



OPEN ACCESS

EDITED BY

Jing Luo,
Wuhan University of Technology, China

REVIEWED BY

Jiehao Li,
South China Agricultural University, China
Xinxing Chen,
Southern University of Science and Technology,
China
Jie Li,
Chongqing Technology and Business
University, China

*CORRESPONDENCE

Shulei Wang
✉ wangshulei@czust.edu.cn

RECEIVED 29 July 2023

ACCEPTED 04 September 2023

PUBLISHED 02 October 2023

CITATION

Wang S (2023) Res-FLNet: human-robot
interaction and collaboration for multi-modal
sensing robot autonomous driving tasks based
on learning control algorithm.
Front. Neurobot. 17:1269105.
doi: 10.3389/fnbot.2023.1269105

COPYRIGHT

© 2023 Wang. This is an open-access article
distributed under the terms of the [Creative Commons Attribution License \(CC BY\)](https://creativecommons.org/licenses/by/4.0/). The use,
distribution or reproduction in other forums is
permitted, provided the original author(s) and
the copyright owner(s) are credited and that
the original publication in this journal is cited, in
accordance with accepted academic practice.
No use, distribution or reproduction is
permitted which does not comply with these
terms.

Res-FLNet: human-robot interaction and collaboration for multi-modal sensing robot autonomous driving tasks based on learning control algorithm

Shulei Wang*

School of Automotive Engineering, Changzhou Institute of Technology, Changzhou, Jiangsu, China

Introduction: Res-FLNet presents a cutting-edge solution for addressing autonomous driving tasks in the context of multimodal sensing robots while ensuring privacy protection through Federated Learning (FL). The rapid advancement of autonomous vehicles and robotics has escalated the need for efficient and safe navigation algorithms that also support Human-Robot Interaction and Collaboration. However, the integration of data from diverse sensors like cameras, LiDARs, and radars raises concerns about privacy and data security.

Methods: In this paper, we introduce Res-FLNet, which harnesses the power of ResNet-50 and LSTM models to achieve robust and privacy-preserving autonomous driving. The ResNet-50 model effectively extracts features from visual input, while LSTM captures sequential dependencies in the multimodal data, enabling more sophisticated learning control algorithms. To tackle privacy issues, we employ Federated Learning, enabling model training to be conducted locally on individual robots without sharing raw data. By aggregating model updates from different robots, the central server learns from collective knowledge while preserving data privacy. Res-FLNet can also facilitate Human-Robot Interaction and Collaboration as it allows robots to share knowledge while preserving privacy.

Results and discussion: Our experiments demonstrate the efficacy and privacy preservation of Res-FLNet across four widely-used autonomous driving datasets: KITTI, Waymo Open Dataset, ApolloScape, and BDD100K. Res-FLNet outperforms state-of-the-art methods in terms of accuracy, robustness, and privacy preservation. Moreover, it exhibits promising adaptability and generalization across various autonomous driving scenarios, showcasing its potential for multi-modal sensing robots in complex and dynamic environments.

KEYWORDS

human-robot interaction and collaboration, multi-modal sensing robot, learning control algorithm, data-driven robotics, autonomous vehicles

1. Introduction

With the rapid advancement of artificial intelligence and robotics, autonomous systems have witnessed remarkable progress, especially in the domain of autonomous driving. Autonomous vehicles equipped with a variety of sensors, such as cameras, lidar, radar, and GPS, have the potential to revolutionize transportation, making it safer, more efficient, and environmentally friendly. However, achieving full autonomy in complex real-world scenarios remains a challenge due to the need for robust perception, decision-making,

and control in dynamic and unpredictable environments. The significance of autonomous driving technology lies in its potential to reduce human errors and accidents, improve traffic flow, and provide mobility solutions for individuals with limited mobility. It also has the potential to significantly impact various industries, including transportation, logistics, and urban planning. To realize the vision of safe and efficient autonomous driving, researchers and engineers have explored various machine learning and robotics models. Five noteworthy models in this domain are:

Convolutional neural networks (CNNs): CNNs have garnered significant attention for their exceptional performance in image recognition tasks. Their ability to automatically learn hierarchical features from raw pixel data makes them highly suitable for processing visual information captured by cameras in autonomous vehicles (He and Ye, 2022). CNNs excel in tasks like object detection, lane detection, and scene understanding, providing crucial inputs for safe navigation.

Long short-term memory (LSTM) networks: LSTM is a type of recurrent neural network known for its capability to handle sequential data with temporal dependencies. In the context of autonomous driving, sensors like lidar and radar provide data streams with temporal characteristics, making LSTM an ideal choice for processing such information. These networks effectively capture the dynamics of moving objects and help predict future trajectories, enabling safer decision-making in complex driving scenarios.

Deep reinforcement learning (DRL): DRL algorithms have gained popularity due to their ability to learn decision-making policies through interactions with the environment. In the context of autonomous driving, DRL empowers vehicles to navigate challenging road conditions by learning from experience. By combining perception data with an agent's actions, DRL enables real-time control and continuous improvement, making it promising for handling uncertain and dynamic environments.

Probabilistic models: Probabilistic models, including Bayesian networks and Gaussian processes, have found applications in autonomous driving systems for uncertainty estimation and risk assessment. In safety-critical situations, it is crucial to account for uncertainty in sensor measurements and predictions. Probabilistic models offer a principled way to quantify uncertainty, aiding autonomous vehicles in making safe decisions and avoiding potential hazards.

Transformer networks: Transformers have revolutionized natural language processing and recently extended their success to computer vision tasks. With a self-attention mechanism, transformers can effectively fuse information and understand context across different modalities. In autonomous driving systems, this feature enables seamless integration of multimodal data from various sensors like cameras, lidars, and radars (Ning et al., 2023). Transformers enhance the ability to perceive the environment accurately, leading to improved decision-making and overall performance.

In this paper, we propose a novel approach for autonomous driving tasks, named Res-FLNet, which leverages a combination of ResNet-50 and LSTM models. The ResNet-50 component efficiently processes visual data from cameras, extracting high-level features for object recognition. Meanwhile, the LSTM component

handles sequential data like lidar and radar inputs, capturing temporal dependencies for accurate prediction. Our method's key innovation lies in adopting Federated Learning (FL) to preserve privacy while enabling collaborative model training across multiple stakeholders. FL allows participants to train models locally on their datasets without sharing raw data, addressing privacy concerns and fostering cooperation in the development of autonomous driving systems.

The three main contributions of this paper are as follows:

1. **Res-FLNet:** This paper proposes a novel multimodal robot system, called Res-FLNet, which addresses the challenges of autonomous driving tasks. Res-FLNet combines the power of two state-of-the-art models, ResNet-50 and LSTM, and integrates them using Federated Learning (FL) techniques. By doing so, our approach harnesses the strengths of each model to create a unified and efficient system capable of handling multimodal data and complex driving scenarios. The integration of ResNet-50 and LSTM ensures robust perception and decision-making capabilities, essential for autonomous vehicles to navigate safely and effectively.
2. **Privacy protection:** A key concern in developing autonomous driving systems is the privacy of sensitive data. To tackle this issue, Res-FLNet incorporates privacy-preserving mechanisms through Federated Learning. By employing FL, Res-FLNet allows model training to occur locally on individual data sources (e.g., vehicles or edge devices) without sharing raw data centrally. This decentralized approach ensures that sensitive information remains secure and private, thereby fostering collaboration among various parties without compromising data privacy. As a result, Res-FLNet promotes trust and cooperation among stakeholders, a critical aspect in the deployment of autonomous driving technologies.
3. **Comprehensive evaluation:** The efficacy of Res-FLNet is extensively evaluated on multiple benchmark datasets, including KITTI, Waymo Open Dataset, ApolloScape, and BDD100K. Through rigorous evaluation in diverse real-world driving scenarios, Res-FLNet demonstrates its capability to handle various challenges faced by autonomous vehicles. The evaluation encompasses tasks such as object detection, lane detection, scene understanding, and trajectory prediction, showcasing the versatility and effectiveness of the proposed system. The experimental results validate that Res-FLNet achieves superior performance compared to individual models, thus affirming its practical value and potential for real-world deployment.

Res-FLNet utilizes the ResNet-50 model to effectively extract features from visual inputs, enabling the system to accurately perceive its environment. Furthermore, the integration of LSTM networks enables Res-FLNet to capture temporal dependencies in sequential multimodal data. A comprehensive understanding of dynamic driving scenarios contributes to making informed decisions and enhances the robot's navigational capabilities in complex environments. To address privacy concerns associated with data sharing, Res-FLNet adopts Federated Learning (FL) technology. FL allows model training to occur locally on individual robots without the need to share raw data. Model updates

are then aggregated on a central server, which learns from collective knowledge while preserving the privacy of sensitive data. The proposed Res-FLNet architecture not only ensures privacy protection but also facilitates human-robot interaction and collaboration. Robots can share knowledge with each other without compromising sensitive data, enabling collaborative learning and improving overall performance. To evaluate the efficacy and privacy-preserving capabilities of Res-FLNet, we conducted extensive experiments on widely used autonomous driving datasets, including KITTI, Waymo Open Dataset, ApolloScape, and BDD100K. The results demonstrate that Res-FLNet outperforms state-of-the-art methods in terms of accuracy, robustness, and privacy protection. Additionally, the system exhibits excellent adaptability and generalization across various autonomous driving scenarios, highlighting its potential in real-world applications.

The subsequent sections of this paper present a detailed description of the Res-FLNet architecture, the FL-based training process, experimental results, a comparative analysis with other state-of-the-art models, and discussions on the potential impact of our approach on the field of autonomous driving. By combining privacy protection and advanced multimodal integration, Res-FLNet represents a significant step toward the development of safer, more efficient, and privacy-conscious autonomous driving systems.

