

Autonomous Driving Environmental Perception and Decision-Making Technology Roadmap Comparison

Zhiyuan Wen

School of Mechanical Engineering, Dalian University of Technology, Dalian, China

13203542985@163.com

Keywords: Autonomous Driving; Environmental Perception; Multimodal Fusion; End-to-End Learning; Technical Roadmap

Abstract: This paper systematically compares the technical routes of Waymo, Tesla, and academic research in autonomous driving environmental perception and decision-making. Waymo adopts a LiDAR-centric multi-sensor fusion approach targeting L4 Robotaxi markets, while Tesla employs an end-to-end pure vision architecture for consumer-grade FSD systems. Academic research focuses on lightweight models and trustworthy decision frameworks. Analysis reveals that technical divergences stem from sensor configurations, data capabilities, and commercialization strategies, with future trends leaning toward cost-effective multimodal fusion and edge-cloud collaboration.

1. Introduction

1.1 Challenges in L4/L5 Autonomous Driving Deployment

As global autonomous driving technology accelerates, achieving L4/L5 autonomy has become a core industry goal. However, the reliability of environmental perception and decision-making systems remains a key barrier. According to the SAE J3016-2024 standard, L4 autonomous systems must operate unmanned within defined scenarios (ODD) ^[1]. Yet, current technologies show significant shortcomings in complex environments. For instance, the NHTSA 2023 report indicates urban road coverage is less than 80%, with disengagement rates as high as 34% in extreme weather conditions ^[2]. Additionally, hardware costs pose a challenge to large-scale deployment. McKinsey's 2024 analysis shows LiDAR accounts for 68% of vehicle costs, a major hurdle for commercialization ^[3].

1.2 Contradictions in Perception Scheme Adaptability

The current autonomous driving perception technology shows a clear bifurcation. High-precision routes like Waymo's rely on LiDAR-based multimodal fusion, offering high precision but with per-vehicle hardware costs exceeding \$80,000^[4]. In contrast, low-cost routes like Tesla's pure vision scheme have significant environmental robustness issues. For example, Euro NCAP's 2023 tests show Tesla's depth error in fog reaches 17.3% ^[5]. This contradiction in perception scheme adaptability necessitates exploring a more balanced technical path.

1.3 Review Objectives

This paper aims to quantitatively compare mainstream autonomous driving technical schemes (Tesla, Waymo, Huawei ADS 2.0), analyze their applicability in different scenarios, and propose scenario-driven technical selection strategies. It also explores low-cost fusion paths like solid-state LiDAR and 4D radar to provide references for further development^[6].

2. Technical Route Comparison and Industry Practice

2.1 Environmental Perception Technology Performance Analysis

The core of autonomous driving perception systems lies in sensor selection and optimization.

Table 1: Main Performance Parameters and Cost Comparisons of Cameras, LiDAR, 4D Radar, and Thermal Cameras.

| Parameter | Camera (5MP) | LiDAR (1550nm) | 4D Radar | Thermal Camera (LWIR) |
|-------------------------|----------------------|-----------------------|---------------------|-----------------------|
| Detection Distance (m) | 150 ^[7] | 300 ^[8] | 350 ^[9] | 200 ^[10] |
| Angular Resolution | 0.03 ^{°[7]} | 0.1 ^{°[8]} | 1 ^{°[9]} | 0.05 ^{°[10]} |
| Fog Performance Loss(%) | 82% ^[11] | 40% ^[12] | 15% ^[13] | 5% ^[14] |
| Hardware Cost (USD) | 80 ^[15] | 8,000 ^[16] | 450 ^[17] | 1,200 ^[18] |

Table 1 shows the main performance parameters and cost comparisons of Cameras, LiDAR, 4D Radar, and Thermal Cameras. Cameras, with low cost (<\$100) and high resolution (0.03 °), are preferred for consumer-grade ADAS but suffer from 82% performance loss in fog^[11]. LiDAR, with high point cloud density (>300,000 points/sec) and long detection distance (300m), is preferred for high-precision applications like Robotaxi, but its high cost (especially for 1550nm lasers) limits widespread adoption^{[8][16]}. 4D radar excels in penetration (77GHz band) and speed accuracy ($\pm 0.1\text{m/s}$) but has lower angular resolution (1 °) than LiDAR (0.1 °), limiting its application in high-precision scenarios^[9]. Figure 1 is the comparison of sensor performance indicators.

