models and long-tail robustness. This paper also exposes that academic literature under-appreciates significant industry challenges right now, such as (1) safety, reward modelling and policy alignment against human expectations and risk, or (2) the significance of establishing a synthetic/real-world data engine for training/validation, which are critical to the success of any machine learning system. I'd love to see more work in these areas.

Great to see @AutoVisionGroup @francislee2020, well done!

Awesome Vision Group @AutoVisionGroup · Sep 18



Yann LeCun 🤣 🙉

A nice survey of end-to-end learning methods for autonomous driving.

Sep 18 Awesome Vision Group @AutoVisionGroup · Sep 18

Why are Tesla @elonmusk and Wayve @alexgkendall @Jamie_Shotton moving towards end-to-end autonomous driving? What is the state-of-the-art in this field? With our friends @francislee2020 we recently wrote an extensive survey paper on this emerging topic: arxiv.org/abs/2306.16927

202339B 端到端自动驾驶

Original 吴双 **吴言吴语** 2023-10-02 05:42

收录于合集 #自动驾驶

20个 >

这周我们读一篇提交到PAMI的端到端自动驾驶的综述论文: SUBMITED TO IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE, JUNE 2023

End-to-end Autonomous Driving: Challenges and Frontiers

(今 吴言吴语 Li Chen, Penghao Wu, Kashyap Chitta, Bernhard Jaeger, Andreas Geiger and Hongyang Li

Arxiv链接: https://arxiv.org/abs/2306.16927

可以看到这篇文章在六月份,好像是CVPR会议期间就挂到了arxiv上,当时眼前一亮随 手放在了桌面,结果回头就忘了,最近SS兄提醒,就给自己安排了周末作业。由于论文 覆盖的内容很多,今天就只聊一聊我个人看到的值得注意或者觉得需要强调的点。

总结:很好的综述,值得看看。

Join Slack Discussions!

...

https://join.slack.com/t/opendrivel ab/shared_invite/zt-244lgu87b-eL onLQzle4wRkg8W8WOUlg

Open riveLab

[1] Chen et al. End-to-end Autonomous Driving: Challenges and Frontiers. arXiv, 2023.





基础方法 / Method

Types of Learning

Supervised Learning:

- Dataset: $\{(x_i, y_i)\}$ (x_i = data, y_i = label) Goal: Learn mapping $x \mapsto y$
- Examples: Classification, regression, imitation learning, affordance learning, etc.

Unsupervised Learning:

- Dataset: $\{(x_i)\}(x_i = data)$ Goal: Discover structure underlying data
- Examples: Clustering, dimensionality reduction, feature learning, etc.

Reinforcement Learning:

- Agent interacts with environment which provides numeric reward design
- Goal: Learn how to take actions in order to maximize reward
- Examples: Learning of manipulation or control tasks (everything that interacts)



Types of Learning for E2E AD



Behavior Cloning

Inverse Optimal Control

Reinforcement Learning

Imitation Learning: Learn a policy from expert demonstrations

- Expert demonstrations are provided
- Supervised learning problem
- Behavior Cloning, Inverse Optimal Control (Inverse Reinforcement Learning)

Reinforcement Learning: Learn a policy through trial-and-error

- No expert demonstrations given
- Agent discovers itself which actions maximize the expected future reward
 - The agent interacts with the environment and obtains reward
 - \circ The agent discovers good actions and improves its policy π



Imitation Learning



Credit to Andreas Geiger, Lecture: Self-Driving Cars



Imitation Learning in a Nutshell



Hard coding policies is often difficult \rightarrow Rather use a data-driven approach!

- **Given**: Demonstrations or demonstrator
- **Goal**: Train a policy to mimic decision
- Variants: Behavior Cloning, Inverse Optimal Control (Inverse Reinforcement Learning)



Imitation Learning: Trajectory or Control?



Control

CIL/CILRS/WoR/...

- End-to-end optimization; elegant structure
- Utilization of advanced RL methods
- Focus on current timestamp; incontinuous output
- Lack of intention-related information
- Coupling with vehicle dynamics, challenging to generalize



- Ctl: Collision Traj: Success
- However, no clear performance gap in existing research papers we think.
- While, trajectory may be more appropriate for practical application.