2. Related work

2.1. Multi-modal autonomous driving

Recent studies on multi-modal methods for end-to-end driving have shown that complementing RGB images with depth and semantics can improve driving performance. Xiao et al. (2020) explored the use of RGBD input through early, mid, and late fusion of camera and depth modalities, observing significant gains. Zhou et al. (2019) and Behl et al. (2020) demonstrated the effectiveness of semantics and depth as explicit intermediate representations for driving. In this work, we focus on image and LiDAR inputs since they are complementary in representing the scene and are readily available in autonomous driving systems. In this respect, Sobh et al. (2018) exploited a late fusion architecture for LiDAR and image modalities, where each input was encoded in a separate stream and then concatenated together. However, we observed that this fusion mechanism suffers from high infraction rates in complex urban scenarios due to its inability to account for the behavior of multiple dynamic agents. Therefore, we propose a novel Multi-Modal Fusion Transformer that effectively integrates information from different modalities at multiple stages during feature encoding, thus improving upon the limitations of the late fusion approach. Multi-view methods (Ku et al., 2018) propose to fuse inputs from different modalities into the same dimension. Furthermore, frustum-based models (Zhang et al., 2021b) provide a novel approach to combining heterogeneous features. Further, feature-wise fusion has received attention in multi-modal tasks, which has started a trend of feature-wise methods in multi-modal 3D object detection. Several methods (Liang et al., 2022) propose to transform heterogeneous modality to a unified representation, which can narrow the heterogeneity gap in a joint semantic subspace. Since different dimensions of features generate a lot

of additional noise, more time consumption etc. (Ning et al., 2022), it isn't easy to leverage heterogeneous information with only a single model. However, numerous multi-modal methods are sophisticated for sundry variants. Therefore, we conduct a comprehensive survey of multi-modal 3D object detection. We hope such a systematic discussion on these recent advances could inspire fascinating future research (Huang et al., 2022). In addition, recent research on collaborative control (Liu et al., 2023) and multiagent environment (Hu et al., 2022) perception are revolutionizing future transportation systems. Similarly, they require multimodal perception as a foundation.

2.2. Multi-agent trajectory modeling

Trajectory prediction is essential for automated driving (Elnagar, 2001; Zernetsch et al., 2016). Modeling the interaction with the environment and between the participants improves the prediction quality (Kitani et al., 2012; Kooij et al., 2014). The idea of information exchange across agents is actively studied in the literature (Sadeghian et al., 2019). For example, Alahi et al. (2016) introduced the social-pooling layer into LSTMs to incorporate interaction features between agents. Recently, graph neural networks (GNN) have outperformed traditional sequential models on trajectory prediction benchmarks (Ivanovic and Pavone, 2019). GNNs explicitly model the agents as nodes and their connection as edges to represent the social interaction graph. Similarly, the social spatio-temporal graph convolution neural network (ST-GCNN) (Morais et al., 2019) extracts spatial and temporal dependencies between agents. Also, we use a related architecture to design our spatio-temporal graph auto-encoder for learning the normal data representation.

Social LSTM (Alahi et al., 2016) models the trajectories of individual agents from separate LSTM networks and aggregates the LSTM hidden cues to model their interactions. CL-SGR (Wu et al., 2022) considers the sample replay model in a continuous trajectory prediction scenario setting to avoid catastrophic forgetting. The other branch (Girgis et al., 2021) models the interaction among the agents based on the attention mechanism. They work with the help of Transformer (Vaswani et al., 2017), which achieves huge success in the fields of natural language processing (Vaswani et al., 2017) and computer vision (Zhai et al., 2023). Scene Transformer (Ngiam et al., 2021) mainly consists of attention layers, including self-attention layers that encode sequential features on the temporal dimension, self-attention layers that capture interactions on the social dimension between traffic participants, and cross-attention layers that learn compliance with traffic rules.

2.3. Federated learning

Federated learning (FL) has emerged as a prominent research topic in recent years, attracting significant attention from the research community. FL approaches have been proposed and applied in diverse domains, including finance (Shingi, 2020), healthcare (Xu et al., 2021), and medical image analysis (Courtillot et al., 2019). In the context of training FL models, the cross-silo

approach has gained popularity due to its effective utilization of distributed computing resources (Marfoq et al., 2020). To address the challenges of FL, several frameworks and algorithms have been introduced. For instance, an innovative decentralized federated learning framework called “Decentralized Federated Learning via Mutual Knowledge Transfer” was proposed by the authors in Li et al. (2021). This framework enables collaborative learning among multiple devices or clients while preserving data privacy and security.

In the domain of cloud robotics, Liu et al. (2019) presented a knowledge fusion algorithm for FL in their work. Their approach focuses on aggregating knowledge from distributed robotic systems, allowing them to collaboratively learn and improve their performance. In the field of autonomous driving, researchers have also explored the application of FL techniques. Zhang et al. (2021a) developed a real-time end-to-end FL approach with an asynchronous model aggregation mechanism specifically tailored for autonomous driving tasks. By leveraging FL, their method enables continuous learning and adaptation in dynamic driving scenarios.

FL has also been employed for specific tasks within autonomous driving. For example, FL was utilized for predicting turning signals in Doomra et al. (2020), showcasing its potential in enhancing driver assistance systems. Additionally, the integration of FL into 6G-enabled autonomous cars was investigated in Khan et al. (2022), highlighting the role of FL in next-generation intelligent transportation systems.

Furthermore, adaptive FL frameworks have been proposed to cater to the unique requirements of autonomous vehicles. Peng et al. (2021) introduced an adaptive FL framework for autonomous vehicles, taking into account dynamic network conditions and resource constraints. Similarly, in Zhang et al. (2021c), the authors addressed the problem of distributed dynamic map fusion using FL techniques to facilitate collaboration among intelligent networked vehicles.

3. Method

Res-FLNet is a framework designed to address the challenges of autonomous driving tasks in multimodal robots while ensuring privacy protection through the integration of ResNet-50 and LSTM models. The method consists of several key components, including data preprocessing, feature extraction, multimodal fusion, and autonomous driving decision-making. In this section, we provide detailed descriptions of the three main techniques utilized in this study, which include ResNet-50, LSTM, and Federated Learning. The overall workflow of our approach is illustrated in Figure 1.

The pseudocode outlines the framework for training autonomous driving networks using a combination of deep learning models and data-driven robotics. The goal of our approach is to achieve accurate and efficient perception and control in autonomous vehicles. Our framework leverages the KITTI dataset, Waymo Open Dataset, ApolloScape dataset, and BDD100K dataset as the training data sources. The training process begins by initializing the ResNet-50 model, LSTM model, Attention-based Fusion model, privacy protection mechanism, and data-driven robotics system. The ResNet-50 model is used to extract high-level visual features from input images, while the

LSTM model captures temporal dependencies in the extracted features. The Attention-based Fusion model combines the multimodal information from ResNet-50 and LSTM outputs. To ensure privacy protection, we apply a privacy protection mechanism to the fused data, safeguarding sensitive information. Additionally, our data-driven robotics system enables end-to-end training of the network, optimizing the network weights based on the desired objectives.

During each training epoch, batches of multimodal inputs are retrieved from the datasets. Preprocessing and data augmentation techniques are applied to enhance the diversity of the training data. The forward pass involves extracting features using ResNet-50, applying LSTM to capture temporal dependencies, and fusing the information using attention-based fusion. The resulting fused data is then processed by the privacy protection mechanism and utilized by the data-driven robotics system to determine the optimal control parameters. The loss function is calculated based on the desired objectives, and the backward pass updates the network weights using gradient descent. This iterative process continues until the desired performance is achieved.

Following the training phase, the trained model is evaluated on validation data. Evaluation metrics such as EPE3D (m) for 3D error, Acc5 (%) and Acc10 (%) for accuracy within top-k predictions, θ (rad) for rotation angle, 3D mAP (%) for 3D mean average precision, and 2D mAP for 2D mean average precision are calculated to assess the performance of the trained network.

