

PNAS-MOT: Multi-Modal Object Tracking with Pareto Neural Architecture Search

Chensheng Peng¹, Zhaoyu Zeng², Jinling Gao², Jundong Zhou²
 Masayoshi Tomizuka¹, Xinbing Wang², Chenghu Zhou², Nanyang Ye²

¹ UC Berkeley ² Shanghai Jiao Tong University

Abstract—Multiple object tracking is a critical task in autonomous driving. Existing works primarily focus on the heuristic design of neural networks to obtain high accuracy. As tracking accuracy improves, however, neural networks become increasingly complex, posing challenges for their practical application in real driving scenarios due to the high level of latency. In this paper, we explore the use of the neural architecture search (NAS) methods to search for efficient architectures for tracking, aiming for low real-time latency while maintaining relatively high accuracy. Another challenge for object tracking is the unreliability of a single sensor, therefore, we propose a multi-modal framework to improve the robustness. Experiments demonstrate that our algorithm can run on edge devices within lower latency constraints, thus greatly reducing the computational requirements for multi-modal object tracking while keeping lower latency. Code is available at <https://github.com/PholyPeng/PNAS-MOT>.

Index Terms—Multiple Object Tracking, Neural Architecture Search

I. INTRODUCTION

MULTIPLE object tracking (MOT) is a fundamental task of consistently assigning a unique ID to each observed object within a video sequence, which holds significant importance across various domains, including motion planning, safe robot navigation, and autonomous driving [20]. The primary challenge inherent to MOT lies in establishing precise associations between tracklets from preceding frames and the object detections within the current frame. To tackle the complexities of multi-object tracking, two main-stream paradigms have emerged: tracking-by-detection [35], [43] and joint-tracking-and-detection [9], [33], [34]. The tracking-by-detection paradigm follows a two-stage process. Initially, a pre-trained detector is employed to procure object detections, after which a tracker undertakes the data association task, assigning a distinct ID to each detected object across successive frames. On the other hand, the joint-tracking-and-detection paradigm endeavors to achieve detection and tracking concurrently, leveraging the benefits of joint optimization strategies. In this paper, our focus lies specifically on the tracking aspect, so we adopt the tracking-by-detection approach due to its inherent efficiency and proven effectiveness in addressing the complexities of object tracking tasks.

In the task of Multiple Object Tracking, two primary sensor modalities, namely images and LiDAR point clouds, are extensively utilized. Various methodologies have been

devised to address MOT challenges leveraging these sensor modalities [2], [10], [35], [37], [42]. Recent studies have revealed the limitations of relying solely on a single modality, often resulting in mismatching issues. For example, cameras offer the advantage of capturing rich texture details but are susceptible to variations in lighting conditions. On the other hand, LiDAR sensors excel in capturing 3D geometry information and maintaining robustness in adverse weather conditions, yet with limitations on perceiving distant object detection owing to the sparsity of points. As a result, methods [15], [29], [38] utilizing the fusion features have been proposed. However, the improvement of accuracy from multi-modal fusion often comes at the cost of significant latency and high energy consumption. These factors present significant challenges when deploying such methods in practical scenarios, particularly within autonomous driving systems.

To address these challenges, we propose the integration of neural architecture search (NAS) techniques to search for efficient deep neural networks (DNNs) [8]. Early NAS methods primarily focused on improving accuracy by searching for network architectures within an expansive search space containing numerous structures. For example, DARTs [18] propose to formulate the searching task in a differentiable manner, focusing solely on the search for minimal network blocks, which are subsequently utilized to construct complete networks. Under the context of the challenging multi-modal MOT, we focus on latency-constrained NAS to search for architectures that exhibit superior performance while maintaining minimal latency. Specifically, we propose a Pareto optimization scheme to find the optimal Pareto frontier for the best trade-off between latency and accuracy. Subsequently, we select the network architectures that are capable of simultaneously achieving the desired levels of latency and accuracy. Moreover, our proposed multi-modal framework integrates information from both sensors to enhance robustness in handling challenging scenarios. Through a learned weighted feature fusion mechanism, the framework assigns greater significance to LiDAR point cloud features when cameras encounter difficulties or fail, and vice versa.

To summarize, our contributions are as follows:

- We propose a constrained neural architecture search (NAS) method that searches for network architectures capable of completing the MOT task within a specified time limit. This is solved by a Pareto frontier searching algorithmic scheme, which finds a suitable latency accuracy trade-off.

- We evaluate the proposed algorithmic framework on the KITTI benchmark and achieve 89.59% accuracy close to the SOTA methods while keeping the latency below 80 milliseconds on different edge devices.

II. RELATED WORK

A. Multiple Object Tracking

Recent MOT methods have made remarkable progress largely due to powerful deep neural networks. Sun *et al.* [30] concentrate on object affinity between different frames and propose an efficient online tracking network named Deep Affinity Network (DAN). Wang *et al.* [35] introduce an attention mechanism to fuse features from multiple modalities. However, object occlusion and overlapping significantly impair the accuracy of localizing and tracking objects. To overcome this issue, Ren *et al.* [26] introduce the Recurrent Rolling Convolution (RRC) architecture. Tracklet-Plane Matching (TPM), proposed by Peng *et al.* [24], is an approach to reduce the influence of noisy or confusing object detection and improve the performance of MOT.

Existing work [4], [36] propose methods for more efficient use of point cloud information for better object tracking performance, and [17] propose methods for improving the robustness of tracking. Wu *et al.* [39] utilize prediction confidence to guide data association and build a more robust tracker for objects temporarily missed by detectors. Dunnhofer *et al.* [7] propose a weakly-supervised adaptation strategy and utilize knowledge distillation to overcome inadequate tracking accuracy in many domains due to distribution shift and overfitting. Jong *et al.* [5] introduce APPLE MOTS, a dataset of homogeneous objects and propose TrackR-CNN and PointTrack architecture for joint detection and tracking. Fu *et al.* [11] propose a novel scale-aware domain adaptation framework, ScaleAwareDA, to solve the gap issue on object scale between the training and inference phases.