[1] Wu et al. Trajectory-guided Control Prediction for End-to-end Autonomous Driving: A Simple yet Strong Baseline. NeurIPS, 2022.

Open 🔁 rive Lab

Imitation Learning - Behavior Cloning

Behavior cloning makes I.I.D. assumption

- Next state is sampled from states observed during expert demonstration
- Thus, next state is sampled independently from action predicted by current policy

What if π_{ρ} makes a mistake?

- Enters new states that haven't been observed before
- New states not sampled from same (expert distribution anymore)
- Cannot recover, catastrophic failure in the worst case

What can we do to overcome this train/test distribution mismatch?





Imitation Learning - Behavior Cloning



Policy

Critical

States

Train

Dataset

Replay Buffe

Rollout

Sample

What can we do to overcome this train/test distribution mismatch? - Data Aggregation (**DAgger**)

- Iteratively build a set of inputs that the final policy is likely to encounter based on previous experience. Query expert for aggregate dataset
- But can easily overfit to main mode of demonstrations
- High training variance (random initialization, order of data)



- Sample critical states from the collected on-policy data based on the utility they provide to the learned policy in terms of driving behavior
- Incorporate a replay buffer which progressively focuses on the high uncertainty regions of the policy's state distribution

Ross, et al. A Reduction of Imitation Learning and Structured Prediction on No-Regret Online Learning. AISTATS, 2011.
Prakash, et al. Exploring Data Aggregation in Policy Learning for Vision-based Urban Autonomous Driving. CVPR, 2020.
Credit to Andreas Geiger, Lecture: Self-Driving Cars

Environment

On-Policy

Data

Open 🗛 riveLab

Imitation Learning - Conditional Behavior Cloning



Idea:

- Condition controller on **navigation command** $c \in \{\text{left, right, straight}\}$
- High-level navigation command can be provided by consumer GPS, i.e., telling the vehicle to **turn left/right** or **go straight** at the next intersection
- This removes the task ambiguity induced by the environment
- State s_+ : current image Action a_+ : steering angle & acceleration

[1] Codevilla et al. End-to-end Driving via Conditional Imitation Learning. ICRA, 2018.



Imitation Learning - Inverse Optimal Control



Inverse Optimal Control (Inverse Reinforcement Learning):

- Agent observes environment state *s*, at time †
- Agent sends action a_t at time t, based on the cost/reward function, to the environment
- Environment returns the new state s_{t+1} to the agent



Imitation Learning - Inverse Optimal Control



Objective cost function:

- Safety Cost: not collide with other detected objects within future periods; not overlap with road boundaries; maintain a safe distance at high velocity
- **Comfort and Progress**: penalize large lateral acc., jerk, or curvature; reward forwarding to designated directions
- Learned Cost Volume: unspecified terms → emergent ability?

[1] Hu et al. ST-P3: End-to-end Vision-based Autonomous Driving via Spatial-Temporal Feature Learning. ECCV, 2022.



Imitation Learning - Summary

Advantages of Imitation Learning, esp. Behavior Cloning

- **Easy** to implement
- Cheap annotations (just driving while recording visual sensor and actions)
- Entire model trained **end-to-end**
- Conditioning removes ambiguity at intersections

Challenges of Imitation Learning, esp. Behavior Cloning

- Behavior cloning uses I.I.D. assumption which is violated in practice
- No memory (can't remember speed signs, etc.)
- Mapping is difficult to interpret ("**black box**"), despite visualization techniques
- More discussion later.



Reinforcement Learning



Reinforcement Learning:

- Agent observes environment state *s*₊ at time †
- Agent sends action a₊ at time t to the environment
- Environment returns the reward r_{t} and its new state s_{t+1} to the agent



Reinforcement Learning



Cart Pole Balancing



Robot locomotion (in mujoco)

- **Objective:** Balance pole on moving cart
- State: Angle, angular vel., position, vel.
- Action: Horizontal force applied to cart
- **Reward**: 1 if pole is upright at time t

State: Position and angle of joints

Action: Torques applied on joints



Atari Games

- **Objective**: Maximize game score
- State: Raw pixels of screen (210x160)
- Action: Left, right, up, down
- Reward: Score increase/decrease at t

Objective: Make robot move forward **Reward**: 1 if upright & forward moving 0462

Self-driving (in gym/CarRacing)

- **Objective:** Lane following
- State: Image (96x96) •

•

- Action: Acceleration, Steering •
- Reward: per frame, + per tile .