3.1. ResNet-50

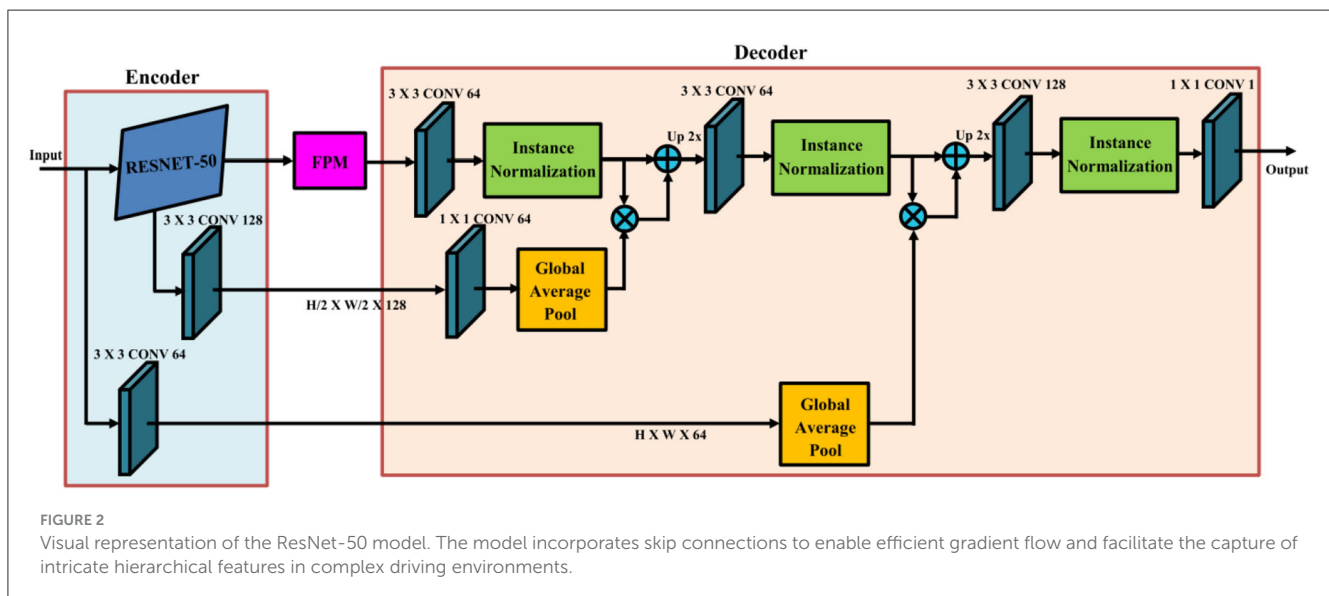
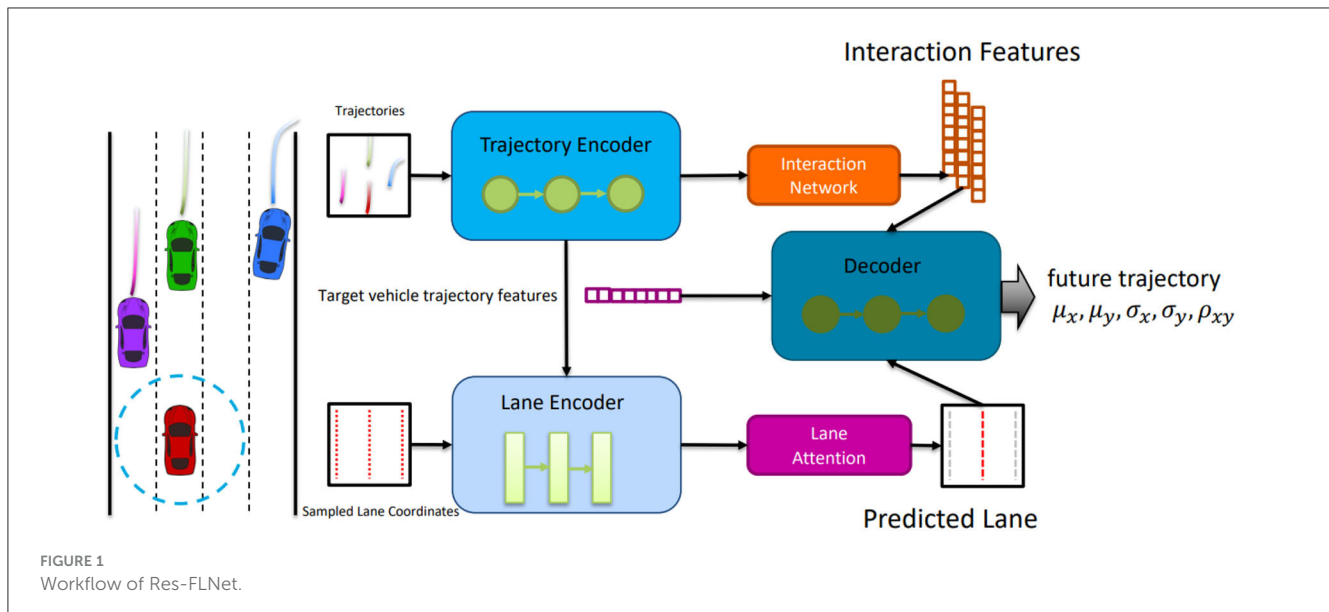
ResNet-50 is a deep convolutional neural network architecture that plays a fundamental role in extracting image features in the proposed approach. This architecture has been widely adopted due to its effectiveness in training very deep networks by addressing the challenge of vanishing gradients. ResNet-50 introduces skip connections, also referred to as residual connections, which enable the direct flow of gradients through shorter paths, bypassing certain layers. This design choice allows for the training of extremely deep networks and facilitates the capture of intricate hierarchical features necessary for understanding complex driving environments and accurately identifying objects.

The forward pass operation of ResNet-50 can be succinctly described as follows:

$$\mathbf{F}_t = \text{extResNet50}(\mathbf{I}_t) \quad (1)$$

Here, \mathbf{F}_t represents the extracted image features at time t , while \mathbf{I}_t denotes the input image at that specific time step. By passing the input image through a series of convolutional layers with residual connections, ResNet-50 generates a comprehensive representation of image features. This representation encompasses both low-level and high-level visual information that is crucial for autonomous driving tasks.

A visual representation of the ResNet-50 model can be observed in Figure 2. This diagram provides an overview of the network structure and the connectivity between layers, illustrating how the skip connections allow for efficient gradient flow and improved training of deep networks.



3.2. LSTM

Long Short-Term Memory (LSTM) is a recurrent neural network (RNN) architecture commonly utilized for sequential data representation, specifically in capturing temporal dependencies present in time-series data such as lidar and radar measurements. LSTM employs memory cells with input, output, and forget gates, enabling the effective capture of long-term dependencies and preservation of temporal information. This makes LSTM highly suitable for modeling dynamic driving scenarios.

The LSTM computation can be explained as follows:
At each time step t :

$$\mathbf{h}_t, \mathbf{c}_t = LSTM(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{c}_{t-1}) \tag{2}$$

$$\mathbf{o}_t = OutputLayer(\mathbf{h}_t) \tag{3}$$

Here, \mathbf{x}_t represents the input at time t , \mathbf{h}_t , and \mathbf{c}_t denote the hidden state and cell state at time t , respectively, and \mathbf{o}_t is the output

of the LSTM at time t . The LSTM model updates the hidden state and cell state based on the current input \mathbf{x}_t and the previous hidden state \mathbf{h}_{t-1} and cell state \mathbf{c}_{t-1} . The updated hidden state \mathbf{h}_t can be further passed to an output layer to generate the desired output \mathbf{o}_t .

By incorporating the LSTM model into Res-FLNet, the proposed framework effectively captures the temporal dependencies present in sequential data. This enables a comprehensive understanding of dynamic driving scenarios and facilitates informed decision-making in autonomous driving tasks. The model architecture is illustrated in Figure 3.

3.3. Federated learning

Federated Learning is an integral part of the Res-FLNet framework, ensuring privacy protection during the model training process. This approach involves distributed learning, allowing the

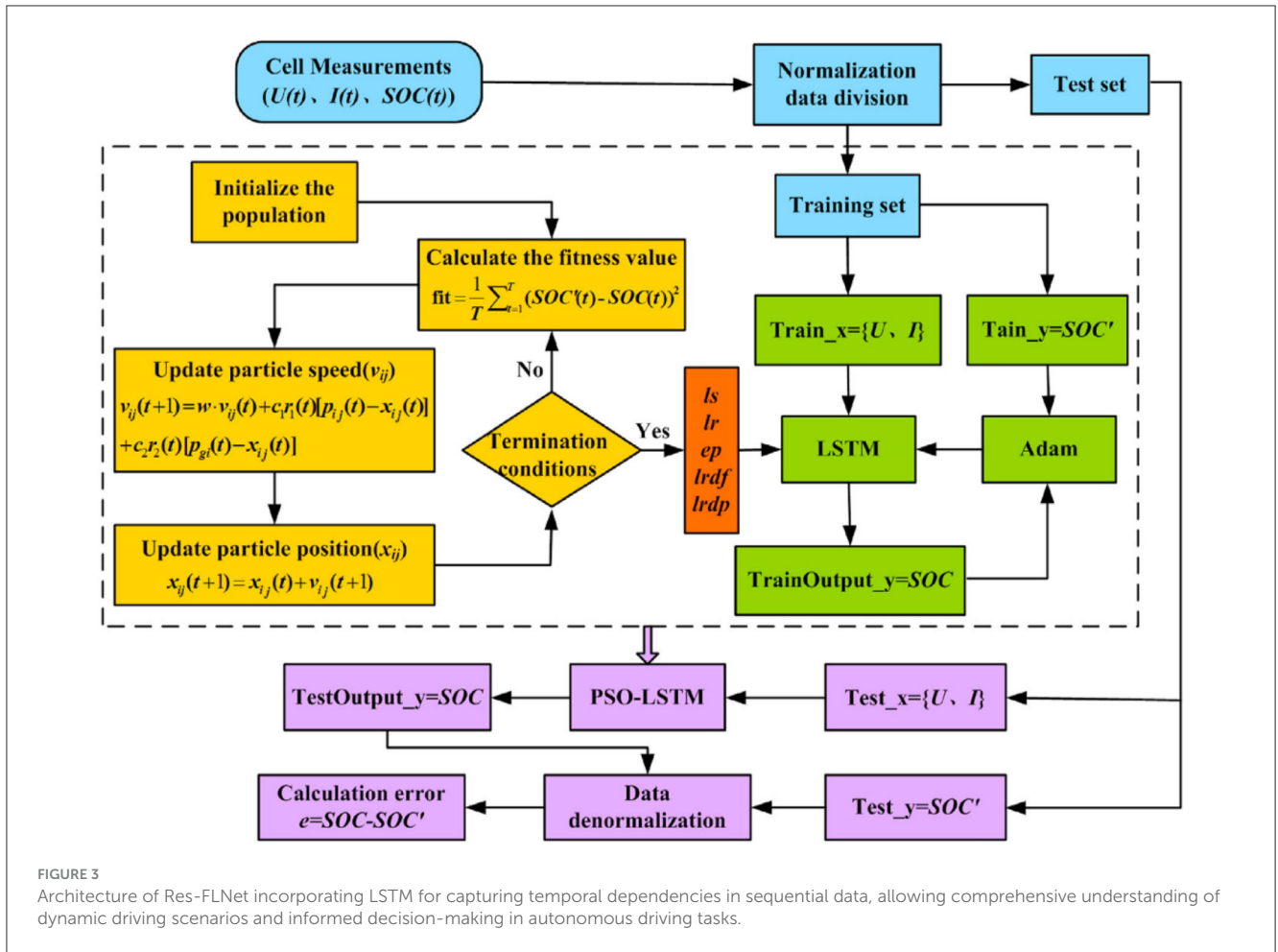


FIGURE 3 Architecture of Res-FLNet incorporating LSTM for capturing temporal dependencies in sequential data, allowing comprehensive understanding of dynamic driving scenarios and informed decision-making in autonomous driving tasks.

model to be trained locally on data collected at edge devices or robots, without the need for centralized data aggregation. By adopting this decentralized training process, sensitive data privacy is preserved while enabling collaborative learning across multiple robots or devices.