B. Multi-modal fusion

Multi-modal fusion remains a burgeoning area of exploration within the domain of autonomous driving [23], [40]. By combining features derived from diverse modalities, it offers a potent strategy to mitigate the limitations inherent in single-modality approaches in Multiple Object Tracking (MOT) tasks, including issues related to mismatching and unreliability. Addressing the challenge of weak-pairing characteristics in multi-modal fusion, Liu *et al.* [19] have made notable strides in enhancing the fusion effectiveness. A variety of modalities and sensors are employed for multimodal fusion regarding the application scenarios. For instance, Qu *et al.* [25] utilize visible, infrared, and hyperspectral sensors simultaneously to increase accuracy and robustness for object detection and tracking tasks. In a similar way, Zhang *et al.* [44] leverage data extracted from paired images and velocities, and then propose an efficient vehicle tracker, underscoring the versatility of multi-modal fusion techniques. EagerMOT, introduced by Kim *et al.* [15], exemplifies the integration of information from depth sensors and cameras, facilitating the identification and localization of distant incoming objects with

heightened precision and reliability. Moreover, Xu *et al.* [41] delve into the semantic-level fusion of 2D RGB images and 3D point clouds, aiming to enhance detection performance through a nuanced integration strategy.

C. Neural Architecture Search

Neural Architecture Search (NAS) is a method for automated neural network search particularly in pursuit of high accuracy. NAS is usually implemented with two categories of approaches, *i.e.*, reward-based methods (*e.g.*, Reinforcement Learning) and gradient-based methods. Zoph *et al.* [45] utilize reinforcement learning to train an RNN network that generates the model descriptions of neural networks. Liu *et al.* [18] introduce a differential search space that allows the use of gradient descent for architecture search. Constrained-NAS are proposed to satisfy relatively insufficient computing resources for limited computing platforms (*e.g.*, Edge-GPU, FPGA). Tan *et al.* [31] introduce latency constraints into the mobile neural architecture search (MNAS) approach and achieve high accuracy on mobile equipment. However, a trade-off between resource consumption and accuracy deteriorates the network accuracy and fails to meet the constraints strictly. To address this, Nayman *et al.* [22] further strengthens the latency constraints and proposes the HardCoRe-NAS method, which satisfies the constraint tightly without sacrificing accuracy.

III. METHODS

A. Problem Statement

In this paper, we follow the widely adopted tracking-by-detection paradigm, where a set of object detections is first obtained from a pre-trained detector. Let the detection results in T consecutive frames of a video sequence be $\mathcal{O} = \{O^1, \dots, O^t, \dots, O^T\}$, where $O^t = \{o_j^t\}$, $j = 0, 1, \dots, N^t$, and $o_j^t = (b_j^t, t)$. N^t is the total number of detections of the t -th frame, and b_j^t is usually a 2D bounding box on image plane for 2D detection results, while a 3D bounding box for 3D detection results, and t is the time stamp. The goal of multiple object tracking is to assign an ID for each instance detection o_j^t . A tracklet is defined as a set of object detections in different frames, $\tau = \{o_{k_1}^{t_1}, o_{k_2}^{t_2}, o_{k_3}^{t_3}, \dots\}$. The objective of consistent ID assignment can be interpreted as finding a set of tracklets $\mathcal{T} = \{\tau_i\}$, $i = 0, 1, \dots, n$, that can best explain the object detections. The problem is defined as predicting the correct tracklets \mathcal{T} by maximizing the conditional probabilities given the set of object detections \mathcal{O} .

$$\mathcal{T}^* = \arg \max_{\mathcal{T}} \mathbb{P}(\mathcal{T} | \mathcal{O}). \quad (1)$$

For a successful tracking process, the ID assignment must be consistent and unique for the same instance in a video sequence. As shown in Figure 2, the three tracked objects whose IDs are 0, 227, and 229 respectively, are a successful tracking procedure, because the assigned IDs are consistent in two consecutive frames $t-1$ and t . However, for the car whose ID is 220 in frame $t-1$, the ID is assigned to the blue vehicle at the left bottom of the image, but it's wrongly assigned to the black vehicle in frame t .

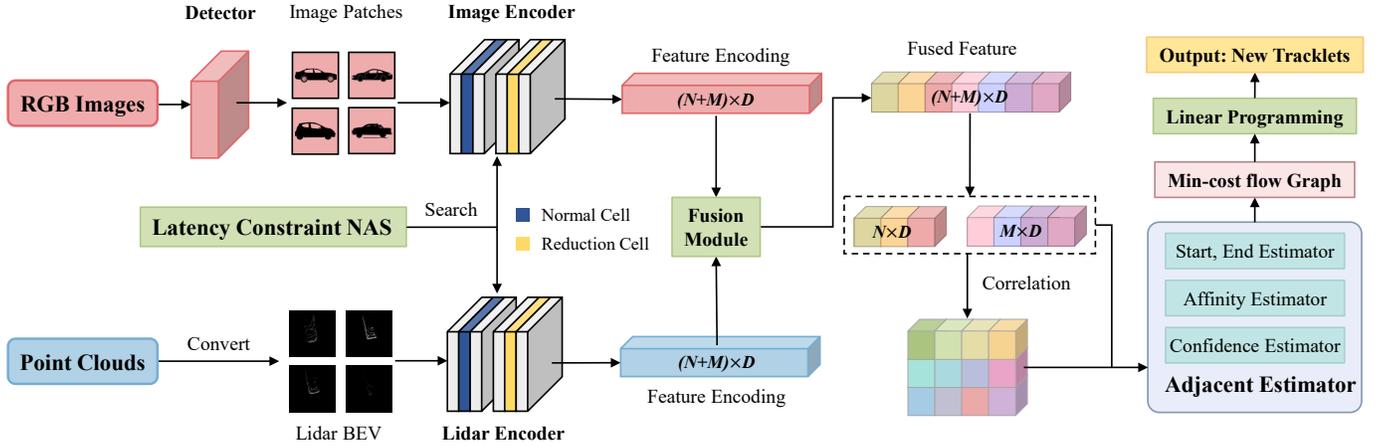


Fig. 1. The structure overview of the tracking framework. We take image patches and LiDAR point clouds as input, which are processed respectively by a corresponding encoder. A fusion module is then used to fuse the multi-modal features. The fused features are split into **tracklet features** and **detection features**, from which the **correlation features** can be calculated. These features are input to the following adjacent estimator to infer confidence and affinity scores. At last, linear programming is used to update the new tracklets from the inferred scores.