Credit to Andreas Geiger, Lecture: Self-Driving Cars



Reinforcement Learning

Advantages of Reinforcement Learning

- Straightforward idea, early attempts
- Easy be trained with privileged simulator information
- Exploration & Exploitation relieve causal confusion

Challenges of Reinforcement Learning

- Large search/state-action space (π) , especially in outdoor driving scenarios; long training times
- Deep Q-Learning
 - \circ Uniform sampling from replay buffer \rightarrow all transitions equally important
 - Simplistic exploration strategy
 - Action space is limited to a discrete set of actions (otherwise, expensive test-time optimization required)
- Trained in simulation mainly. **Sim2Real gap**.
- More discussion later.







验证与评测 / Benchmarking

Benchmarking

Real-world



On-field tests (c.f. Waymo)

Online/Closed-loop



nuPlan

Route Completion * **T** Infraction Penalty

The world's first benchmark for autonomous veh planning



Driving Score =

the virtual world.

Metric:

•

Autonomous Racing

Mcity, UMich





Offline/Open-loop

NUSCENES

Metric:

- L2 = | GT Pred |
- Collision rate

Open 🕰 riveLab

Benchmarking - Closed-loop Simulation

Parameter Init.

Procedural Generation

- Sample probabilistic distribution of simulation properties with algorithms
- Hand-tuned rules & params.

Data-Driven

- Sample from logs
- Generate by models



Traffic Sim.

3-Point Turn U-Turn Non-Compliant

Complex Interactions

Rule-Based

 Intelligent Driver Model (IDM) (proves to be effective)

Data-Driven

- Generate by models
- Waymo Sim Agent

Sensor Sim.

Graphic-Based

- 3D models with assets, plus physical rendering process and sensors modeling
- Physical occlusion, shadows, reflections, etc.

Data-Driven

- NeRF & 3D Gaussian Splatting
- GAN & Diffusion



Vehicle Dyn. Sim.



Simplified Vehicle Model

- Unicycle / Bicycle model
- Multi-body system

Data-Driven

 NN-based learned from data



Benchmarking - Open-loop Evaluation

Current prevailing end-to-end autonomous driving methods commonly use **nuScenes** for **open loop evaluation of their planning behavior.**

However:

• **NuScenes** dataset, characterized by **relatively simple driving scenarios**, leads to an underutilization of perception information in end-to-end models.



[1] Li et al. Is Ego Status All You Need for Open-Loop End-to-End Autonomous Driving? CVPR, 2024.



EgoStatus | Motivation

Current prevailing end-to-end autonomous driving methods commonly use **nuScenes** for **open loop evaluation of their planning behavior.**

However:

- **NuScenes** dataset, characterized by **relatively simple driving scenarios**, leads to an underutilization of perception information in end-to-end models.
- **AD-MLP** paper recently points out that a simple MLP network can also achieve state-of-the-art planning results, **relying solely on the ego status information**.


EgoStatus | Motivation

Current prevailing end-to-end autonomous driving methods commonly use **nuScenes** for **open loop evaluation of their planning behavior.**

However:

- **NuScenes** dataset, characterized by **relatively simple driving scenarios**, leads to an underutilization of perception information in end-to-end models.
- **AD-MLP** paper recently points out that a simple MLP network can also achieve state-of-the-art planning results, **relying solely on the ego status information**.

Is Ego Status All You Need for Open-Loop End-to-End Autonomous Driving?



EgoStatus | Conclusion

- The planning performance of existing open-loop autonomous driving models based on nuScenes is highly affected by ego status
- Existing planning **metrics** fall short of fully capturing the true performance of models.

The development of more appropriate *datasets* and *metrics* represents a more critical and urgent challenge to tackle

[1] Li et al. Is Ego Status All You Need for Open-Loop End-to-End Autonomous Driving? CVPR, 2024.







关键研究内容与挑战 / Frontiers and Challenges

Key Challenges



Modalities





Dilemma over Sensing and Input Modalities



Early Fusion

- Combine sensory information before feeding it into the feature extractor
- Concatenation (imgs, or point-painting)

Middle Fusion

- Separately encode inputs and then combining them at the feature level
- Transformer attention
- Fusion under BEV space

Late Fusion

• Combine multiple results from multi-modalities

Challenge: Various sensors possess distinct perspectives, data distributions (similar issues in perception tasks). Some unique inputs such as vehicle states and navigation signals exist for E2E driving.