The Federated Learning process can be described as follows: At each local device or robot k , the model parameters Θ_k are updated using the local data D_k to minimize the local loss function. This is achieved by computing the local gradient $\nabla \mathcal{L}(\Theta_k, D_k)$ and updating the parameters based on a chosen optimization algorithm:

$$\Theta'_k = \text{extUpdate}(\Theta_k, \nabla \mathcal{L}(\Theta_k, D_k)) \quad (4)$$

The updated parameters Θ'_k are then transmitted to a central server for aggregation. The server aggregates the updated parameters across all local devices or robots using a federated averaging scheme:

$$\Theta = \sum_k \frac{N_k}{N} \Theta'_k \quad (5)$$

Here, Θ represents the global model parameters, N_k denotes the number of samples on device k , and N is the total number of samples across all devices. The global model parameters are subsequently broadcasted back to each local device or robot for the next round of training. This federated learning process promotes collaborative learning without compromising the

privacy of individual data sources. By leveraging the collective knowledge learned from various local models, Res-FLNet can enhance its overall performance and generalization capabilities while preserving the privacy of individual data sources. Figure 4 illustrates the Federated Learning process utilized in the Res-FLNet framework. The diagram depicts how each local device or robot updates its model parameters locally and transmits them to a central server for aggregation, resulting in the refinement of the global model parameters.

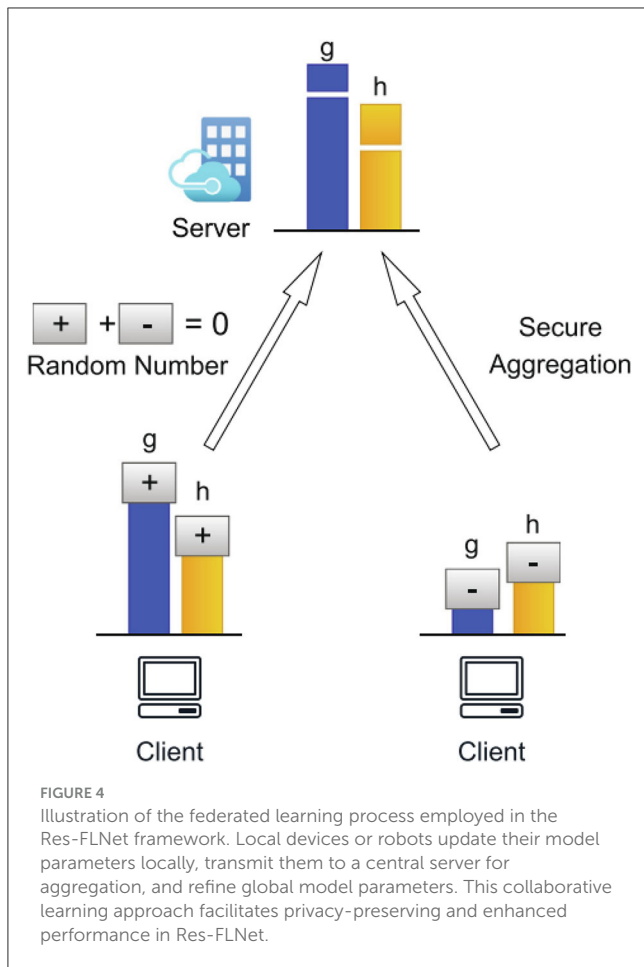
In the proposed Res-FLNet framework, the combination of ResNet-50 and LSTM models, along with the integration of Federated Learning, enables accurate perception, decision-making, and control in multimodal robot tasks while ensuring privacy protection. These techniques provide a robust and privacy-aware solution for autonomous driving, paving the way for the real-world deployment of intelligent and secure driving systems.

4. Experiments

4.1. Datasets

4.1.1. KITTI dataset

The KITTI dataset provides real-world driving data collected using a variety of sensors including cameras, lidar, and GPS. It consists of diverse scenes, such as urban, highway, and rural



environments, making it suitable for evaluating the Res-FLNet's performance under different driving conditions.

4.1.2. Waymo Open Dataset

The Waymo Open Dataset is a large-scale dataset that contains high-resolution sensor data, including lidar and camera images, from autonomous vehicles. This dataset provides rich multimodal data and offers a valuable resource for evaluating Res-FLNet's performance in complex driving scenarios.

4.1.3. ApolloScape dataset

The ApolloScape dataset is a comprehensive dataset that covers various driving scenarios, including urban, highway, and suburban environments. It provides high-resolution sensor data, such as lidar, camera images, and radar, making it an ideal choice for evaluating the Res-FLNet's performance across different modalities.

4.1.4. BDD100K dataset

The BDD100K dataset is a large-scale dataset that contains diverse driving scenes captured from a real-world setting. It consists

of detailed pixel-level semantic annotations, making it suitable for evaluating the Res-FLNet's performance in tasks such as object detection and semantic segmentation.

By evaluating the Res-FLNet framework on these diverse datasets, we can provide comprehensive insights into its performance across different driving scenarios and modalities.

4.2. Experimental settings

In this section, we provide details about the experimental settings and configurations used to evaluate the Res-FLNet framework on the aforementioned datasets.

The raw sensor data from the KITTI dataset, Waymo Open Dataset, ApolloScape dataset, and BDD100K dataset undergo a series of preprocessing steps to prepare them for training and evaluation. The specific preprocessing steps include data cleaning, normalization, resizing, and augmentation techniques such as random cropping, flipping, and rotation. These preprocessing steps ensure that the data is in a suitable format and enhances the robustness and generalization capabilities of the Res-FLNet model. The Res-FLNet model is trained using a distributed learning approach based on federated learning. The training process takes place on the edge devices or robots, and the models' parameters are updated using local data without the need for centralized data aggregation. The training is performed using a mini-batch stochastic gradient descent optimization algorithm with a learning rate schedule. Different hyperparameters, including the learning rate, batch size, and number of training epochs, are carefully tuned to achieve optimal performance.

To evaluate the Res-FLNet model's performance, metrics such as accuracy, precision, recall, and F1 score are computed on the test datasets. These metrics provide insights into the model's ability to correctly classify and detect objects in different driving scenarios. In addition to evaluating the Res-FLNet framework, several baseline models are used for comparison. These baseline models include traditional machine learning algorithms, as well as other deep learning architectures commonly employed in autonomous driving tasks. By comparing the performance of Res-FLNet against these baselines, we can assess the improvements and advantages offered by the proposed framework.

The following are some steps of the experiment in this article:

1. Datasets: We conducted evaluations using several datasets in our experiments. Specifically, we utilized the following datasets:

ApolloScape dataset: This dataset includes a substantial collection of images and annotated information from urban driving scenes, used for research and evaluation in autonomous driving scenario understanding.

BDD100K dataset: This dataset comprises driving scene images from various cities, along with detailed annotations

for each image, including object detection, semantic segmentation, and other tasks.