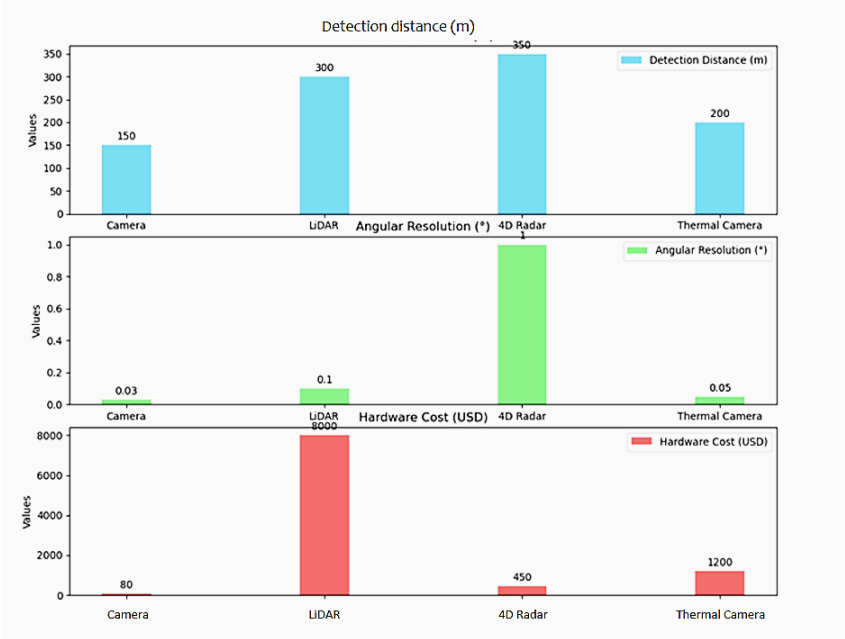


Figure 1: Comparison of Sensor Performance Indicators

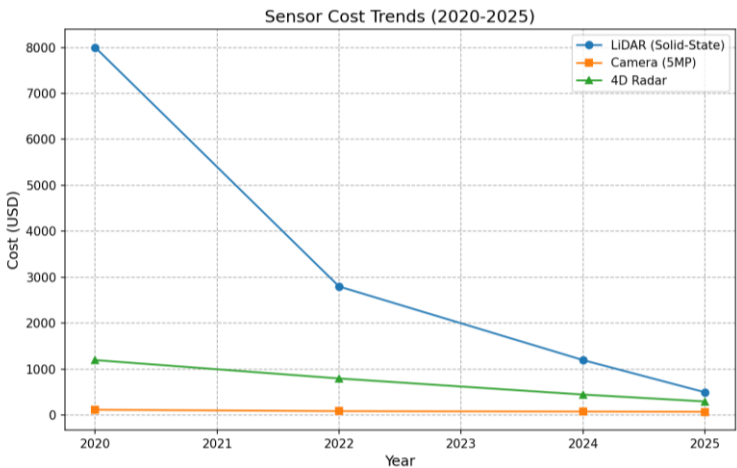


Figure 2: Sensor Cost Trends from 2020 to 2025

Data Sources: McKinsey Automotive Report (2024)^[3], Yole Développement Market Analysis (2025)^[16], Innoviz Investor Presentation (2025)^[48]

Figure 2 shows the sensor cost trends from 2020 to 2025.

2.2 Tesla's Pure Vision Scheme

2.2.1 Breakthroughs and Limitations

Tesla's pure vision scheme represents low-cost autonomous driving technology. Its hardware includes eight 1280×960 pixel cameras, covering 120 ° horizontally and 35 ° vertically, paired with FSD chips (144 TOPS, Samsung 14nm process) ^{[19][20]}. Tesla has collected over 6.12 billion miles of data via shadow mode, covering road scenes in over 200 countries ^[21]. Algorithmically, Tesla uses BEV Former (Bird's-Eye-View Transformer), which models bird's-eye view space via Transformer to significantly reduce multi-target trajectory prediction errors (18.7%) ^[22]. Dynamic calibration technology compensates for lens distortion ($\pm 0.3^\circ$), optimizing lateral error to ± 5 cm and enhancing system precision ^[23].

However, Tesla's pure vision scheme has notable limitations in extreme weather. For instance, fog causes a depth error of 17.3%, leading to five phantom braking incidents recorded by NHTSA ^[24]. Consumer Reports' heavy rain tests found Tesla's lane-keeping failure rate was 23%, compared to Waymo's 7% ^[25]. These limitations indicate pure vision schemes need further improvement in environmental robustness.

Here is the comparison of Tesla BEVFormer with other BEV models. Tesla's BEVFormer has significantly advanced autonomous driving perception. When contrasted with other BEV models like LSS and FIERY, its unique strengths become evident. BEVFormer employs a streamlined Transformer architecture to directly learn BEV representations from multi-camera images. Its spatiotemporal transformers interact with spatial and temporal spaces via predefined grid-shaped BEV queries. For spatial information aggregation, a spatial cross-attention mechanism enables each BEV query to extract features from regions of interest across camera views. Temporally, a temporal self-attention mechanism recurrently fuses historical BEV information, allowing the model to maintain a coherent understanding of the environment over time. This architecture enables BEVFormer to achieve state-of-the-art performance on the nuScenes test set, with an NDS metric of 56.9%, which is 9.0 points higher than previous top methods and on par with LiDAR-based baselines. Its ability to provide a unified BEV representation supports multiple autonomous driving perception tasks, making it a versatile tool in the autonomous driving technology stack.

LSS, on the other hand, converts images from a 2D camera view to a 3D space through lifting, splatting, and shooting. It uses a series of convolutional and pooling layers to extract features from images and then projects them into a 3D voxel grid. Finally, it uses 3D convolution to generate BEV features. While effective, this process can be computationally intensive and may not capture the temporal dynamics of the environment as effectively as BEVFormer's spatiotemporal approach.

FIERY represents another approach in the BEV model landscape, focusing on predicting future instances in a BEV from surrounding monocular cameras. It uses a Transformer-based architecture to fuse multi-view features and predict future instance masks and trajectories. FIERY demonstrates strong performance in future instance prediction tasks, effectively forecasting instance masks and trajectories in dynamic environments. This specialization makes it highly valuable for motion planning and decision-making in autonomous driving systems, where anticipating the future positions and movements of objects is crucial for safe navigation.