Dilemma over Sensing and Input Modalities - Language



Instruction: Head upstairs and walk past the piano through an archway directly in front. Turn right when the hallway ends at pictures and table. Wait by the moose antlers hanging on the wall.

Vision Language Navigation, CVPR'18



LINGO-2, arXiv'24



Insight from 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

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

- 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?

Dilemma over Sensing and Input Modalities - Language

Outdoor Navigation



Turn and go with the flow of traffic. At the first traffic light turn left. Go past the next two traffic light, As you come to the third traffic light you will see a white building on your left with many American flags on it. Touchdown is sitting in the stars of the first flag.

TOUCHDOWN, CVPR'19

- Discrete action space (Forward/Left/Right/Stop)
- Crowdsourcing routes or collected from Google Maps

Linguistic Instruction



- **Grounding** objects/areas with human advice, and predict control signals with attention
- Human annotated videos



- Encode perception/prediction/etc results into texts (or tokens) and **prompt with LLMs** for planning and behavior explanation
- OR: construct visual question answering dataset to train VLMs

Existed Opportunity: Language can provide certain **high-level instructions**, enhancing human-machine-interaction. **Current Challenge**: LLMs (or VLMs) require **long inference** time, lack of **quantitative accuracy** and stability of outputs.

Dependence on Visual Abstraction



Key function:

- Compact intermediate representations
- Pre-trained visual encoders accelerate downstream training, especially for RL methods (sample efficiency)



Dependence on Visual Abstraction - Representation Design

Representation forms of visual inputs

- CNN vs Transformer:
 - CNN still dominates in E2E driving than ViT. Why?
- Bird's-eye-view (BEV) & Occupancy grid:
 - BEV is good for sensor/temporal fusion and facilitates downstream tasks
 - 3D occupancy can capture irregular objects and used for collision avoidance
- Maps:
 - BEV segmentation, vectorized lanelines, centerlines and topology, lane segments, etc.



Transfuser++, ICCV'23



Scene as Occupancy, ICCV'23



VAD, ICCV'23

However, if explicit representations are necessary for performance improvement is unclear. This is a debate about modular E2E vs. scaling up "black box" models.

Open 🔁 rive Lab

Dependence on Visual Abstraction - Representation Learning



Current methods first pre-train the visual encoder of the network using proxy pre-training tasks.

There inevitably exist possible **information bottlenecks** in the learned representation, and redundant information unrelated to driving decisions may be included.



Dependence on Visual Abstraction - Representation Learning

You Tube

ACO, ECCV'22

Traditional RL methods' pretraining: abstracted information, also helpful for Sim2Real

- Semantic segmentation
- Depth map
- VAE (Variational autoencoder)

Pre-training from large-scale vision data







Self-supervised Geometric Modeling

(a) Self-supervised Visuomotor Policy Pre-training

(b) Downstream Tasks



PPGeo, ICLR'23

Open 🔁 rive Lab

95

Representation Learning - ViDAR



Method	Detection		Tracking			Map	oping	Motion Forecasting			Future Occupancy Prediction				Planning	
	NDS \uparrow	$mAP\uparrow$	AMOTA↑	AMOTP↓	IDS↓	IoU-lane↑	IoU-road↑	minADE↓	minFDE↓	MR↓	IoU-n.↑	IoU-f.↑	VPQ-n.↑	VPQ-f.↑	avg.L2↓	avg.Col.↓
UniAD	49.36	37.96	38.3	1.32	1054	31.3	69.1	0.75	1.08	0.158	62.8	40.1	54.6	33.9	1.12	0.27
ViDAR	52.57	42.33	42.0	1.25	991	33.2	71.4	0.67	0.99	0.149	65.4	42.1	57.3	36.4	0.91	0.23

By fine-tuning on UniAD, ViDAR boosts perception, prediction, and planning at the same time by a large margin, providing a new solution for using large-scale unlabeled data.