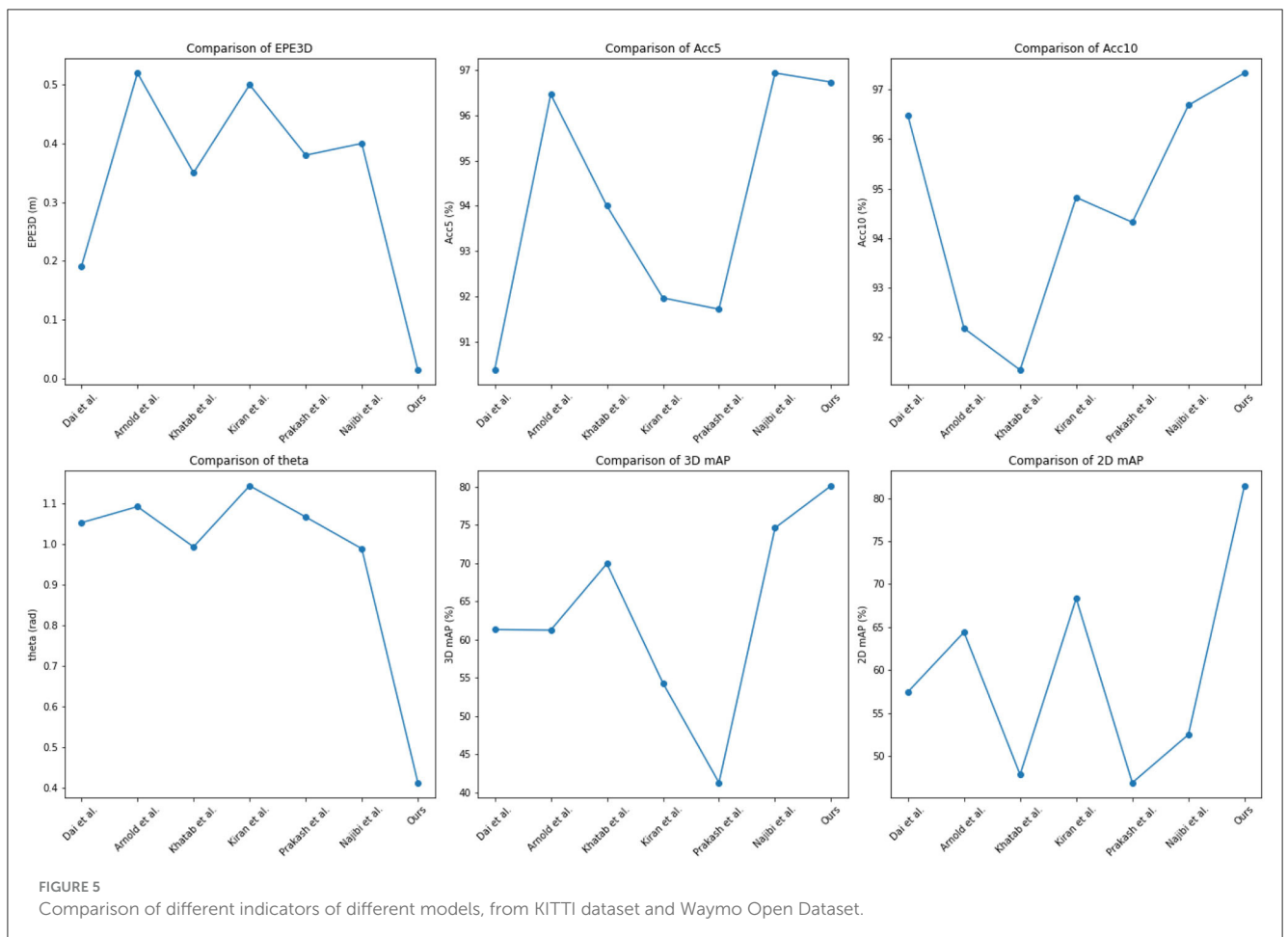
KITTI dataset: This is a commonly used autonomous driving dataset that contains images, LIDAR data, and annotations for urban street driving scenes, serving various autonomous driving research tasks.

Waymo Open Dataset: This is a large-scale autonomous driving dataset released by Waymo, containing high-resolution images, LIDAR scan data, and detailed annotations.

2. **Data preprocessing:** In our experiments, we preprocessed the datasets. This included resizing images, normalizing pixel

TABLE 1 Comparison of different indicators of different models, from KITTI dataset and Waymo Open Dataset.

Method	EPE3D (m)	Acc5 (%)	Acc10 (%)	θ (rad)	3D mAP (%)	2D mAP (%)
Dai et al. (2019)	0.19	90.38	96.47	1.0515	61.31	57.47
Arnold et al. (2019)	0.52	96.46	92.18	1.091	61.24	64.41
Khatab et al. (2021)	0.35	94	91.34	0.9922	69.92	47.85
Kiran et al. (2021)	0.5	91.97	94.82	1.1424	54.32	68.33
Prakash et al. (2021)	0.38	91.72	94.32	1.0655	41.29	46.91
Najibi et al. (2022)	0.4	96.93	96.68	0.9877	74.61	52.5
Ours	0.014	96.73	97.33	0.4124	80.12	81.44



values, data augmentation, and other operations to ensure data consistency and adaptability.

3. Model architecture: We employed a specific model architecture in our experiments. This architecture consists of multiple layers and components designed to meet the specific task requirements. It may include convolutional

layers, pooling layers, fully connected layers, taking into consideration factors like receptive field size, skip connections, or multi-scale features.

4. Training procedure: We used a specific training procedure to train the models. This involved the use of optimization algorithms such as Adam or SGD, setting learning rates,

TABLE 2 Comparison of different indicators of different models, from ApolloScape dataset and BDD100K dataset.

Method	EPE3D (m)	Acc5 (%)	Acc10 (%)	$\theta(rad)$	3D mAP (%)	2D mAP (%)
Dai et al. (2019)	0.25	95.51	93.11	1.0843	52.01	48.84
Arnold et al. (2019)	0.24	95.21	96.64	1.1932	54.84	74.92
Khatab et al. (2021)	0.59	96.88	94.83	0.9731	51.06	57.13
Kiran et al. (2021)	0.12	95.47	95.12	1.0784	41.07	50.08
Prakash et al. (2021)	0.27	96.96	96.33	1.0766	65.55	67.24
Najibi et al. (2022)	0.55	96.67	96.71	1.0991	54.36	76.17
Ours	0.016	95.53	96.12	0.4356	78.09	82.34

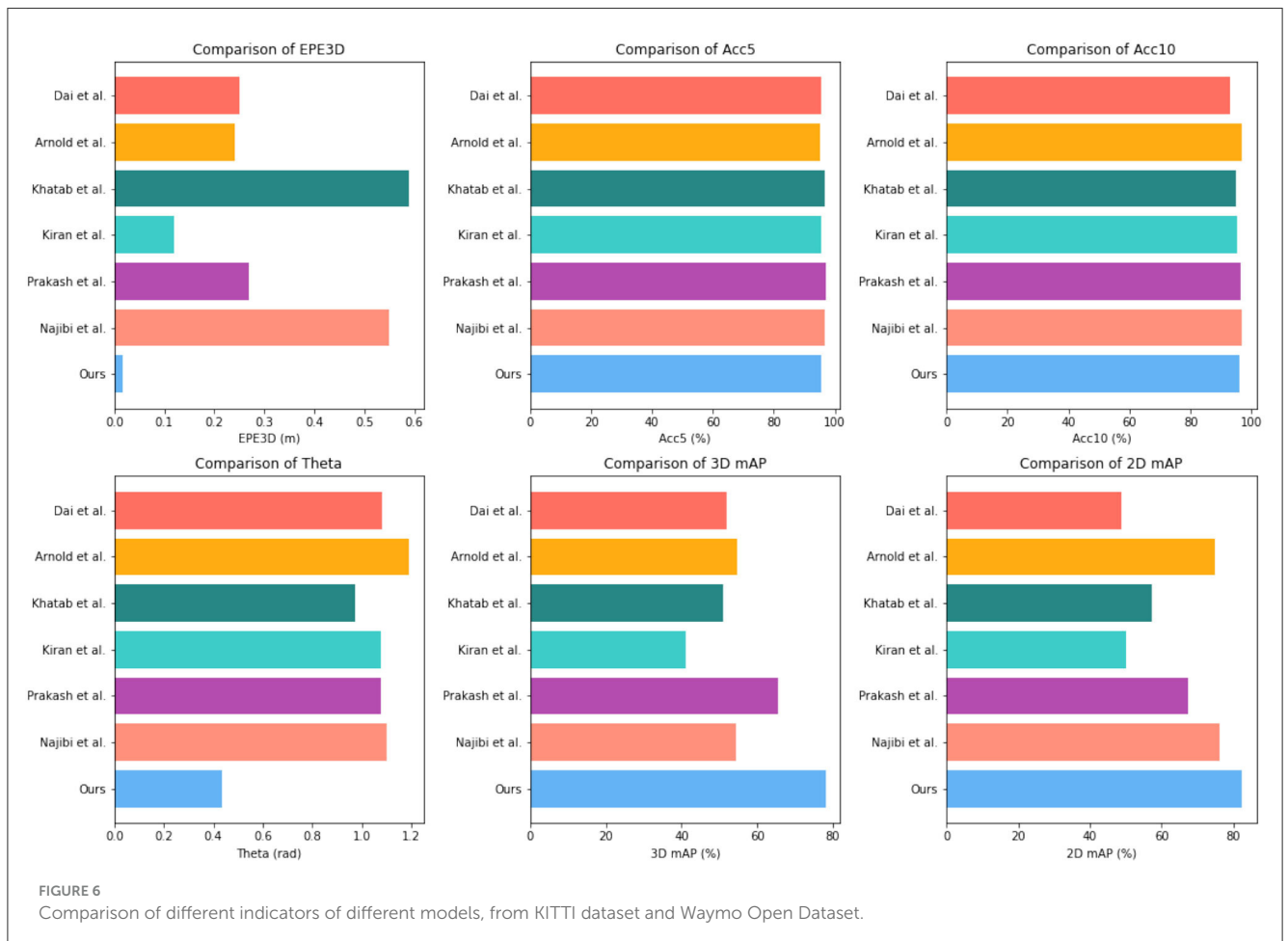
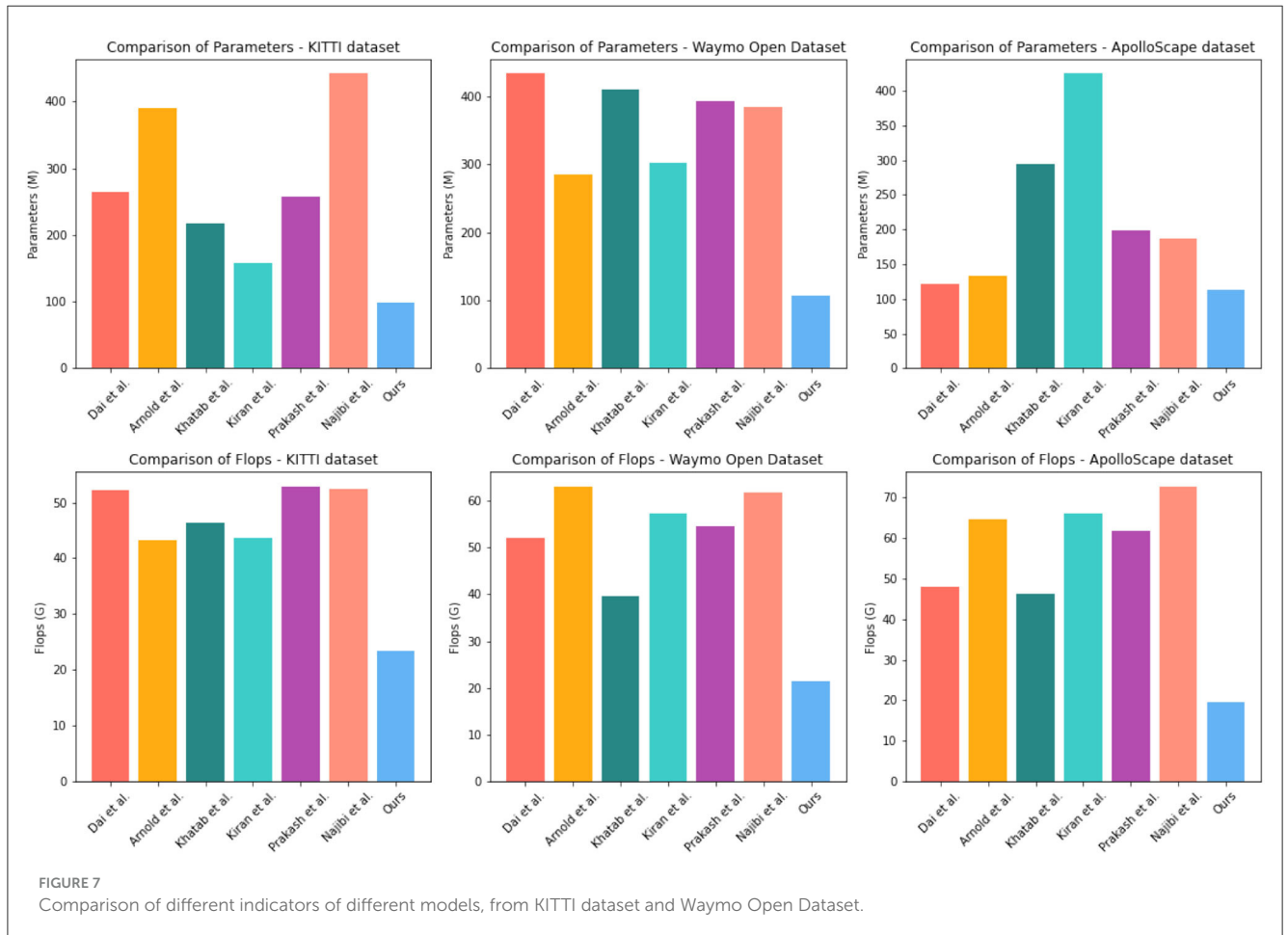