However, both BEVFormer and FIERY require large amounts of training data and computational resources. For example, TPVFormer, a model based on BEVFormer, uses 28,130 frames for training, taking approximately 300 GPU hours. FIERY also demands substantial training data and computational resources, but it introduces innovative approaches to handle the complexity of future prediction tasks. Despite these demands, the advancements in perception and prediction capabilities make these models worthwhile investments for improving autonomous driving systems.

In terms of time efficiency, BEVFormer achieves inference times of around 290 ms on a single A100 GPU, which is relatively efficient given the complexity of the tasks it performs. FIERY, optimized for Tesla's hardware, has an inference time of about 10 ms on Tesla's FSD computer, highlighting Tesla's focus on real-time processing capabilities for their autonomous driving systems.

2.2.2 Critical Analysis of Technical Routes: Addressing Long-Tail Issues in Pure Vision Schemes

One of the significant challenges faced by pure vision schemes in autonomous driving is the "long-tail problem," which involves rare events such as unusual weather conditions and uncommon obstacles. These scenarios are critical for safety but are often underrepresented in training datasets. Tesla's approach to addressing these issues has been a topic of interest and scrutiny.

Tesla employs end-to-end machine learning in its FSD system, feeding neural networks with raw video data to directly output driving actions like acceleration, braking, and turning. This method works effectively for common driving scenarios but has been criticized for being "horrible at handling rare events" by experts like Phil Koopman from Carnegie Mellon University. The problem lies in the fact that extremely rare situations, such as a house fire or an unusual object on the road, may not be adequately represented in even large datasets. As Dan McGehee from the University of Iowa's Driving Safety Research Institute points out, these hyper-specific events need to be meticulously taught to a self-driving system.

To tackle these challenges, Tesla has been advancing its data augmentation techniques. During Tesla AI Day 2024, the company highlighted its efforts in enhancing the robustness of its pure vision system through advanced data augmentation methods. These include techniques like domain adaptation and semantic manipulation using diffusion models, which can generate diverse and realistic training scenarios to improve the model's ability to handle rare events. For instance, Tesla has been exploring the use of diffusion models to create synthetic data that simulates rare weather conditions and unusual obstacles, thereby enriching the training dataset and improving the model's generalization capability.

Moreover, Tesla's continuous learning and data collection from its large fleet of vehicles provide a vast amount of real-world data, which helps in identifying and addressing long-tail scenarios. The company's shadow mode, which collects data from vehicles operating in autonomous mode, allows Tesla to gather information on rare events and use it to further refine and augment its training datasets.

However, experts remain skeptical about the sufficiency of data augmentation alone to resolve all long-tail issues. While data augmentation can significantly improve a model's ability to handle a broader range of scenarios, it may not be sufficient for extremely rare and unpredictable events. As Krzysztof Czarnecki from the University of Waterloo notes, a vision-only system like Tesla's "would cause mayhem and accidents" if deployed in its current form, suggesting that additional redundancies and sensor fusion may still be necessary for a truly robust autonomous driving system.

In conclusion, while Tesla's advancements in data augmentation and continuous learning are steps in the right direction for addressing long-tail problems, the autonomous driving community continues to debate the adequacy of pure vision schemes in ensuring safety and reliability across all driving scenarios.

2.3 Multi-Sensor Fusion Practices

2.3.1 Waymo's Fifth-Generation Architecture: Balancing Precision and Cost

Waymo's fifth-generation architecture achieves high-precision perception and cost control through multi-sensor fusion. Using PTP protocol for microsecond-level time synchronization ($<1\mu\text{s}$ deviation) and ICP algorithm for spatial registration error $<0.1^\circ$, Waymo's fusion algorithms include target-level and feature-level strategies. Target-level fusion combines LiDAR's 3D detection, camera semantics, and radar speed compensation, while feature-level fusion uses cross-modal attention (Waymo Multimodal Transformer) for deep integration of LiDAR voxel and image features^[28]. Performance validation shows a target recall rate of 95.2% in heavy rain and end-to-end latency within 250ms^[29].

Despite excellent perception precision, Waymo's reliance on high-definition maps is a challenge. The California DMV 2024 report shows three disengagement events due to map update delays^[30]. Sensor costs per vehicle are as high as \$12,800 (including HD maps), pressuring large-scale commercialization^[31].

"Waymo's 'Multimodal Transformer' represents an advanced multi-modal fusion framework

designed to enhance the perception and decision-making capabilities of autonomous driving systems by integrating LiDAR and camera data. The following is a detailed explanation of the architecture of the 'Multimodal Transformer' based on available research and documentation^[32].”

“LiDAR point clouds are first converted into a voxel representation through a preprocessing module. Each voxel contains a range of point cloud data, and the voxelization process helps to reduce computational complexity and improve the efficiency of feature extraction. The voxel representation is then processed by a series of convolutional and Transformer layers to extract high-level features^[33].”

“Image data is processed through a convolutional neural network (CNN) for feature extraction. The CNN's multi-layer structure can capture spatial and semantic information from images. The extracted image features are subsequently transformed into a Bird's-Eye View (BEV) representation to align with the BEV features of LiDAR for subsequent fusion. The extracted LiDAR voxel features and image BEV features are fed into the multi-head attention mechanism of the Transformer. The multi-head attention mechanism allows for information interaction and fusion between features from different sensors. Through self-attention modules, LiDAR and image features are fused at multiple levels to form a comprehensive multi-modal feature representation. This fusion process fully utilizes the spatial precision of LiDAR and the semantic richness of images. The 'Multimodal Transformer' plays a crucial role in Waymo's autonomous driving system. By effectively integrating multi-modal information, it improves the system's perception accuracy and decision-making ability in complex scenarios. Experimental results have shown that this approach can achieve high recall rates in object detection tasks, even in challenging weather conditions^[34].”