B. Tracking Network Structure

For the tracking network, we take the image and point clouds as input. Through a detector, we obtain the object detections in frame t . To reduce the amount of data, we crop the instance detections from the full image and then resize them to image patches whose size is 256×256 . In 3D spaces, we remove the point clouds outside the bounding box and convert the remaining points to Birds Eye View (BEV) images whose resolution is 256×256 .

Next, as shown in Figure 1, in an online setting, tracklets $T_{t-1} = \{o_n^{t-1}\}$, where $n = 1, 2, \dots, N$, with N being the number of tracklets in frame $t-1$ and object detections $O^t = \{o_m^t\}$, where $m = 1, 2, \dots, M$, with M being the number of detections in frame t , are fed to the feature extraction network, which consists of two branches, an image encoder, and a LiDAR encoder, each comprised of a normal cell and a reduction cell searched from the latency constraint NAS method. We can get the feature encoding from each modality and then add them together to obtain the fused feature, whose size is $(N+M) \times D$.

During the following correlation process, the fused feature is split to a $N \times D$ vector and a $M \times D$ vector. Then the two vectors are used to calculate the correlation features between existing tracklets and current detections. In this work, the difference calculation is defined as $\mathbf{D}_{n,m} = |f_n^{t-1} - f_m^t|$, where f_n^{t-1} and f_m^t are the features of the n -th object in frame $t-1$ and the m -th in frame t respectively.

Following previous works [10], [16], [28], [43], MOT can be solved as a min-cost flow problem. Together with the fused features, correlation features are then fed to the adjacency estimator adapted from [43], which infers the corresponding new, end, confidence, and affinity scores. The predicted scores are used to construct a graph with each edge representing the cost. Then, we use the mixed integer programming (MIP) solver provided by Google OR-Tools¹ to find the optimal solution in the min-cost graph, generating the updated tracklets

by adding the connected detections to their corresponding tracklets or initiating a new tracklet.

During the data association process, for the matched tracklets and detections, the tracklet will be updated by appending the detection to the tracklet. The object detection will obtain the ID from the corresponding tracklet through the ID propagation process. For the unmatched detection, it will not be instantly initialized as a new tracklet, instead, a birth check procedure is used to check whether it's a wrong detection or not. Only when the detection appeared in t_{birth} consecutive frames, the detection will be initialized. Similarly, for the unmatched tracklets, we won't delete them immediately, instead, we check if they disappeared in t_{death} frames.

C. Hardware latency function

To enable efficient NAS, we need to estimate the latency of the network architecture on specific hardware platforms. We test the latency of different operations. We first list all the possible *in_channels*, *out_channels*, *resolutions*, and other parameters used in our architectures. Every set of parameters is defined as a configuration. We set a random input with the corresponding input size of every configuration and pass it to the operation. The time spent is recorded by a preset timer. After hundreds of repetitions of such processes and an averaging procedure, we obtain a dictionary of the average latency per iteration for every operation with different configurations.

The architecture encoding α_i for each operation op_i , where $i = 1, 2, \dots, n$ with n being the number of operations, can be obtained through a Softmax function, then the estimated latency can be obtained from a weighted sum:

$$\text{Lat}(\alpha) = \sum_{i=1}^n \alpha_i \cdot \text{Lat}(op_i),$$

where the $\text{Lat}(op_i)$ can be easily acquired through a table-lookup process.

¹<https://developers.google.com/optimization>

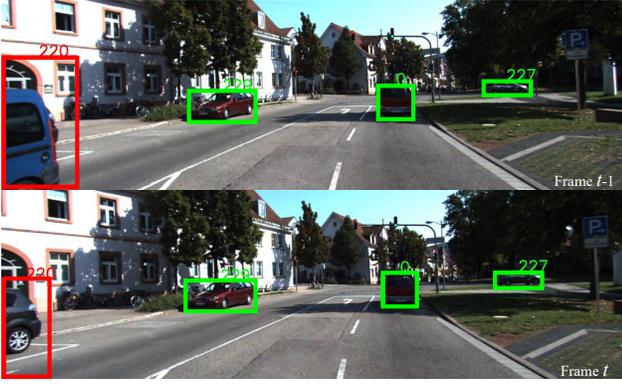


Fig. 2. Successful and failed examples of multi-object tracking. **Green bounding box:** The assigned IDs 0, 227, and 229 are consistent in two consecutive frames $t-1$ and t , which are successful trackings. **Red bounding box:** The ID 220 is wrongly assigned to the black vehicle in frame t , which is a failure tracking. Under the tracking-by-detection setting, Multi-object tracking can be regarded as a downstream detection task. It is hard to track dynamic objects under occlusion, which is common in real driving scenarios and significant for safe autonomous driving.

D. Pareto Optimization

Fig. 3 shows the structure overview of the MOT framework. We aim to find an optimal DNN that suits the latency accuracy trade-off. To achieve this trade-off, we design a two-stage Pareto optimization scheme.