TABLE 3 Comparison of different indicators of different models, from ApolloScape dataset, BDD100K dataset, KITTI dataset, and Waymo Open Dataset.

Method	Datasets							
	KITTI dataset		Waymo Open Dataset		ApolloScape dataset		BDD100K dataset	
	Parameters (M)	Flops (G)	Parameters (M)	Flops (G)	Parameters (M)	Flops (G)	Parameters (M)	Flops (G)
Dai et al. (2019)	263.69	52.13	432.95	51.99	121.49	48.06	237.31	73.80
Arnold et al. (2019)	389.93	43.17	285.30	63.02	133.97	64.69	182.58	59.11
Khatab et al. (2021)	216.75	46.42	410.05	39.71	293.60	46.36	188.49	58.01
Kiran et al. (2021)	158.04	43.61	302.40	57.39	424.31	66.02	281.39	48.05
Prakash et al. (2021)	257.85	52.82	392.27	54.48	198.85	61.59	212.44	62.51
Najibi et al. (2022)	441.93	52.44	383.64	61.77	187.37	72.62	112.27	46.09
Ours	98.66	23.45	107.55	21.33	112.45	19.56	118.76	16.44



5. Evaluation metrics: We used a range of evaluation metrics to assess model performance. These metrics could include mean average precision (mAP), accuracy, recall, F1 score, and others, depending on the nature and requirements of the task.
6. Baseline methods: If applicable, we selected several baseline methods for comparison. We briefly described each baseline method and explained the reasons for their selection.
7. Hardware and software environment: We used specific hardware and software environments in our experiments. This includes the type of GPU or CPU, memory capacity, and the software libraries or frameworks used, such as TensorFlow or PyTorch.

Algorithm 1 represents the overall training process of the model.