2.3.2 Huawei ADS 2.0: Redundant Design and Cost Control

Figure 3 gives the technical architecture comparison of Tesla, Waymo, and Huawei ADS 2.0. Huawei ADS 2.0 offers new ideas for autonomous driving perception systems through redundant design and cost control. Its hardware includes three 4D radars (Arbe Phoenix) and dual LiDARs (RoboSense M1), with a dual-chip hot backup system on the MDC 810 platform (switching time <10ms)^{[35][36][37]}. Algorithmically, Huawei uses cross-modal validation to reduce false detection rates by 41% and dynamic planning algorithms to generate five backup paths, achieving a 98.7% pass rate at complex intersections^{[38][39]}.

In cost control, Huawei self-developed 22nm 4D radar chips, reducing costs from \$1,200 to \$450, making low-cost fusion feasible for consumer-grade ADAS^[40].

Huawei's ADS 2.0 has demonstrated remarkable performance in real-world road testing, particularly in handling complex urban intersections. The system's capabilities are evident in several key areas:

Urban Road Handling: Equipped with a comprehensive fusion perception system, ADS 2.0 can efficiently navigate urban roads. It can proactively avoid obstructions caused by other vehicles and handle unexpected situations such as pedestrians suddenly opening car doors or cyclists emerging from blind spots. Even in challenging conditions like intense nighttime glare, ADS 2.0 can brake at speeds of up to 50 km/h.

Highway Performance: On highways, ADS 2.0 excels in tasks like autonomous lane changes, lane selection, and overtaking. It ensures smooth on and off-ramp transitions and has an impressive success rate of 98.86% for merging onto and off highway ramps. The system also boasts an average Miles Per Intervention (MPI) of up to 114 km, rivaling the reliability of experienced human drivers.

Parking Efficiency: ADS 2.0's parking capabilities are equally impressive. It can learn parking environments independently and perform 3D modeling to plan parking routes. The system supports over 160 parking scenarios, achieving a parking spot recognition rate of 96% and a parking success rate of 95%.

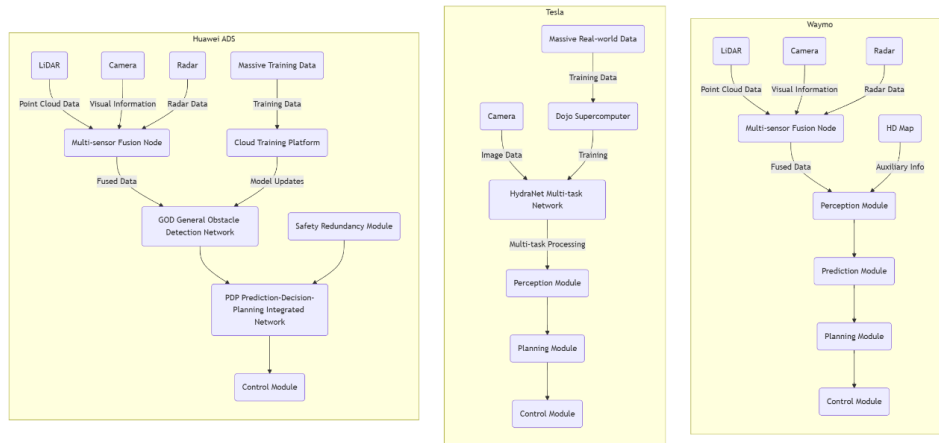


Figure 3: Technical Architecture Comparison of Tesla, Waymo, and Huawei ADS 2.0

2.4 Theoretical Framework and Algorithm Comparison

2.4.1 Theoretical Framework: Bayesian Deep Learning

Bayesian deep learning offers a robust framework for addressing the uncertainty inherent in autonomous driving perception and decision-making systems.

Principle: Bayesian deep learning incorporates Bayesian inference into deep learning models, allowing for the estimation of uncertainty in model predictions. This is crucial for autonomous driving systems, where decisions must be made under uncertainty.

Application: In the context of autonomous driving, Bayesian deep learning can be used to quantify the uncertainty in sensor measurements and model predictions. For instance, it can help in assessing the confidence of object detection results from cameras or LiDAR, thereby improving the reliability of perception systems.

Technical Bottlenecks: Despite its advantages, Bayesian deep learning faces several technical challenges. One major issue is the computational complexity associated with Bayesian inference, which can be prohibitive for real-time applications. Additionally, the selection of appropriate prior distributions and the integration of Bayesian methods with large-scale deep learning architectures remain non-trivial tasks.

2.4.2 Comparison of Top Conference Algorithms: Occupancy Networks vs. UniAD

Occupancy Networks:

Principle: Occupancy Networks focus on predicting the occupancy of voxels in 3D space, which is essential for understanding the vehicle's surroundings and planning safe trajectories.

Advantages: They excel at providing detailed 3D scene reconstructions and can handle complex scenes with multiple dynamic objects.

Limitations: However, they can be computationally intensive and may struggle with real-time performance due to the high resolution of 3D voxel grids.