Complexity of World Modeling





Complexity of World Modeling

A Path Towards Autonomous Machine Intelligence Version, Yann Lecun

Task / Objective:

- Represent the world & Learn to predict and re-act
 - Simulate the world without **REAL** interaction with the world.





configurator

Short-tern

Trending: Recent Work on World Model







Trending: Recent Work on World Model



World model to generate videos of the driving scenario. **Then what?** Is it useful for downstream tasks? (To be validated)

Complexity of World Model



States Cost / Reward Success/Fail Ego agent --**RL Gyms** Other objects (static) Intermediate Reward --**Background environment** -**Ego-vehicle** Collision --Autonomou S Other vehicles, pedestrians, Comfort -_ Driving cyclists, etc (moving) Forward -**Background environment** -- etc Hard to define! **Complicated!** A video predictor?

Complexity of World Model

Forms of world model

- Images
- Point cloud (LiDAR) and Occupancy
- BEV maps
- Integrated Motion Prediction and Planning
- Latent model (LSTM/GRU-based, etc)

Application of world model

- Reward for RL/Sampling/etc
- Decode planning with Inverse Dynamics Model
- Serve as pre-training (fine-tuned for planning)



Complexity of World Model | GenAD

Summary: Training a **billion-scale video prediction model** on **web-scale driving videos**, to enable its **generalization across** a wide spectrum of **domains and tasks**.



Complexity of World Model | GenAD

- GenAD (5.9B) = SDXL (2.7B) + Temporal Reasoning Blocks (2.5B) + CLIP-Text (0.7B)
- Tuning the image generation model (SDXL) into a highly-capable video prediction model



Tasks | Zero-shot Generalization (Video Prediction)



Zero-shot video prediction on unseen datasets including Waymo, KITTI and Cityscapes



Tasks | Language-conditioned Prediction







"Drive slowly down at intersection, several barriers beside the road"





"Turn right, some parked cars, a parking lot"



Tasks | Action-conditioned Prediction (Simulation)

Method	Condition	nuScenes Action Prediction Error (↓)				
Ground truth	-	0.9				
GenAD	text	2.54				
GenAD-act	text + traj.	2.02				

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

Simulate the future differently conditioned on **future trajectory.**





Tasks | Planning



Mathad	# Trainable	nuScenes				
Method	Params.	ADE (\downarrow)	FDE (\downarrow)			
ST-P3* [20]	10.9M	2.11	2.90			
UniAD* [22]	58.8M	1.03	1.65			
GenAD (Ours)	0.8M	1.23	2.31			

Table 5. Task on Planning. A lightweight MLP with *frozen* GenAD gets competitive planning results with $73 \times$ fewer trainable parameters and front-view image alone. *: multi-view inputs.

Training process **speeds up by 3400 times** compared to UniAD (CVPR Best Paper).



Vista: Generalized action conditions



Control with traj & angles (translated to commands for vis)



Forward

Left

Right

Stop

[1] Gao et al. Vista: A Generalizable Driving World Model with High Fidelity and Versatile Controllability. arXiv, 2024.





ViDAR - World Model in Driving

The First Multimodal World Model

Visual Inputs



ViDAR - World Model in Driving

Summary: Training multimodal world model by **Visual Point Cloud Forecasting** and boosting **End-to-End Autonomous Driving**.



ViDAR | Future Prediction Experiments







ViDAR | Different Ego Control Experiments



Go

Turn Right





ViDAR | Downstream Experiments

Method	Detection		Tracking			Map	ping	Motion Forecasting			Future Occupancy Prediction				Planning	
	NDS \uparrow	mAP↑	AMOTA↑	AMOTP↓	IDS↓	IoU-lane↑	IoU-road↑	minADE↓	minFDE↓	MR↓	IoU-n.↑	IoU-f.↑	VPQ-n.↑	VPQ-f.↑	avg.L2↓	avg.Col.↓
UniAD	49.36	37.96	38.3	1.32	1054	31.3	69.1	0.75	1.08	0.158	62.8	40.1	54.6	33.9	1.12	0.27
ViDAR	52.57	42.33	42.0	1.25	991	33.2	71.4	0.67	0.99	0.149	65.4	42.1	57.3	36.4	0.91	0.23

By fine-tuning on UniAD,

ViDAR boosts perception, prediction, and planning at the same time by a large margin, providing a new solution for using large-scale unlabeled data.





Reliance on Multi-task Learning

Multi-task learning (MTL): Jointly perform several related tasks based on a shared representation through separate branches/heads.