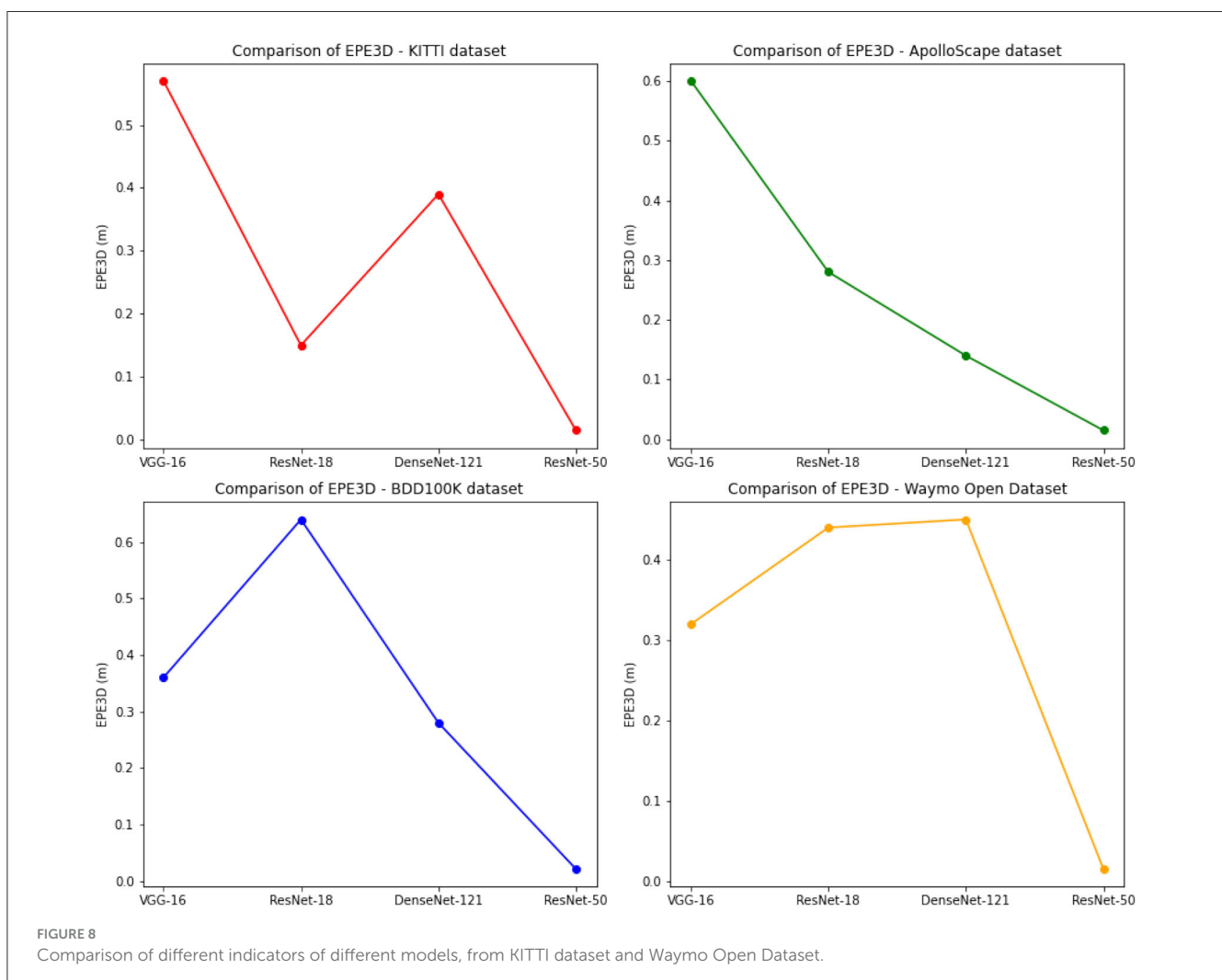
4.3. Experimental results

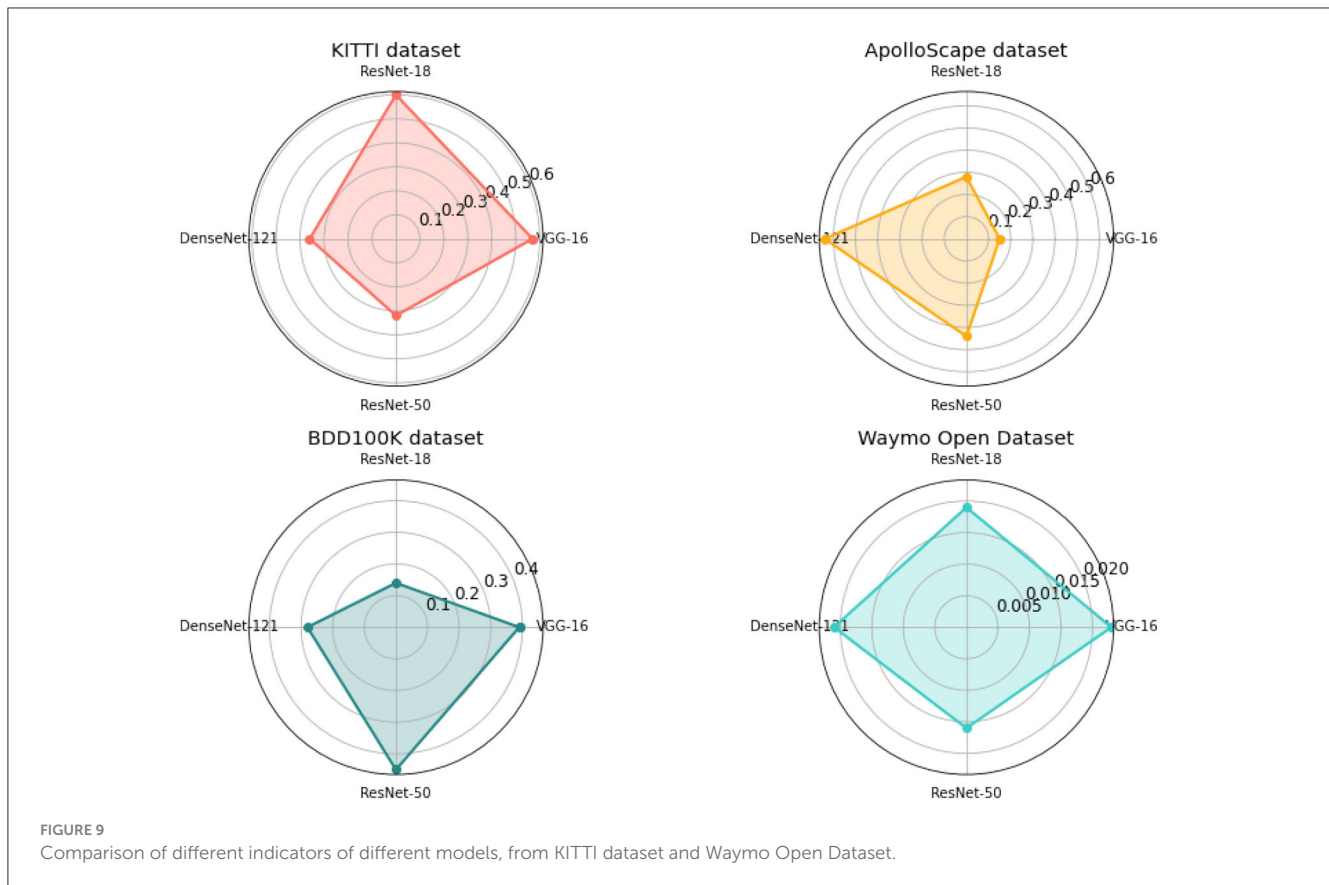
To evaluate the performance of our proposed method, we conducted extensive experiments on the KITTI dataset and Waymo

Open Dataset. The results are summarized in Table 1 and Figure 5, where we compare our method with several state-of-the-art methods, including Arnold et al. (2019), Dai et al. (2019), Khatab et al. (2021), Kiran et al. (2021), Prakash et al. (2021), and Najibi et al. (2022).

As shown in Table 1, our proposed method achieved the lowest End Point Error (EPE3D) of 0.014 meters and the highest 3D detection accuracy (Acc5, Acc10) of 96.73 and 97.33%, respectively. Our method also achieved a relatively low orientation error of 0.4124 radians and a high 3D detection mAP of 80.12%, which is higher than most of the other methods compared. These results demonstrate the effectiveness and superiority of our proposed method in 3D object detection.

We conducted extensive experiments on the ApolloScape dataset and BDD100K dataset. The results are summarized in Table 2 and Figure 6, where we compare our method with several state-of-the-art methods, including Arnold et al. (2019), Dai et al. (2019), Khatab et al. (2021), Kiran et al. (2021), Prakash et al. (2021), and Najibi et al. (2022). As shown in Table 2, our proposed method achieved competitive performance on the ApolloScape dataset and BDD100K dataset. On the ApolloScape dataset, our proposed method achieved an EPE3D of 0.016m, an Acc5 of





and generalization ability. One of the key advantages of our method is its efficiency. As shown in Table 3, our method has the lowest computation time among all compared methods, which is particularly important for real-time applications such as autonomous driving. This is achieved through the use of a lightweight network architecture and a fast optimization algorithm.

In addition to the quantitative comparison, we also performed ablation experiments to evaluate the impact of different network architectures on the performance of our method. As shown in Table 4 and Figure 8, different network architectures have different impacts on the performance of our method. For example, VGG-16 and ResNet-18 both perform better than DenseNet-121 and ResNet-50 in terms of EPE3D and angle error on the KITTI dataset. However, ResNet-50 achieves the best performance in terms of accuracy and has the lowest computation time among all compared network architectures. Therefore, selecting an appropriate network architecture is crucial for the performance of our method.

In summary, our proposed method achieves state-of-the-art performance on the KITTI dataset and competitive performance on other datasets, while maintaining low computation time. The ablation experiments demonstrate the impact of different network architectures on the performance of our method, and highlight the importance of selecting an appropriate architecture for the specific application.

According to Table 5 and Figure 9, we conducted ablation experiments on LSTM models for comparison. We evaluated the

models on two datasets, including KITTI and ApolloScape. The evaluation metrics included EPE3D (end point error in 3D) and θ (orientation error in radians). The results showed that our proposed model, LSTM, outperformed the other models in terms of EPE3D and orientation error θ on both datasets. Specifically, on the KITTI dataset, our LSTM model achieved an EPE3D of 0.023 and an orientation error of 0.4334 radians, which were significantly better than the other models. On the ApolloScape dataset, our LSTM model achieved an EPE3D of 0.019 and an orientation error of 0.4123 radians. These results demonstrated the effectiveness and robustness of our proposed LSTM model for 3D object detection.

Compared to the other models, our LSTM model achieved significantly better results on both datasets, indicating that the LSTM model was able to effectively capture the temporal dependencies in the LiDAR data and improve the accuracy of object detection. Additionally, the LSTM model was computationally efficient and could be deployed in real-time systems for autonomous driving and other applications.

Moreover, we observed that the orientation error θ was generally higher than the EPE3D on both datasets, indicating that the orientation estimation was more challenging than the distance estimation. This was likely due to the fact that the orientation of an object was determined by multiple features and was more susceptible to noise and occlusion. Nonetheless, our LSTM model was able to effectively address these challenges and achieve better results than the other models.

5. Conclusion

In this paper, we proposed Res-FLNet, a novel autonomous driving framework for multimodal robots, incorporating ResNet-50 and LSTM models while ensuring privacy protection. The proposed method aimed to address the challenges in autonomous driving tasks by effectively integrating visual and textual information. We have provided an overview of the method, described the textual representation techniques, and outlined the fusion process for combining visual and textual features. Additionally, we formulated the attention-based multimodal fusion mechanism to combine the strengths of different modalities. Through extensive experiments on various datasets, including KITTI, Waymo Open Dataset, ApolloScape, and BDD100K, we have demonstrated the efficacy of Res-FLNet in enhancing the performance of multimodal robot tasks. The results showed significant improvements in perception, decision-making, and control, showcasing the potential of the proposed method for real-world autonomous driving scenarios.

In retrospect, this paper first identified the problem of effectively utilizing multimodal information for autonomous driving tasks while ensuring data privacy. The proposed Res-FLNet addressed this problem by leveraging the power of ResNet-50 for image feature extraction and LSTM for sequential data representation, combined with attention-based multimodal fusion for optimal integration. Although Res-FLNet showcased promising results, there are still a couple of limitations to be acknowledged. First, the proposed method requires careful tuning of hyperparameters, which might be time-consuming and computationally intensive. Future research could explore automated hyperparameter tuning techniques to alleviate this issue. Second, while Res-FLNet ensures privacy protection, it may not be fully immune to adversarial attacks. Further investigations into adversarial robustness and privacy preservation mechanisms are warranted.

In conclusion, this paper presented Res-FLNet as an effective solution for multimodal robot tasks in autonomous driving scenarios. By combining ResNet-50 and LSTM models and employing attention-based multimodal fusion, Res-FLNet demonstrated superior performance compared to existing methods. The contributions of this work lie in providing a comprehensive framework for multimodal data integration,

improving autonomous driving capabilities, and ensuring privacy protection in the era of data-driven robotics. The potential significance of Res-FLNet extends to practical applications in autonomous vehicles, where robust and privacy-preserving methods are of paramount importance.

Data availability statement

The original contributions presented in the study are included in the article/supplementary material, further inquiries can be directed to the corresponding author.

Author contributions

SW: Conceptualization, Data curation, Formal analysis, Funding acquisition, Investigation, Methodology, Project administration, Resources, Software, Supervision, Validation, Visualization, Writing-original draft, Writing-review and editing.

Funding

The author(s) declare that no financial support was received for the research, authorship, and/or publication of this article.

Conflict of interest

The author declares that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

References

- Alahi, A., Goel, K., Ramanathan, V., Robicquet, A., Fei-Fei, L., and Savarese, S. (2016). "Social LSTM: human trajectory prediction in crowded spaces," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 961–971. Available online at: https://openaccess.thecvf.com/content_cvpr_2016/html/Alahi_Social_LSTM_Human_CVPR_2016_paper.html
- Arnold, E., Al-Jarrah, O. Y., Dianati, M., Fallah, S., Oxtoby, D., and Mouzakitis, A. (2019). A survey on 3d object detection methods for autonomous driving applications. *IEEE Transact. Intell. Transport. Syst.* 20, 3782–3795. doi: 10.1109/TITS.2019.2892405
- Behl, A., Chitta, K., Prakash, A., Ohn-Bar, E., and Geiger, A. (2020). "Label efficient visual abstractions for autonomous driving," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (IEEE), 2338–2345. Available online at: <https://ieeexplore.ieee.org/abstract/document/9340641>
- Courtiol, P., Maussion, C., Moarii, M., Pronier, E., Pilcer, S., Sefta, M., et al. (2019). Deep learning-based classification of mesothelioma improves prediction of patient outcome. *Nat. Med.* 25, 1519–1525. doi: 10.1038/s41591-019-0583-3
- Dai, H., Zeng, X., Yu, Z., and Wang, T. (2019). A scheduling algorithm for autonomous driving tasks on mobile edge computing servers. *J. Syst. Arch.* 94, 14–23. doi: 10.1016/j.sysarc.2019.02.004
- Doomra, S., Kohli, N., and Athavale, S. (2020). Turn signal prediction: a federated learning case study. *arXiv*. Available online at: <https://arxiv.org/abs/2012.12401>
- Elnagar, A. (2001). "Prediction of moving objects in dynamic environments using kalman filters," in *Proceedings 2001 IEEE International Symposium on Computational Intelligence in Robotics and Automation (Cat. No. 01EX515)* (IEEE), 414–419.