For the first stage, we search for efficient network structures with the proposed constrained-NAS method. The backbone architecture to be searched is parameterized as α , and the search space is denoted as \mathcal{S} . Therefore, the complete network parameters are $\zeta = (\alpha, \theta)$, consisting of the architecture α and the corresponding weights θ . Thus, the optimization problem is formulated as follows with tracking accuracy and latency objectives,

$$\min_{\alpha \in \mathcal{S}, \theta} \{\mathbb{E}_{\mathcal{O}, \mathcal{T} \sim \mathcal{D}_{val}} [\mathcal{L}(\mathcal{O}, \mathcal{T} | \theta, \alpha)]\} \quad (2)$$

where \mathcal{D}_{val} denotes the distribution of the validation dataset, and $\mathcal{L} = \mathcal{L}_{track} + \lambda \mathcal{L}_{latency}$ denotes the overall loss function which consists of the performance loss \mathcal{L}_{track} , the latency function $\mathcal{L}_{latency}$ and the weight coefficient λ . To solve this multi-objective optimization problem, we implement a Pareto optimization procedure, consisting of two alternating stages.

Our goal for the first stage is to find the best architecture α^* of the backbone on the training dataset \mathcal{D}_{train} .

$$\alpha^* = \operatorname{argmin}_{\alpha} \mathbb{E}_{\mathcal{O}, \mathcal{T} \sim \mathcal{D}_{train}} [\mathcal{L}(\mathcal{O}, \mathcal{T} | \theta, \alpha)] \quad (3)$$

Note that we set $\lambda > 1$ since we mainly focus on reducing the latency in the first stage. We avoid solely using the latency loss function because α can influence both latency and performance.

For the second stage, we attempt to find the best model ζ^* with the highest accuracy using the searched architecture α^* from stage I. The optimization problem can be formulated as,

$$\min_{\theta} \mathbb{E}_{\mathcal{O}, \mathcal{T} \sim \mathcal{D}_{val}} [\mathcal{L}_{track}(\mathcal{O}, \mathcal{T} | \theta, \alpha^*)] \quad (4)$$

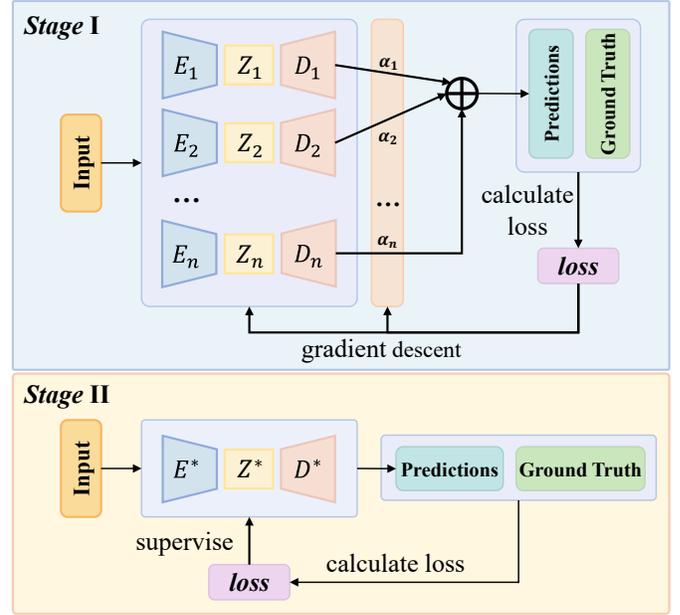


Fig. 3. The searching and training process is divided into two stages. We obtain the optimal network structures in the first stage through constrained NAS and get the optimal network parameters in the second training stage.

Algorithm 1 Stage 1: search for feature extraction backbone

Input: Object detections \mathcal{O} , Ground truth tracklets \mathcal{T}

Output: Optimal network structure encoding α^*

Initialize architecture $\alpha = \alpha_0$

for each epoch i **do**

if converged **then**

 Generate the optimal architecture encoding α^* ;

else

 Update the architecture encoding α

 Calculate the estimated latency from α

for each iteration t **do** training

 backpropagation with loss function

$\mathcal{L}_{track} + \lambda \mathcal{L}_{latency}$

 Update the network parameters θ

end for

 Evaluation on the validation dataset

end if

end for

After several iterations of the two-stage searching on the Pareto front of latency and accuracy, we can obtain the optimal parameters ζ^* .

$$\begin{aligned} \theta^* &= \operatorname{argmin}_{\theta} \mathbb{E}_{\mathcal{O}, \mathcal{T} \sim \mathcal{D}_{train}} [\mathcal{L}(\mathcal{O}, \mathcal{T} | \theta, \alpha^*)] \\ \zeta^* &= (\alpha^*, \theta^*) \end{aligned} \quad (5)$$

E. Two-stage Neural Architecture Search

We divide the training process into two stages. In the first stage named **train search**, we first update the architecture parameters $\alpha = \alpha_0$. Then for each epoch, we use the architecture from the last epoch to build our tracking network. We train the tracking network for T iterations, after which the model will be evaluated on a validation dataset to test

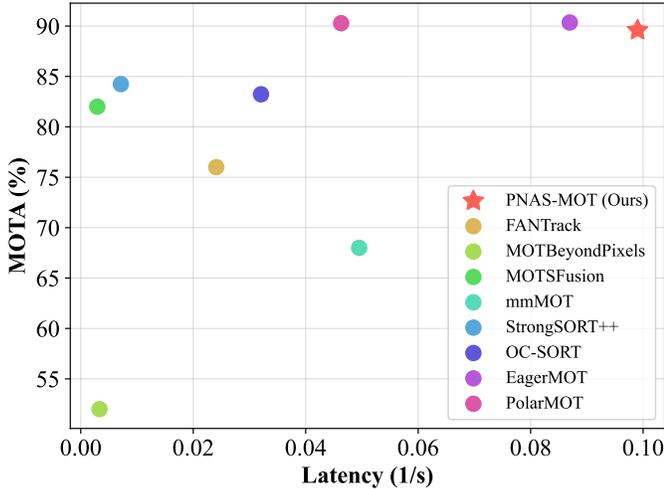


Fig. 4. Comparison with other methods on the KITTI testing dataset. The x-axis represents the reciprocal of the latency. As shown in the figure, our method minimally sacrifices the MOTA metrics and achieves outstanding performance on latency, with less than half the latency of mmMOT.

its performance. If the result does not converge, we continue to update the architecture parameters α for the next training epoch. In order to search for an efficient network structure with a low inference time, we add a latency constraint to the loss function for the first stage.