Pros

- Significant computational cost reduction
- Related domain **knowledge is shared** within the shared model

Challenges

- The **optimal combination** of auxiliary tasks and the appropriate weighting of their losses
- Construct large-scale datasets with multiple types of aligned and high-quality annotations

Typical multi-tasks in E2E driving:

- Early works:
 - Semantic segmentation: high-level understanding
 - **Depth estimation**: 3D geometry
- Recent works:
 - 3D object detection
 - BEV segmentation
 - Visual affordance, including traffic light states, distance to opposite lanes, etc



Inefficient Experts and Policy Distillation

The popular "Teacher-Student" IL Paradigm



- Ease training and better generalization (with privileged information)
- Can be queried with any states, instead of logged states only



Inefficient Experts and Policy Distillation

• Expert: Ground Truth (GT) to action

Gap

• Student: Img to action

- What for or How to Distillation
 - Critical features
 - Input gap Casual confusion

Expert (by RL/IL/hand-rule, gt input):

• Not/Can't perfect, even for a certain benchmark

Method	Input	Driving Score \uparrow
Transfuser [39, 8]	Camera + LiDAR	31.0
LAV [3]	Camera + LiDAR	46.5
Student Model + Frozen Roach	Camera + LiDAR	8.9
Roach [55]	Privileged Info.	74.2
Roach + Rule [50]	Privileged Info.	87.0

Student (IL etc, sensor input):

• Not/Can't perfect, even with large-scale data and have visually great representations





BEVFusion + Mask2Former 2M training data

Privileged Input

Perception Result

DriveAdapter

How to balance the efficiency and causal reasoning ability?



Utilize the strong RL-based privileged teacher model!

- Train a Teacher Model for Planning by RL
- End-to-End Connected by Adapter
- Train a Student Model for Perception

[1] Jia et al. DriveAdapter: New Paradigm for End-to-End Autonomous Driving to Alleviate Causal Confusion. ICCV, 2023.

Open 🔁 rive Lab
Lack of Interpretability

Summary of the different forms of interpretability



They aid in human comprehension on the:

- Decision-making processes of end-to-end models
- Perception failures
- Reliability of the outputs



Lack of Safety Guarantees

Safety of the intended functionality (SOTIF)

• Applicable design, verification, and validation measures

Modular driving stacks

- Safety-related constraints or optimizations, within motion planning or speed prediction modules
- Integrated into E2E models as post-process steps or safety checks
- Detection and motion prediction results can be used in post-processing procedures



Causal Confusion



- Driving is a task that exhibits **temporal smoothness**, which makes past motion a reliable predictor of the next action.
- However, methods trained with **multiple frames** can become overly reliant on this shortcut. This is referred to as the **copycat problem** and is a manifestation of **causal confusion**.



Causal Confusion

Current Solutions

- Adversarial model predicts ego's past action → min-max optimization trains the model to eliminate its past from intermediate layers
- Random dropout
- Upweighting keyframes in the training loss, where a decision change occurs
- Predict action residuals instead of actions
- Use stacked LiDAR points

Still challenging



Fighting copycat agents in behavior cloning from observation histories, NeurIPS' 20

- Problem: F learns to rely heavily on the information about a_{t-1} in e_t to predict a_t
- Target: Remove information about a_{t-1} in e_t
- Method:
 - Adversarial (to remove information) network
 D predicts a₊₁ from e₊
 - E maximizes the conditional entropy $H(a_{t-1}|e_t)$
- Issue: removing all maybe counterproductive; the copycat problem arises only when a_t and a_{t-1} are highly correlated

Open AriveLab

Lack of Robustness





Lack of Robustness - Long-tail Distribution

Current Solutions

- Hand-crafted scenarios for more diverse data in simulation, especially for OOD perception
- Non-ego agents' prediction to promote data diversity (philosophy of integrated prediction and planning)
- Importance-sampling to accelerate evaluation of rare-event probabilities
- Adversarial attacks
 - Bayesian Optimization
 - Policy gradient for generation
 - Hand-crafted modification on agents' trajectories



Open 🔁 rive Lab





未来可能与讨论 / Future Trends and Discussions

Foundation Models



Open AriveLab

Foundation Models (cont'd)



Open PriveLab

Trending in E2EAD | Driving + Language





Trending in E2EAD | Driving + Language



Open PriveLab

DriveLM | Introduction

Generalization and Interactivity in Autonomous Driving

- Generalized to unseen sensor configuration and objects
- Regional / Global (e.g. European) regulations require explainability through interaction

Recent success in Vision Language Models

- Good **reasoning** ability, enabled by LLM
- No BEV representation, since human do not rely on BEV

Why VLM in AD?

