# MO-YOLO: End-to-End Multiple-Object Tracking Method with YOLO and Decoder

Liao Pan[1,2], Yang Feng[1,2], Wu Di[1,2], Liu Bo[1,2], Zhang Xingle[1,2]

[1] Northwestern Polytechnical University
[2] School of Automation

**Abstract.** The recent emergence of Transformer-based end-to-end models, exemplified by MOTR, has showcased remarkable multi-object tracking (MOT) performance on datasets like DanceTracker. However, the computational complexities associated with these models pose challenges in both training and deployment. Drawing insights from successful models such as GPT, our proposed MO-YOLO stands out as an efficient and computationally frugal end-to-end MOT solution. MO-YOLO integrates principles from YOLO and RT-DETR, adopting a decoder-centric architecture alongside other complementary structures. By leveraging the RT-DETR decoder and architectural components from YOLOv8, MO-YOLO achieves high-speed performance, shorter training times, and proficient MOT capabilities. On the Dancetrack dataset, MO-YOLO not only achieves parity with MOTR in tracking performance but also surpasses it in terms of speed (MOTR 9.5 FPS, MO-YOLO 19.6 FPS). Moreover, MO-YOLO demonstrates significantly reduced training times and lower hardware requirements compared to MOTR. This research introduces a promising paradigm for efficient end-to-end MOT, highlighting enhanced performance and resource efficiency in a succinct and impactful manner.

**Keywords:** Multiple-Object Tracking, Transformer ,YOLO, End-to-End

## 1 Introduction

Multi-Object Tracking (MOT) is a crucial task in computer vision, aiming to predict object trajectories across consecutive images [2]. The tracking-by-detection approach, relying on object detectors for individual target identification [1,9,22], has gained prominence. This method involves associating detections using algorithms like the Kalman filter [3] to form object trajectories. While demonstrating remarkable performance, the intrinsic two-step process—commencing with initial object detection and followed by object tracking—introduces heightened computational complexity and the potential for error propagation.

In recent times, the field has witnessed the emergence of end-to-end MOT models based on Transformer architectures [38,40,42], boasting superior performance on datasets such as DanceTracker [29] compared to tracking-by-detection methods. However, the computational demands associated with Transformer

models present challenges in terms of both training and deployment, necessitating formidable hardware resources.

In order to address these challenges more effectively, a retrospective analysis of the origins of the Transformer model becomes imperative. The Transformer architecture was initially proposed for the domain of Natural Language Processing (NLP) [31]. As researchers continuously refined this model, the GPT [4] emerged as a prominent highlight in the NLP field over the past few years. GPT, exclusively utilizing the decoder component of the Transformer, achieves pre-training and fine-tuning for various language tasks by stacking multiple decoder layers to construct a deep neural network.

Inspired by the success of GPT, this paper adopts a network structure combining the principles of You Only Look Once (YOLO) [25] and RT-DETR [21], incorporating tracking concepts from MOTR [42]. The resultant model, named MO-YOLO, follows a decoder-based end-to-end MOT approach.

MO-YOLO leverages the decoder from RT-DETR and the block and neck components from YOLOv8 [14], offering advantages such as high speed and shorter training times. This architecture demonstrates proficient performance in multitarget tracking tasks.
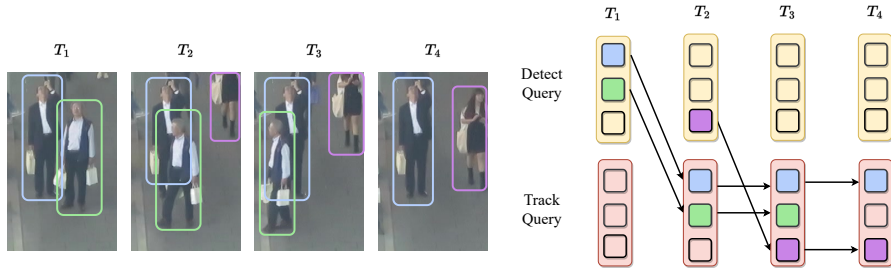
This paper pioneers a fusion of concepts from YOLO and the end-to-end MOTR model, resulting in a comprehensive end-to-end MOT model. The contributions of this work can be encapsulated in four fundamental advancements:

1) The primary contribution involves the development of a novel end-to-end tracking network. This network amalgamates structural components derived from YOLO, RT-DETR, and the conceptual framework of MOTR, establishing a cohesive and efficient model for Multiple Object Tracking.
2) A second key innovation introduces a unique training strategy characterized by a three-stage process. This strategy accelerates model convergence, reducing overall training time and enhancing training process efficiency.
3) The third contribution introduces the Tracking Box Selection Process (TBSP) strategy, strategically filtering bounding boxes during training to expedite model convergence and enhance training process efficiency.
4) Lastly, we substantiate the feasibility of a decoder based model akin to GPT in the domain of MOT in computer vision.