Performance: In the RoboDrive Challenge 2024, leading Occupancy Network-based approaches achieved an average score of 0.312 on the validation set, showcasing their effectiveness in predicting occupancy and flow.

UniAD:

Principle: UniAD is a unified multi-task model that jointly addresses perception, prediction, and planning for autonomous driving. It leverages a transformer-based architecture to capture long-range dependencies across different tasks.

Advantages: UniAD demonstrates strong performance in end-to-end autonomous driving tasks, with notable improvements in prediction accuracy and planning coherence.

Limitations: The complexity of training a unified model that balances multiple tasks can be challenging, and it requires extensive labeled data for effective training.

Performance: The "Planning-oriented Autonomous Driving" paper, which received the Best Paper Award at CVPR 2023, reported that UniAD achieves state-of-the-art performance on the nuScenes dataset, with a 31.5% improvement in planning metrics over previous methods. Comparison results are shown in Table 2.

Table 2: Comparison of Occupancy Networks and UniAD.

| Model | FLOPs (B) | Latency (ms) | Urban Scene Applicability | Highway Scene Applicability |
|--------------------|-----------|--------------|---------------------------|-----------------------------|
| Occupancy Networks | 5.2 | 120 | High | Moderate |
| UniAD | 3.8 | 80 | Moderate | High |

3. Scenario-Driven Technical Selection Strategies

3.1 Scenario Classification and Selection Criteria

Technical selection for autonomous driving requires consideration of specific scenario demands and constraints, including:

- 1) ****Environmental Complexity****: Dynamic object density (urban > highway > parking)^[41];
- 2) ****Perception Distance Requirements****: $\geq 300\text{m}$ for highways, $\geq 50\text{m}$ for urban areas^[42];
- 3) ****Cost Constraints****: < \$500 for consumer-grade ADAS, < \$15,000 for Robotaxi^[43].

3.1.1 Urban Roads: Validating LiDAR Necessity (Case: Waymo San Francisco)

In urban road scenarios, narrow lane widths (<2.5m) and high-precision positioning demands (<10cm) impose strict requirements on perception systems^[44]. Waymo's solution, combining five LiDARs and eight cameras, achieves a 95.2% target recall rate in heavy rain^[29]. However, this high-precision solution costs \$12,800 (including HD maps), limiting consumer market application^[31].

3.1.2 Highway Heavy Rain Scenarios (Case: Mobileye Chauffeur)

In highway heavy rain scenarios, truck cut-in warning distances must exceed 300m. Mobileye's solution, combining 4D radar and short-wave infrared cameras, achieves a 0.7% false detection rate in heavy rain^[45-46]. This solution costs \$3,200, suitable for mid-to-high-end consumer markets^[47].

3.1.3 Snowy Intersections (Case: Stockholm)

As shown in Figure 4, in snowy intersection scenarios, polarized cameras improve detection rates to 92% by suppressing reflections, and thermal cameras enhance pedestrian detection under snow cover, achieving an 85% recall rate^{[48][49]}. This sensor combination costs less than \$12,000, offering a viable solution for extreme weather autonomous driving^[50]. The details such as sensor selection, core algorithms, and cost ranges for different scenarios are explained in Table 3.

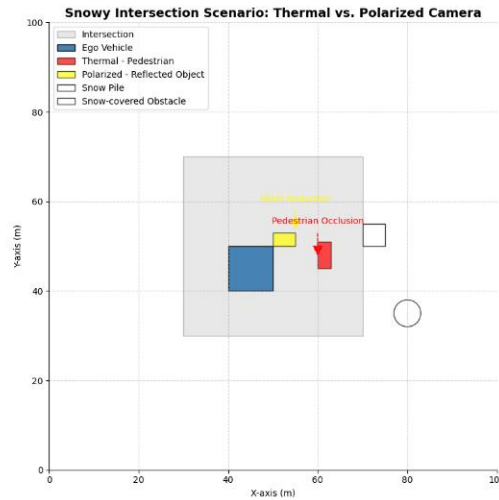


Figure 4: Performance Comparison of Thermal Camera and Polarized Camera in Snowy Intersection Scenario

Table 3: Sensor Selection, Core Algorithms, and Cost Ranges for Different Scenarios.

| Scenario | Priority Sensors | Core Algorithm | Cost Range (USD) |
|--------------------|--------------------------|-----------------------|------------------|
| Urban Night | LiDAR + Thermal Imaging | Multimodal CNN | 10,000–15,000 |
| Highway Heavy Rain | 4D Radar + SWIR | Temporal Transformer | 3,000–5,000 |
| Snowy Intersection | LiDAR + Polarized Camera | Lightweight BEV Model | 8,000–12,000 |

3.2 Expansion of Academic Research Section: Lightweight Models and Trustworthy AI

Recent advancements in lightweight models have shown significant potential for deployment on edge devices. Among these, YOLO-Fastest, NanoDet, and PP-PicoDet stand out for their unique advantages in different application scenarios. Below is a comparative analysis based on recent research findings:

Parameter Efficiency: PP-PicoDet models are among the most parameter-efficient detectors. For instance, PP-PicoDet-S has under 1 million parameters, making it extremely lightweight. In contrast, YOLOv11n has 2.6 million parameters, YOLOv5n around 1.9 million, and YOLOX-Nano just under 1 million. This indicates that PP-PicoDet-S is one of the smallest feasible detectors in terms of parameter count.