The second stage is the further training stage. In stage one, the operations between each node are mixed operations, each with a different weight ranging in $(0, 1)$, which causes the model will be extremely heavy. Therefore, we will remove the operations with low weight to reduce the size of the model. Pruning the low-weight operation will lead to a gap between the model from the first stage and the second stage. As a result, we continue to train the pruned model to get the best tracking network. In the second stage, the latency constraint is removed because the architecture will not be updated anymore. We focus on the improvement of accuracy in this stage.

Algorithm 2 Stage 2: Training with the searched architecture

Input: Object detections O , Ground truth tracklets \mathcal{T}

Output: Optimal network weights θ^*

Initialize architecture $\alpha = \alpha^*$

for each iteration t **do**

if not converged **then**

 Backpropagation using loss function \mathcal{L}_{track}

 Update the weights θ for the tracking network

else

 Save the model ζ^* with best performance

end if

 Evaluation on validation dataset at a fixed interval

end for

IV. EXPERIMENTS

A. Accuracy on KITTI Benchmark

Our model is trained and evaluated on the KITTI Benchmark [12], which consists of 29 testing sequences and 21

training sequences. Since KITTI does not provide an official validation dataset, we choose sequences 00, 02, 05, 07, 10, 11, 14, 16, 17, 18, and 19 as the validation dataset and the rest as the training dataset.

We submit the tracking results obtained from our model with the best performance on the validation dataset to the KITTI server, to evaluate our performance on the testing dataset. The comparison of runtime and accuracy can be observed in Figure 4. We compared our performance with other methods [1], [3], [6], [14], [15], [21], [29], [43]. Though our method may not come state of the art in accuracy (MOTA), our model can maintain relatively high accuracy with a short inference time. The quantitative results are demonstrated in Table I, and the qualitative results are illustrated in Figure 5.

B. Latency measured on different devices

We randomly chose sequence 19 from the KITTI dataset as a validation set for our evaluation, containing 1059 data, images, and point cloud files. We use the validation data to calculate the average inference time.

As Table II shows, we evaluated our algorithm on different GPUs, including Jetson Nano, Jetson Orin Nano, GTX, Quadro, and TITAN. We conduct experiments on both edge devices and effective GPUs. The results show that our algorithm can run on Jetson Nano, whose memory size is only 2GB, and the inference time is under 80ms. When running on other high-performance GPUs, the inference time can be reduced to less than 8ms at most.

During the search process, we made a trade-off between accuracy and latency. As can be observed from Table III, we can relax the latency constraint by reducing the parameter λ in the loss function of the first stage, resulting in a higher accuracy with higher latency.

C. Ablation studies

To evaluate the effectiveness of the latency constraint neural architecture search method. We conduct an ablation study on the KITTI dataset [12]. We compare the performance of three different search modes. For search mode 0, we only search for the image branch of the feature extraction backbone, and we use ResNet18 as the LiDAR branch. Similar to search mode 1, we use ResNet34 as the image branch and solely search for the LiDAR branch. Besides, we search for both image and LiDAR branches simultaneously under the setting of search mode 2. The experiments are conducted on Quadro RTX 6000. We choose sequence 0,2,5,7,10,11,14,16,18,19 of the KITTI dataset to evaluate the MOTA and sequence 19 to evaluate the latency of our searched network architectures.

As shown in Table IV, compared with search mode 0 and 1, at search mode 2, the latency can be reduced significantly, and the MOTA can be improved, especially the memory size during inference can be reduced by 22.3% and 14.8%, respectively. The model trained from search mode 2 can run on Jetson Nano, while the other two cannot since the Jetson Nano only has a 2GB memory size.