## 2   Related work

### 2.1   MOTR Series

MOTR [42] is a groundbreaking model that achieves end-to-end multi-object tracking. It extends the concept of object queries from DETR [7] to track queries, allowing for dynamic tracking in video sequences. By employing a learnable positional embedding that interacts with the decoder and feature maps through multi-scale deformable attention, MOTR predicts target positions and classes in each frame. Notably, it establishes implicit temporal associations without the need for explicit data association or post-processing, seamlessly integrating

**Fig. 1:** MOTR's Streamline the process of updating detect queries and track queries under typical MOT scenarios. The track query set, initially empty, is dynamically adjusted in length. Detect queries play a role in identifying new objects.

detection and tracking to improve accuracy and robustness. The introduction of Track Query Aggregation with Learned Attention (TALA) is a key innovation, which models long-term temporal relations using a query memory bank and multi-head attention. Additionally, MOTR introduces innovative concepts such as Collective Average Loss (CAL) and Temporal Aggregation Network (TAN) to address target conflicts and enhance the integration of temporal information. These design elements collectively position MOTR as a significant milestone in multi-object tracking, elevating both its performance and applicability.

MOTRv2 [47] enhances MOTR's detection performance by incorporating the pre-trained object detector YOLOX [11] for detection priors. This involves anchoring object queries and utilizing YOLOX-generated proposal boxes as anchors to initialize the object queries. This strategy not only improves detection but also mitigates conflicts between detection and association tasks while preserving query propagation characteristics.

MeMOTR [10] introduces track queries, modeling target trajectories throughout the video and predicting them iteratively in each frame. It implements a tracklet-aware label assignment method to achieve one-to-one assignment between track queries and target trajectories, considering newborn object queries.

MOTRv3 [40] and CO-MOT [37] share a common focus on addressing unfair label assignment between detection and tracking queries in MOTR, aiming to balance training samples and improve end-to-end multi-object tracking performance. MOTRv3 employs a release-fetch supervision strategy for balanced label assignment, while CO-MOT uses a coopetition label assignment strategy to increase positive samples for detection queries and introduces "shadows" to enrich query quantity and diversity.

## 2.2 YOLO

YOLO [25] has transformed object detection into a regression model, achieving efficient target detection. It partitions the image into grids, with each grid responsible for predicting the target box and class information within its corresponding region. The object detection task is treated as a regression problem

for coordinates and dimensions. This unique approach simplifies the detection process, enabling YOLO to simultaneously process multiple targets in a single forward pass, thus achieving outstanding speed and efficiency. The concept of this regression model provides a solid foundation for real-time object detection.

Due to its relatively high efficiency and accuracy, the SORT [3] initially applied YOLO in the field of MOT. However, with the proliferation of such models, YOLO has become widely used in the MOT domain. Nevertheless, this approach still introduces additional processing layers, potentially leading to some errors. Therefore, in this paper, we construct an end-to-end tracking model based on the YOLO framework.

### 2.3  RT-DETR

DETR [50] revolutionizes object detection by introducing the Transformer architecture, transforming the task into a multi-class classification problem. It excels in precise object localization and modeling object relationships, making it particularly suitable for applications like multi-object tracking. Despite its effectiveness with limited annotated data, DETR comes with the drawback of high computational complexity, demanding significant computational resources.

RT-DETR [21] addresses DETR's computational challenges, presenting the first real-time usable DETR model. It tackles complexity with an efficient hybrid encoder that manages multi-scale features through smart intra-scale interactions and cross-scale fusion decoupling. This significantly reduces computational costs, achieving real-time object detection.
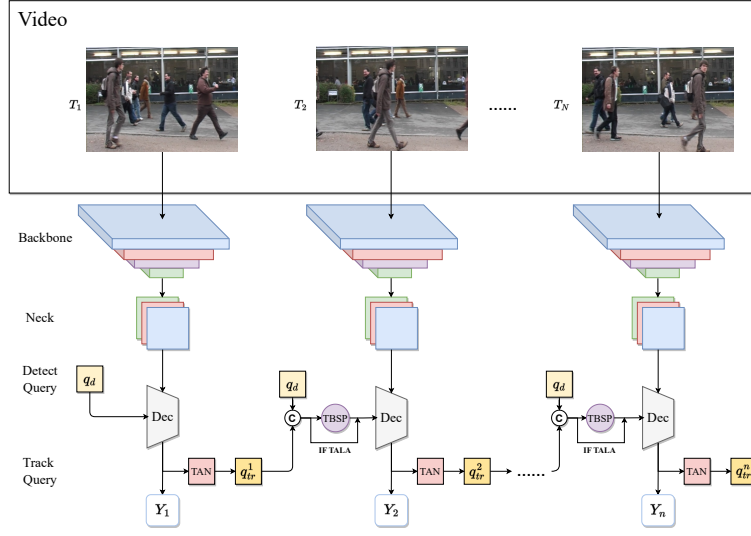
RT-DETR's decoder module incorporates deformable attention, aligning with attention mechanisms in various DETR models [36,50,51]. By dynamically generating attention sampling points, it adapts to diverse feature distributions and target shapes, enhancing adaptability. Unlike some DETR models [8,36,43], RT-DETR skips the use of a reconstruction network, relying on the decoder's output for final predictions. This simplifies the architecture while maintaining efficiency and performance.