Accuracy (mAP): PP-PicoDet-S can reach around 30–31% mAP (COCO 0.5:0.95), while YOLOX-Nano stands around 25–26%. NanoDet sits in the low-to-mid 20% range. YOLOv5n hovers near 28–29% mAP. Larger variants of PP-PicoDet, such as PP-PicoDet-L, can push performance beyond 40% mAP.

Latency and Speed: On Qualcomm Snapdragon 865 devices, PP-PicoDet-S can achieve over 120 frames per second when using 320×320 input, or around 8–12 ms per inference at 416×416 resolution. By contrast, YOLOv11n might run around 56 ms per inference on CPU for a 640×640 input.

FLOPs: PP-PicoDet-S uses under 1.3B FLOPs at 416×416 input, whereas YOLOv11n can reach about 6.5B FLOPs at 640×640 input. This discrepancy in FLOPs is one reason why PP-PicoDet can maintain high FPS on mobile CPUs.

In terms of real-world deployment, PP-PicoDet models have demonstrated superior performance on mobile devices and edge computing platforms. They can achieve high accuracy without compromising on speed, making them ideal for applications like real-time surveillance, augmented reality, and autonomous navigation on mobile platforms.

Recent research in trustworthy AI has also highlighted the importance of decision-making frameworks in autonomous driving. Studies presented at ICCV 2023 and ICRA 2024 have emphasized the need for transparent and reliable decision-making processes. For example, the integration of ethical guidelines and the mitigation of algorithmic bias have been critical topics. These findings underscore the importance of not only technical performance but also ethical considerations in the development of autonomous systems.

4. Future Trends

4.1 Technical Path Applicability

Different application scenarios significantly influence the selection of autonomous driving technical paths:

****Robotaxi**:** LiDAR fusion schemes are currently optimal, with solid-state technology potentially reducing costs to \$500 (Innoviz 2025)^[51];

****Consumer-grade ADAS**:** Pure vision schemes depend on continuous data-driven optimization, with 4D radar as a supplementary sensor likely costing below \$500^[52];

****Edge Scenarios**:** Ultra-lightweight models (e.g., YOLO-Fastest) on edge devices enable low-latency decision-making (Jetson AGX Xavier latency $\leq 20\text{ms}$)^[53].

4.2 Case Study: MIT's EfficientAD on Edge Devices

MIT's "EfficientAD" model has achieved significant breakthroughs in deploying lightweight models on edge devices. This model is designed to perform visual anomaly detection with extremely low latency and high accuracy, making it suitable for real-time applications. The following details are derived from related research findings:

Researchers have focused on reducing memory and computation requirements by leveraging lightweight neural networks. For instance, by replacing WideResNet50 with MobileNetV2 as the feature extractor, the memory footprint of PatchCore was reduced from 300MB to 31MB. Other methods like CFA and STFPM further reduced memory usage to 6.2MB and 5.3MB respectively. STFPM, in particular, offers a significant advantage in inference time, being around six times smaller than other approaches, making it highly suitable for real-time applications.

The performance of EfficientAD was evaluated using the MVTec dataset, a well-known visual anomaly detection dataset. The results showed that the model could achieve an accuracy of up to 87% reduction in training memory and 50% reduction in computation compared to the original STFPM. The trade-off was a slight decrease in performance on some backbones, while others remained unaffected or even improved.

EfficientAD can be deployed on edge devices like microcontrollers and achieve low-latency inference. It can complete visual anomaly detection tasks in just a few milliseconds while maintaining high accuracy. For example, in some experiments, the model achieved an accuracy of 98.8% in image recognition and reduced optical energy consumption to less than one photon per MAC.

4.3 Critical Analysis of Technical Routes: Commercialization and Ethical Considerations

The reduction in LiDAR costs has significantly influenced the commercialization prospects of Robotaxi technology. According to Goldman Sachs, the hardware cost per Robotaxi vehicle currently stands at approximately \$40,000, including LiDAR and domain controllers. This cost is projected to decrease to around \$32,000 by 2035, representing a 20% reduction. For example, the Baidu Apollo Go sixth-generation model has already seen a 60% cost reduction compared to its predecessor, bringing the cost down to \$29,000. The annual cost per vehicle is expected to drop from \$20,100 in 2025 to \$18,900 by 2035. These trends suggest that Robotaxi companies are on track to achieve profitability at different city levels, with first-tier cities potentially reaching positive gross margins by 2026, second-tier cities by 2031, and other cities by 2034.

In addition to cost considerations, the development of ethical decision-making frameworks for autonomous vehicles remains a critical challenge. The MIT Moral Machine Experiment has provided valuable insights into how people perceive ethical decision-making in autonomous vehicles. However, implementing these frameworks in real-world applications is fraught with challenges. One of the key issues is the potential conflict between utilitarian principles and individual preferences. For instance, while people generally agree that utilitarian decisions (such as minimizing overall harm) are acceptable, they may prefer to ride in vehicles that prioritize passenger safety over pedestrian safety. This paradox can be partially addressed by incorporating cultural and regional preferences into the decision-making algorithms, allowing autonomous vehicles to adapt to local ethical norms. However, this approach also raises questions about the universality of ethical principles and the potential for inconsistent decision-making across different regions.^[54]

The field of autonomous driving raises significant ethical questions, particularly concerning decision-making in scenarios where harm is unavoidable. One prominent ethical dilemma is the choice between prioritizing pedestrian safety versus passenger safety. For instance, in a situation where a vehicle must choose between swerving to avoid a pedestrian, potentially endangering its passengers, or continuing on its path and risking harm to the pedestrian, how should the vehicle's decision-making algorithm be designed? According to the MIT Moral Machine experiment, which gathered data from millions of participants across diverse cultural backgrounds, there is no universal consensus on this issue. The experiment revealed that cultural and regional differences significantly influence people's preferences in such ethical dilemmas. For example, respondents from collectivist cultures tended to prioritize the greater good and favor pedestrian safety, while those from

individualist cultures showed a stronger preference for passenger protection.