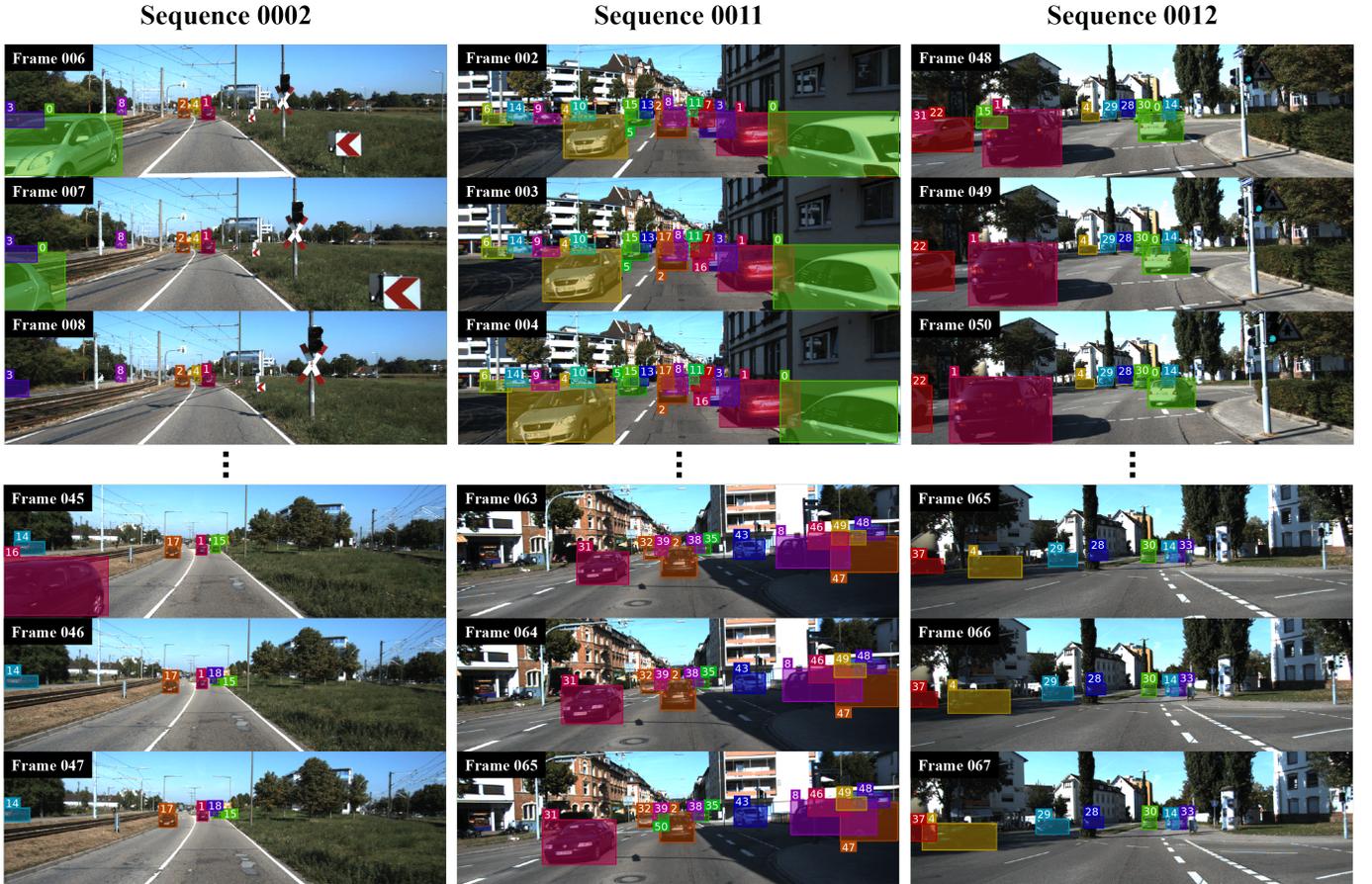


Fig. 5. Tracking results of our multi-object tracking model. We illustrate the tracking results of different sequences from the testing dataset of the KITTI benchmark. These are some of the extreme conditions under highly complex environmental challenges. Yet our model accurately keeps track of every object despite the interference of illumination, shadows, and occlusions. This suggests our multi-object tracking model has high accuracy and robustness concerning complex background environments.

TABLE I
COMPARISON ON THE TESTING DATASETS OF KITTI.

Method	MOTA \uparrow	MOTP \uparrow	IDS \downarrow	HOTA \uparrow	DetA \uparrow	AssA \uparrow	FN \downarrow	FP \downarrow	Frag \downarrow	MT \uparrow	ML \downarrow
DSM [10]	73.94	83.5	939	60.05	64.09	57.18	637	7388	737	59.38	8.46
extraCK [13]	79.29	82.06	520	59.76	65.18	55.47	675	5929	750	62.31	5.85
FANTrack [1]	75.84	82.46	743	60.85	64.36	58.69	1305	6262	701	62.77	8.77
MOTBeyondPixels [29]	82.68	85.50	934	63.75	72.87	56.4	741	4283	581	72.61	2.92
PMBM [27]	79.23	81.58	485	59.12	65.43	54.28	1024	5634	554	62.77	6.46
JCSTD [32]	80.24	81.85	173	65.94	65.37	67.03	405	6217	700	57.08	7.85
mmMOT [43]	83.23	85.03	733	62.05	72.29	54.02	752	4284	570	72.92	2.92
Ours	89.59	85.44	751	67.32	77.69	58.99	568	2261	276	86.59	2.46

TABLE II
LATENCY ON DIFFERENT DEVICES. THE PERFORMANCES OF OUR ALGORITHM ON DIFFERENT DEVICES STRICTLY SATISFY LATENCY REQUIREMENTS AND ARE FAR BETTER THAN THE MINIMUM LIMIT.

Devices	Latency (milliseconds)	Floating-point performance (GFLOPS)	compute capability	Memory Size (GB)
Jetson Nano	78	235.8	5.3	2
Jetson Orin Nano	58	1280	8.7	8
GTX 1050 Ti	18	2138	6.1	4
GTX 2080 Ti	10	13450	7.5	11
Quadro RTX 6000	8	16310	7.5	24
Quadro RTX 8000	8	16310	7.5	48
TITAN RTX	9	16310	7.5	24
TITAN V	8	14900	7.0	12

TABLE III
MOTA-LATENCY TRADE-OFF ON THE VALIDATION DATASET

λ	Latency (ms)	MOTA (%)	HOTA (%)	IDSW
10.0	6.2	89.42	71.02	334
1.0	8.3	90.48	73.85	235
0.1	10.1	90.91	75.84	188
0.01	21.0	91.01	76.82	183

TABLE IV
COMPARISON BETWEEN DIFFERENT SEARCH MODE

Search Mode	Latency (ms)	MOTA %	Memory Size (MB)
0	12	90.3	2462
1	10	90.6	2246
2	8	90.9	1913

D. Implementation Details

We test the latency of our model on different devices. The best model is trained on an RTX 8000, and the batch size we use is 1. We use ADAM as our optimizer with a learning rate of $3e^{-6}$. The search space is adapted from DARTs [18] while modified to fit the MOT tasks. The candidate operations are as follows: none, identity, 3×3 , 5×5 , 7×7 separable convolutions, 3×3 , 5×5 dilated convolution, 3×3 max pooling, and 3×3 average pooling.