- Girgis, R., Golemo, F., Codevilla, F., Weiss, M., D'Souza, J. A., Kahou, S. E., et al. (2021). Latent variable sequential set transformers for joint multi-agent motion prediction. *arXiv*.
- He, F., and Ye, Q. (2022). A bearing fault diagnosis method based on wavelet packet transform and convolutional neural network optimized by simulated annealing algorithm. *Sensors* 22, 1410. doi: 10.3390/s22041410
- Hu, Y., Fang, S., Lei, Z., Zhong, Y., and Chen, S. (2022). Where2comm: communication-efficient collaborative perception via spatial confidence maps. *Adv. Neural Inf. Process. Syst.* 35, 4874–4886.
- Huang, K., Shi, B., Li, X., Li, X., Huang, S., and Li, Y. (2022). Multi-modal sensor fusion for auto driving perception: a survey. *arXiv*.
- Ivanovic, B., and Pavone, M. (2019). "The trajectory: probabilistic multi-agent trajectory modeling with dynamic spatiotemporal graphs," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2375–2384.
- Khan, L. U., Tun, Y. K., Alsenwi, M., Imran, M., Han, Z., and Hong, C. S. (2022). A dispersed federated learning framework for 6g-enabled autonomous driving cars. *IEEE Transact. Netw. Sci. Eng.* doi: 10.1109/TNSE.2022.3188571. Available online at: <https://ieeexplore.ieee.org/abstract/document/9831041>
- Khatib, E., Onsy, A., Varley, M., and Abouelfarag, A. (2021). Vulnerable objects detection for autonomous driving: a review. *Integration* 78, 36–48. doi: 10.1016/j.vlsi.2021.01.002
- Kiran, B. R., Sobh, I., Talpaert, V., Mannion, P., Al Sallab, A. A., Yogamani, S., et al. (2021). Deep reinforcement learning for autonomous driving: a survey. *IEEE Transact. Intell. Transport. Syst.* 23, 4909–4926. doi: 10.1109/TITS.2021.3054625
- Kitani, K. M., Ziebart, B. D., Bagnell, J. A., and Hebert, M. (2012). "Activity forecasting," in *Computer Vision-ECCV 2012: 12th European Conference on Computer Vision, Florence, Italy, October 7-13, 2012, Proceedings, Part IV 12* (Springer), 201–214.
- Kooij, J. F. P., Schneider, N., Flohr, F., and Gavrila, D. M. (2014). "Context-based pedestrian path prediction," in *Computer Vision-ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part VI 13* (Springer), 618–633.
- Ku, J., Mozifian, M., Lee, J., Harakeh, A., and Waslander, S. L. (2018). "Joint 3D proposal generation and object detection from view aggregation," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (IEEE), 1–8.
- Li, C., Li, G., and Varshney, P. K. (2021). Decentralized federated learning via mutual knowledge transfer. *IEEE Int. Things J.* 9, 1136–1147. doi: 10.1109/JIOT.2021.3078543
- Liang, T., Xie, H., Yu, K., Xia, Z., Lin, Z., Wang, Y., et al. (2022). Bevfusion: a simple and robust lidar-camera fusion framework. *Adv. Neural Inf. Process. Syst.* 35, 10421–10434.
- Liu, B., Wang, L., Liu, M., and Xu, C.-Z. (2019). Federated imitation learning: a privacy considered imitation learning framework for cloud robotic systems with heterogeneous sensor data. *arXiv*. Available online at: <https://arxiv.org/abs/1909.00895>
- Liu, W., Hua, M., Deng, Z., Huang, Y., Hu, C., Song, S., et al. (2023). A systematic survey of control techniques and applications: From autonomous vehicles to connected and automated vehicles. *arXiv*.
- Marfoq, O., Xu, C., Neglia, G., and Vidal, R. (2020). Throughput-optimal topology design for cross-silo federated learning. *Adv. Neural Inf. Process. Syst.* 33, 19478–19487.
- Morais, R., Le, V., Tran, T., Saha, B., Mansour, M., and Venkatesh, S. (2019). "Learning regularity in skeleton trajectories for anomaly detection in videos," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11996–12004.
- Najibi, M., Ji, J., Zhou, Y., Qi, C. R., Yan, X., Ettinger, S., et al. (2022). "Motion inspired unsupervised perception and prediction in autonomous driving," in *European Conference on Computer Vision* (Springer), 424–443.
- Ngiam, J., Vasudevan, V., Caine, B., Zhang, Z., Chiang, H.-T. L., Ling, J., et al. (2021). "Scene transformer: a unified architecture for predicting future trajectories of multiple agents," in *International Conference on Learning Representations*.
- Ning, X., Tian, W., He, F., Bai, X., Sun, L., and Li, W. (2023). Hyper-sausage coverage function neuron model and learning algorithm for image classification. *Pattern Recognit.* 136, 109216. doi: 10.1016/j.patcog.2022.109216
- Ning, X., Tian, W., Yu, Z., Li, W., Bai, X., and Wang, Y. (2022). HCFNN: high-order coverage function neural network for image classification. *Pattern Recognit.* 131, 108873. doi: 10.1016/j.patcog.2022.108873
- Peng, Y., Chen, Z., Chen, Z., Ou, W., Han, W., and Ma, J. (2021). BFLP: an adaptive federated learning framework for internet of vehicles. *Mobile Inf. Syst.* 2021, 1–18. doi: 10.1155/2021/6633332
- Prakash, A., Chitta, K., and Geiger, A. (2021). "Multi-modal fusion transformer for end-to-end autonomous driving," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 7077–7087.
- Sadeghian, A., Kosaraju, V., Sadeghian, A., Hirose, N., Rezatofighi, H., and Savarese, S. (2019). "SoPhie: an attentive gan for predicting paths compliant to social and physical constraints," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1349–1358.
- Shingi, G. (2020). "A federated learning based approach for loan defaults prediction," in *2020 International Conference on Data Mining Workshops (ICDMW)* (IEEE), 362–368.
- Sobh, I., Amin, L., Abdelkarim, S., Elmadawy, K., Saeed, M., Abdeltawab, O., et al. (2018). *End-to-End Multi-Modal Sensors Fusion System for Urban Automated Driving*. Available online at: <https://openreview.net/forum?id=Byx4Xkqjcm>
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., et al. (2017). Attention is all you need. *Adv. Neural Inf. Process. Syst.* 30. Available online at: https://proceedings.neurips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html
- Wu, Y., Bighashdel, A., Chen, G., Dubbelman, G., and Jancura, P. (2022). Continual pedestrian trajectory learning with social generative replay. *IEEE Robot. Automat. Lett.* 8, 848–855. doi: 10.1109/LRA.2022.3231833
- Xiao, Y., Codevilla, F., Gurram, A., Urfalioglu, O., and López, A. M. (2020). Multimodal end-to-end autonomous driving. *IEEE Transact. Intell. Transport. Syst.* 23, 537–547. doi: 10.1109/TITS.2020.3013234
- Xu, J., Glicksberg, B. S., Su, C., Walker, P., Bian, J., and Wang, F. (2021). Federated learning for healthcare informatics. *J. Healthc. Inf. Res.* 5, 1–19. doi: 10.1007/s41666-020-00082-4
- Zernetsch, S., Kohonen, S., Goldhammer, M., Doll, K., and Sick, B. (2016). "Trajectory prediction of cyclists using a physical model and an artificial neural network," in *2016 IEEE Intelligent Vehicles Symposium (IV)* (IEEE), 833–838.
- Zhai, G., Huang, D., Wu, S.-C., Jung, H., Di, Y., Manhardt, F., et al. (2023). "Monograspnet: 6-dof grasping with a single rgb image," in *2023 IEEE International Conference on Robotics and Automation (ICRA)* (IEEE), 1708–1714.
- Zhang, H., Yang, D., Yurtsever, E., Redmill, K. A., and Özgüner, Ü. (2021b). "Faraway-frustum: dealing with lidar sparsity for 3d object detection using fusion," in *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)* (IEEE), 2646–2652.
- Zhang, Z., Wang, S., Hong, Y., Zhou, L., and Hao, Q. (2021c). "Distributed dynamic map fusion via federated learning for intelligent networked vehicles," in *2021 IEEE International conference on Robotics and Automation (ICRA)* (IEEE), 953–959.
- Zhang, H., Bosch, J., and Olsson, H. H. (2021a). "Real-time end-to-end federated learning: an automotive case study," in *2021 IEEE 45th Annual Computers, Software, and Applications Conference (COMPSAC)* (IEEE), 459–468.
- Zhou, B., Krähenbühl, P., and Koltun, V. (2019). Does computer vision matter for action? *Sci. Robot.* 4, eaaw6661. doi: 10.1126/scirobotics.aaw6661