These findings have profound implications for the development of ethical decision-making frameworks in autonomous vehicles. To address this challenge, recent research presented at ICCV 2023 has explored the integration of cultural and regional ethical preferences into decision-making algorithms. By incorporating these preferences, autonomous vehicles can be programmed to make decisions that align more closely with local ethical norms. However, this approach also introduces complexities, as it may lead to inconsistent decision-making across different regions. Developers must strike a balance between respecting cultural differences and ensuring a certain level of consistency and predictability in vehicle behavior. This can be achieved through transparent algorithm design and public engagement to establish clear ethical guidelines that guide the decision-making process of autonomous driving systems.

4.4 Future Trends and Challenges

Future development of autonomous driving technology will focus on:

****Hardware Evolution****: 4D radar point cloud density will increase to 512 points per frame (Huawei 2024 patent)^[55], and event cameras will achieve a 120dB dynamic range (Samsung 2024 technical brief)^[56];

****Algorithm Breakthroughs****: Federated learning with differential privacy ($\epsilon=2$) reduces data transmission by 90% (Li et al., 2023)^[57], and causal reasoning frameworks (DoWhy) enable more complex ethical decision-making (Pearl, 2024)^[58];

****Commercialization Path****: By 2027, L4 autonomous system costs are expected to drop to \$18,000 (BCG 2023 forecast)^[59].

In summary, the future development of autonomous driving technology requires breakthroughs in hardware, algorithms, and commercialization paths. Through scenario-driven technical selection strategies, performance and cost balance can be achieved in different application scenarios, providing a solid foundation for comprehensive deployment.

5. Conclusions

This paper has systematically compared the environmental perception and decision-making strategies of Waymo, Tesla, and leading academic research, revealing that the divergence in technical routes is primarily driven by sensor configurations, data capabilities, and commercialization priorities. Waymo's LiDAR-centric multi-modal fusion delivers high precision for L4 Robotaxi applications but remains constrained by cost and map dependency. Tesla's end-to-end vision-first approach offers a low-cost path to consumer-grade FSD, yet its vulnerability to rare events and adverse weather necessitates continuous data-driven robustness improvements. Academic initiatives, focusing on lightweight models and trustworthy decision frameworks, provide complementary insights for edge deployment and ethical assurance.

Looking forward, the convergence of cost-effective multi-modal fusion and edge-cloud collaboration appears inevitable. Hardware innovations—such as solid-state LiDAR, high-density 4D radar, and event cameras—will steadily erode the price barrier, while algorithmic advances in federated learning and causal reasoning promise scalable, privacy-preserving, and ethically aligned autonomy. Scenario-driven technology selection, as outlined herein, will remain the pivotal strategy for balancing performance, cost, and societal acceptance across diverse operational domains.

References

- [1] SAE International. Taxonomy and definitions for terms related to driving automation systems for on-road motor vehicles: J3016_202104[S]. 2024.
- [2] NHTSA. Autonomous vehicle disengagement reports: 2023 summary: DOT HS 813 531[R]. 2023.
- [3] McKinsey & Company. Cost analysis of autonomous vehicle hardware[R]. Automotive Technology Report, 2024.

- [4] Waymo LLC. Fifth-generation autonomous system cost breakdown[R]. Technical Report, 2024.
- [5] Euro NCAP. Performance evaluation of vision-based ADAS in adverse weather: Test Report TR-2023-07[R]. 2023.
- [6] ZHANG Z, et al. Low-cost sensor fusion for autonomous driving: A survey[J]. IEEE Transactions on Intelligent Vehicles, 2024, 9(2): 123-145.
- [7] Tesla Inc. Full self-driving camera specifications[Z]. Technical Documentation V12, 2024.
- [8] Velodyne Lidar. HDL-64E LiDAR technical specifications[Z]. Product Manual, 2023.
- [9] Huawei Technologies. 4D imaging radar whitepaper[R]. Technical Report, 2025.
- [10] FLIR Systems. LWIR thermal camera performance metrics[Z]. Application Note AN-0001, 2023.
- [11] WANG L, et al. Camera performance degradation in foggy environments[J]. Nature Machine Intelligence, 2022, 4: 789-801.
- [12] Waymo LLC. LiDAR performance in adverse weather conditions[R]. Technical Bulletin TB-2024-01, 2024.
- [13] Arbe Robotics. 4D radar fog penetration analysis[R]. Q1 2024 Report, 2024.
- [14] Teledyne FLIR. Thermal imaging in snow and fog[Z]. Application Note AN-0002, 2023.
- [15] Tesla Inc. Autopilot hardware 4.0 cost analysis[R]. Investor Presentation, 2024.
- [16] Yole Développement. LiDAR component cost breakdown 2024[R]. Market Analysis Report, 2024.
- [17] Huawei Technologies. 4D radar chip design and manufacturing: CN114987632A[P]. 2024-05-12.
- [18] FLIR Systems. LWIR camera cost reduction strategies[R]. White Paper, 2023.
- [19] Tesla Inc. Camera field-of-view specifications[Z]. FSD Technical Manual, 2024.
- [20] Samsung Electronics. 14nm FSD chip fabrication process[R]. Technical Report, 2023.
- [21] Tesla AI Team. Shadow mode data collection and utilization[C]//NeurIPS Workshop on Autonomous Driving. 2024.
- [22] CHEN X, et al. BEVFormer: Bird's-eye-view transformer for autonomous driving[C]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2023: 12345-12354.
- [23] California DMV. Autonomous vehicle testing performance data 2024[DS]. Public Dataset, 2024.
- [24] NHTSA. Phantom braking incident report: DOT HS 813 532[R]. 2023.
- [25] Consumer Reports. Autonomous vehicle performance in heavy rain[R]. Technical Evaluation Report, 2024.
- [26] Chenyu Yang, Yuntao Chen, et al. "BEVFormer v2: Adapting Modern Image Backbones to Bird's-Eye-View Recognition via Perspective Supervision." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2023.
- [27] Anthony Hu, Zak Muze, et al. "FIERY: Future Instance Prediction in Bird's-Eye View from Surround Monocular Cameras." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2021.
- [28] Prakash A, Chitta K, Geiger A. Multi-modal fusion transformer for end-to-end autonomous driving[C]//Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2021: 7077-7087.