- Reasoning ability helps generalization
- Language output provides interactivity



















DriveLM | Data



- To ensure the data quality, we introduce human annotation with multi-round quality check in nuScenes
- To scale up annotation, we adopt auto-labelling in CARLA



DriveLM | Data







- To ensure the **data quality**, we introduce human annotation with multi-round quality check in nuScenes
- To scale up annotation, we adopt auto-labelling in CARLA

Diversity matters, spanning from perception to prediction and planning

Open 🔁 rive Lab

DriveLM | Agent



- 🧩 General and scalable VLM architecture
- 🌏 Web-scale pre-training

- 🋠 Fine-tuned end-to-end for planning
 - Interpretable and interactive

_



DriveLM - Experiments

Method	Behavior Context	Motion Context	Behavior (B)			Motion (M)	
			Acc. \uparrow	Speed \uparrow	Steer \uparrow	ADE↓	$FDE \downarrow$
Command Mean	-	-	-	-	-	7.98	11.41
UniAD-Single BLIP-RT-2	-	-	-	-	-	4.16	9.31
	-	-	-	-	-	2.78	6.47
DriveLM-Agent	None	В	35.70	43.90	65.20	2.76	6.59
	Chain	B	34.62	41.28	64.55	2.85	6.89
	Graph	В	39.73	54.29	70.35	2.63	6.17

Conclusion:

• Trained on DriveLM-Data (nuScenes-based), DriveLM-Agent (ours) gains **better zero-shot ability** on Waymo scenarios, overpassing other methods by a large margin.



DriveLM - Experiments

Method	Behavior	Motion	Behavior (B)			Motion (M)	
	Context	Context	Acc. \uparrow	Speed \uparrow	Steer \uparrow	ADE \downarrow	$FDE \downarrow$
Command Mean	-	-	-	-	-	7.98	11.41
UniAD-Single BLIP-RT-2	-	-	-	-	-	4.16	9.31
	-	-	(=)	-	3 -	2.78	6.47
DriveLM-Agent	None	В	35.70	43.90	65.20	2.76	6.59
	Chain	B	34.62	41.28	64.55	2.85	6.89
	Graph	В	39.73	54.29	70.35	2.63	6.17

Conclusion:

 Trained on DriveLM-Data (nuScenes-based), DriveLM-Agent (ours) gains better zero-shot ability on Waymo scenarios, overpassing other methods by a large margin.

Conclusion:

• Qualitative result shows that DriveLM-Agent does **understand the unseen scenarios** in some way.



Open AriveLab

DriveLM - Limitation



Driving-specific Inputs

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



Closed-loop Planning

DriveLM-Agent 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, DriveLM-Agent suffers from long inference time, which may impact practical implementation.



One-page Takeaway

End-to-end Autonomous Driving

- Challenge: Generalization & Explainability
- Recent trend: use vision language model to embed "world knowledge" to solve challenges

DriveLM: Driving with Graph Visual Question Answering

- Use Graph VQA as a proxy task to mimic human's driving logic
- Some good result under zero-shot setting, but still far from claiming good generalization



System 1&2



System 1: fast, automatic, frequent, emotional, stereotypic, unconscious

- Determine that an object is at a greater distance than another
- Drive a car on an empty road

System 2: Slow, effortful, infrequent, logical, calculating, conscious

- Look for the woman with the grey hair
- Park into a tight parking space

[1] Daniel Kahneman. Thinking, Fast and Slow. 2011.



System 1&2

DriveVLM



Idea

- Incorporate results from the traditional pipeline to LVLM, as prompts
- Chain-of-thought VLM
- Take the low frequency VLM planning output as a reference to refine traditional trajectory planning (selectively attend to the additional information)

[1] Tian et al. DriveVLM: The Convergence of Autonomous Driving and Large Vision-Language Models. arXiv, 2024.







Q&A

End



 \bigcirc

2

0

Shanghai Al Lab