In our designed architecture, there are two branches in the feature extraction backbone, one is image modality, and the other one is LiDAR modality. For each branch, the network consists of two different kinds of cell, normal cell and reduction cell, where the former one does not change the feature channels while the latter one does. Each cell consists of N nodes, where each edge between two nodes represents an operation in a pre-defined search space. Our goal is to find the best architecture of normal and reduction cells, including the connection between each node and the operation of each edge. Notably, the searched architecture parameters α are shared by cells of both branches.

V. LIMITATION

We employ DARTs [18] as the backbone for Neural Architecture Search, where the network structure is parameterized by weights assigned to edges between nodes. The searched structure is obtained by discarding edges with weights falling below a certain threshold. Consequently, the structure of Stage II deviates slightly from that of Stage I. This discrepancy introduces a minor deviation in the searched structure from the optimal configuration.

VI. CONCLUSION

In this paper, we introduce a latency-constrained multiple modalities fusion neural architecture search method for MOT tasks. Numerical experiments have demonstrated the superiority of our proposed scheme. We have achieved 89.59% accuracy close to the SOTA methods while keeping the latency below 80 milliseconds. This methodology may serve as a foothold for future efficient autonomous driving research.

REFERENCES

- [1] Erkan Baser, Venkateshwaran Balasubramanian, Prarthana Bhat-tacharyya, and Krzysztof Czarnecki. Fantrack: 3d multi-object tracking with feature association network. In *2019 IEEE Intelligent Vehicles Symposium (IV)*, pages 1426–1433, 2019. 5, 6
- [2] Philipp Bergmann, Tim Meinhardt, and Laura Leal-Taixé. Tracking without bells and whistles. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 941–951, 2019. 1
- [3] Jinkun Cao, Jiangmiao Pang, Xinshuo Weng, Rawal Khirodkar, and Kris Kitani. Observation-centric sort: Rethinking sort for robust multi-object tracking. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9686–9696, 2023. 5
- [4] Yubo Cui, Jiayao Shan, Zuoxu Gu, Zhiheng Li, and Zheng Fang. Exploiting more information in sparse point cloud for 3d single object tracking. *IEEE Robotics and Automation Letters*, 7(4):11926–11933, 2022. 2
- [5] Stefan de Jong, Hilmy Baja, Karsjen Tamminga, and João Valente. Apple mots: Detection, segmentation and tracking of homogeneous objects using mots. *IEEE Robotics and Automation Letters*, 7(4):11418–11425, 2022. 2
- [6] Yunhao Du, Zhicheng Zhao, Yang Song, Yanyun Zhao, Fei Su, Tao Gong, and Hongying Meng. Strongsort: Make deepsort great again. *IEEE Transactions on Multimedia*, 2023. 5
- [7] Matteo Dunnhofer, Niki Martinel, and Christian Micheloni. Weakly-supervised domain adaptation of deep regression trackers via reinforced knowledge distillation. *IEEE Robotics and Automation Letters*, 6(3):5016–5023, 2021. 2
- [8] Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. Neural architecture search: A survey. *Journal of Machine Learning Research*, 20(1):1997–2017, 2019. 1
- [9] Christoph Feichtenhofer, Axel Pinz, and Andrew Zisserman. Detect to track and track to detect. In *2017 IEEE/CVF International Conference on Computer Vision*, pages 3038–3046, 2017. 1
- [10] Davi Frossard and Raquel Urtasun. End-to-end learning of multi-sensor 3d tracking by detection. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 635–642, 2018. 1, 3, 6
- [11] Changhong Fu, Teng Li, Junjie Ye, Guangze Zheng, Sihang Li, and Peng Lu. Scale-aware domain adaptation for robust uav tracking. *IEEE Robotics and Automation Letters*, 2023. 2
- [12] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3354–3361. IEEE, 2012. 5
- [13] Gültekin Gündüz and Tankut Acarman. A lightweight online multiple object vehicle tracking method. In *2018 IEEE Intelligent Vehicles Symposium (IV)*, pages 427–432, 2018. 6
- [14] Aleksandr Kim, Guillem Brasó, Aljoša Ošep, and Laura Leal-Taixé. Polarmot: How far can geometric relations take us in 3d multi-object tracking? In *2022 European Conference on Computer Vision (ECCV)*, pages 41–58. Springer, 2022. 5
- [15] Aleksandr Kim, Aljoša Ošep, and Laura Leal-Taixé. Eagermot: 3d multi-object tracking via sensor fusion. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11315–11321, 2021. 1, 2, 5
- [16] Philip Lenz, Andreas Geiger, and Raquel Urtasun. Followme: Efficient online min-cost flow tracking with bounded memory and computation. In *2015 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4364–4372, 2015. 3
- [17] Jian Li, Yisen Huang, Xue Zhang, Ke Xie, Yitian Xian, Xiao Luo, Philip Wai Yan Chiu, and Zheng Li. An autonomous surgical instrument tracking framework with a binocular camera for a robotic flexible laparoscope. *IEEE Robotics and Automation Letters*, 2023. 2
- [18] Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: Differentiable architecture search. In *International Conference on Learning Representations (ICLR)*, 2019. 1, 2, 7
- [19] Huaping Liu, Yupei Wu, Fuchun Sun, Bin Fang, and Di Guo. Weakly paired multimodal fusion for object recognition. *IEEE Transactions on Automation Science and Engineering*, 15(2):784–795, 2017. 2
- [20] Jiuming Liu, Guangming Wang, Zhe Liu, Chaokang Jiang, Marc Pollefeys, and Hesheng Wang. Regformer: an efficient projection-aware transformer network for large-scale point cloud registration. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 8451–8460, 2023. 1
- [21] Jonathon Luiten, Tobias Fischer, and Bastian Leibe. Track to reconstruct and reconstruct to track. *IEEE Robotics and Automation Letters*, 5(2):1803–1810, 2020. 5