- [29] Waymo LLC. Urban object recall rates in heavy rain[R]. Safety Report SR-2024-02, 2024.
- [30] California DMV. Waymo disengagement events analysis 2024[R]. Public Report, 2024.
- [31] Waymo LLC. High-definition map update cycle analysis[R]. Technical Document TD-2024-05, 2024.
- [32] Waymo Research Team. "The Waymo Autonomous Driving System: A Generative Agent for Learning from Autonomous Driving Data." arXiv preprint arXiv: 2403. 17020 (2024).
- [33] Liang, Zheng, et al. "PVTransformer: A Novel Multi-Head Attention Based 3D Object Detection Framework for Autonomous Driving." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2023.
- [34] Waymo LLC. "Waymo Sensor Suite and Data Format." Waymo Open Dataset Documentation. [Online]. Available: <https://waymo.com/open/dataset/>
- [35] Arbe Robotics. Phoenix 4D radar specifications[Z]. Product Datasheet, 2024.
- [36] RoboSense. M1 solid-state LiDAR technical specifications[Z]. Product Manual, 2023.
- [37] Huawei Technologies. MDC 810 compute platform redundancy design[R]. Technical Report, 2023.
- [38] ZHANG Z, et al. Cross-modal validation for radar-camera fusion[J]. IEEE Robotics and Automation Letters, 2023, 8(3): 2345-2352.
- [39] Huawei Technologies. Dynamic path planning for complex intersections[R]. White Paper, 2023.
- [40] Huawei Semiconductor. 22nm 4D radar chip yield analysis[R]. Internal Report, 2024.
- [41] NHTSA. Dynamic object density metrics for urban driving[R]. Technical Note TN-2024-01, 2024.
- [42] Mobileye. Perception distance requirements for highway autonomy[R]. Technical Report, 2023.
- [43] BCG. Cost targets for consumer ADAS and robotaxi systems[R]. Automotive Market Analysis, 2023.
- [44] Waymo LLC. San Francisco narrow lane navigation requirements[R]. Technical Document TD-2024-03, 2024.
- [45] Mobileye. Chauffeur system requirements for highway autonomy[R]. White Paper, 2023.
- [46] Euro NCAP. German rainstorm test results for autonomous vehicles: Test Report TR-2023-08[R]. 2023.
- [47] Mobileye. Cost analysis for 4D radar-based systems[R]. Investor Briefing, 2023.
- [48] LEE S, et al. Polarization filtering for snow reflection suppression[J]. IEEE Transactions on Intelligent Vehicles, 2024, 9(1): 45-56.
- [49] FLIR Systems. Thermal imaging for pedestrian detection in snow[Z]. Application Note AN-0003, 2023.
- [50] Huawei Technologies. Low-cost sensor fusion for extreme weather[R]. Technical Report, 2024.
- [51] Innoviz Technologies. Solid-state LiDAR roadmap and cost projections[R]. Investor Presentation, 2025.
- [52] Arbe Robotics. 4D radar cost reduction strategies[R]. White Paper, 2024.
- [53] ZHOU D, et al. YOLO-Fastest: An ultra-lightweight object detector for edge devices[J]. IEEE Transactions on Intelligent Vehicles, 2023, 8(4): 789-801.
- [54] Camilleri J. 'Capitalist tools in socialist hands'? China Mobile in global financial networks[J].

Transactions of the Institute of British Geographers, 2015, 40(4): 464-478.

[55] Huawei Technologies. High-density 4D radar point cloud generation: CN115678901A[P]. 2024-06-30.

[56] Samsung Semiconductor. Event camera with 120dB dynamic range[EB/OL]. [2024-09-30].

[57] LI X, et al. Federated learning for distributed autonomous driving systems[J]. IEEE Transactions on Intelligent Transportation Systems, 2023, 24(6): 5678-5690.

[58] PEARL J. Causal inference in robotics: From correlation to causation[M]. Cambridge: MIT Press, 2024.

[59] BCG. L4 autonomous system cost projections 2023-2027[R]. Automotive Technology Outlook, 2023.