- [22] Niv Nayman, Yonathan Aflalo, Asaf Noy, and Lihi Zelnik. Hardcorenas: Hard constrained differentiable neural architecture search. In *International Conference on Machine Learning (ICML)*, pages 7979–7990. PMLR, 2021. 2
- [23] Chensheng Peng, Guangming Wang, Xian Wan Lo, Xinrui Wu, Chenfeng Xu, Masayoshi Tomizuka, Wei Zhan, and Hesheng Wang. Delflow: Dense efficient learning of scene flow for large-scale point clouds. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 16901–16910, 2023. 2
- [24] Jinlong Peng, Tao Wang, Weiyao Lin, Jian Wang, John See, Shilei Wen, and Erui Ding. Tpm: Multiple object tracking with tracklet-plane matching. *Pattern Recognition*, 107:107480, 2020. 2
- [25] Yufu Qu, Guirong Zhang, Zhaofan Zou, Ziyue Liu, and Jiansen Mao. Active multimodal sensor system for target recognition and tracking. *Sensors*, 17(7):1518, 2017. 2
- [26] Jimmy Ren, Xiaohao Chen, Jianbo Liu, Wenxiu Sun, Jiahao Pang, Qiong Yan, Yu-Wing Tai, and Li Xu. Accurate single stage detector using recurrent rolling convolution. In *2017 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5420–5428, 2017. 2
- [27] Samuel Scheidegger, Joachim Benjaminsson, Emil Rosenberg, Amrit Krishnan, and Karl Granström. Mono-camera 3d multi-object tracking using deep learning detections and pmbm filtering. In *2018 IEEE Intelligent Vehicles Symposium (IV)*, pages 433–440, 2018. 6
- [28] Samuel Schuster, Paul Vernaza, Wongun Choi, and Manmohan Chandraker. Deep network flow for multi-object tracking. In *2017 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6951–6960, 2017. 3
- [29] Sarthak Sharma, Junaid Ahmed Ansari, J Krishna Murthy, and K Madhava Krishna. Beyond pixels: Leveraging geometry and shape cues for online multi-object tracking. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3508–3515, 2018. 1, 5, 6
- [30] ShiJie Sun, Naveed Akhtar, HuanSheng Song, Ajmal Mian, and Mubarak Shah. Deep affinity network for multiple object tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(1):104–119, 2019. 2
- [31] Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V Le. Mnasnet: Platform-aware neural architecture search for mobile. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2820–2828, 2019. 2
- [32] Wei Tian, Martin Lauer, and Long Chen. Online multi-object tracking using joint domain information in traffic scenarios. *IEEE Transactions on Intelligent Transportation Systems*, 21(1):374–384, 2019. 6
- [33] Ba Tuong Vo, Chong Meng See, Ning Ma, and Wee Teck Ng. Multi-sensor joint detection and tracking with the bernoulli filter. *IEEE Transactions on Aerospace and Electronic Systems*, 48(2):1385–1402, 2012. 1
- [34] Paul Voigtlaender, Michael Krause, Aljosa Osep, Jonathon Luiten, Berin Balachandar Gnana Sekar, Andreas Geiger, and Bastian Leibe. MOTs: Multi-object tracking and segmentation. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7942–7951, 2019. 1
- [35] Guangming Wang, Chensheng Peng, Yingying Gu, Jinpeng Zhang, and Hesheng Wang. Interactive multi-scale fusion of 2d and 3d features for multi-object vehicle tracking. *IEEE Transactions on Intelligent Transportation Systems*, 2023. 1, 2
- [36] Pan Wang, Liangliang Ren, Shengkai Wu, Jinrong Yang, En Yu, Hangcheng Yu, and Xiaoping Li. Implicit and efficient point cloud completion for 3d single object tracking. *IEEE Robotics and Automation Letters*, 8(4):1935–1942, 2023. 2
- [37] Xinshuo Weng, Jianren Wang, David Held, and Kris Kitani. 3d multi-object tracking: A baseline and new evaluation metrics. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 10359–10366. IEEE, 2020. 1
- [38] Xinshuo Weng, Yongxin Wang, Yunze Man, and Kris M Kitani. Gnn3dmot: Graph neural network for 3d multi-object tracking with 2d-3d multi-feature learning. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6499–6508, 2020. 1
- [39] Hai Wu, Wenkai Han, Chenglu Wen, Xin Li, and Cheng Wang. 3d multi-object tracking in point clouds based on prediction confidence-guided data association. *IEEE Transactions on Intelligent Transportation Systems*, 23(6):5668–5677, 2021. 2
- [40] Yichen Xie, Chenfeng Xu, Marie-Julie Rakotosaona, Patrick Rim, Federico Tombari, Kurt Keutzer, Masayoshi Tomizuka, and Wei Zhan. Sparsefusion: Fusing multi-modal sparse representations for multi-sensor 3d object detection. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 17591–17602, 2023. 2
- [41] Shaoqing Xu, Dingfu Zhou, Jin Fang, Junbo Yin, Zhou Bin, and Liangjun Zhang. Fusionpainting: Multimodal fusion with adaptive attention for 3d object detection. In *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, pages 3047–3054. IEEE, 2021. 2
- [42] Tianwei Yin, Xingyi Zhou, and Philipp Krahenbuhl. Center-based 3d object detection and tracking. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11784–11793, 2021. 1
- [43] Wenwei Zhang, Hui Zhou, Shuyang Sun, Zhe Wang, Jianping Shi, and Chen Change Loy. Robust multi-modality multi-object tracking. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2365–2374, 2019. 1, 3, 5, 6
- [44] Yue Zhang, Bin Song, Xiaojiang Du, and Mohsen Guizani. Vehicle tracking using surveillance with multimodal data fusion. *IEEE Transactions on Intelligent Transportation Systems*, 19(7):2353–2361, 2018. 2
- [45] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2017. 2