## 3   MO-YOLO

### 3.1   MO-YOLO Architecture

MO-YOLO's overall architecture is illustrated in the provided Fig.2. In this setup, a video sequence is input to YOLOv8's backbone, neck to extract features. For the first frame, where there is no tracking information available, a fixed-length learnable detect query (referred to as $q_d$ in the figure) is input to the RT-DETR decoder. The learning pattern of $q_d$ is similar to RT-DETR [21].

For subsequent frames in the video sequence, the tracking queries from the previous frame and the learnable detect queries are concatenated and undergo filtering before being input to the decoder. These queries interact with the image features in the decoder to generate hidden states for bounding box predictions. Additionally, a feature aggregation process using MOTR's TAN is performed to generate trajectory queries for the next frame.
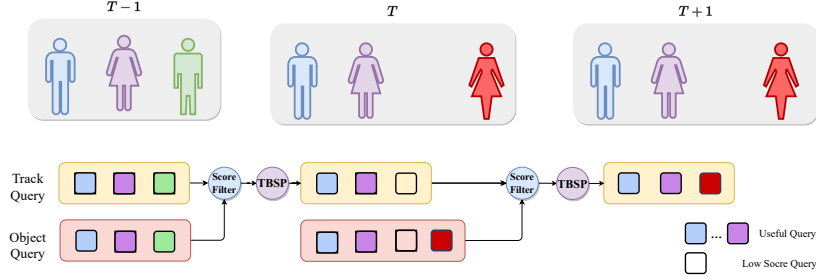
**Fig. 2:** The overall structure of MO-YOLO consists of an architecture that combines YOLOv8's backbone and neck, facilitating the extraction of image features for each frame. By merging the detect queries ($q_d$) and track queries ($q_{tr}$), the decoder (Dec) processes them to generate hidden states. When the model is trained using the TALA strategy, the TBSP can be skipped.

## 3.2 Tracking Box Selection Process

When training a detection model using the detection method and subsequently employing an end-to-end tracking approach, we observed a gradual increase in the number of tracking boxes for the same object over time. Upon examination, this phenomenon arises from the fact that each newly generated detect query is consistently assigned to an already tracked object. This is due to the nature of DETR, where a detect query can potentially be assigned to any object in the image, as label assignment is determined through binary matching of all detect queries and ground truth.
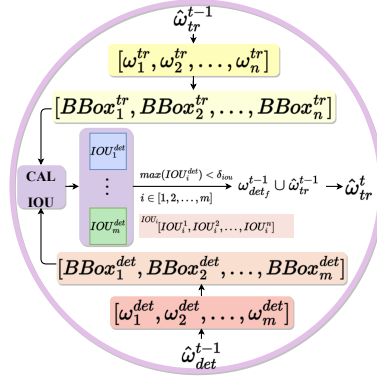
In MOTR [42], a solution to this issue is introduced through the utilization of TALA. This strategy gradually suppresses the generation of detect queries similar to track queries during the training process, effectively addressing the identified problem. However, this process increases the number of parameters that the network needs to automatically adjust, resulting in a prolonged convergence time for the model. Therefore, this paper proposes a simple yet effective approach named tracking box selection process (TBSP). Experimental results demonstrate the efficacy of TBSP, revealing that this method allows the model to achieve comparable performance to fully trained models even without extensive training. Additionally, the simplicity of TBSP ensures that it introduces minimal additional computational burden.

**Fig. 3:** Refinement process for MO-YOLO on partially trained track queries. When MO-YOLO undergoes comprehensive training with the TALA strategy, TBSP will no longer be utilized. At that point, the track query generation method of MO-YOLO aligns with MOTR, as illustrated in the Fig. 1.

We know that both track query and detect query contain information about the object's bounding box(BBox). Let's denote the BBox in $\hat{\omega}_{tr}^{t-1}$ as $[BBox_1^{tr}, BBox_2^{tr}, ..., BBox_n^{tr}]$, where $n$ is the number of track queries at time $t-1$. Each query in $\hat{\omega}_{tr}^{t-1}$ is denoted as $[\omega_1^{det}, \omega_2^{det}, ..., \omega_m^{det}]$, and its corresponding BBox is denoted as $[BBox_1^{det}, BBox_2^{det}, ..., BBox_m^{det}]$.

For each element in $[BBox_1^{det}, BBox_2^{det}, ..., BBox_m^{det}]$, calculate the Intersection over Union (IOU) with each element in $[BBox_1^{tr}, BBox_2^{tr}, ..., BBox_n^{tr}]$. If the IOU is greater than $\delta_{iou}$, the BBox corresponding to that detect Query will be discarded. Finally, the remaining detect queries are denoted as $\omega_{det_f}^{t-1}$.



**Fig. 4:** The specific process of TBSP is represented by a formula. The specific meaning of the formulas in the figure refers to Section 3.2.

Similarly, assuming the track query at frame $t$ (where $t > 2$) is denoted as $\omega_{tr}^t$, it can be expressed as:

$$\hat{\omega}_{tr}^t = \omega_{det_f}^{t-1} \cup \hat{\omega}_{tr}^{t-1} \tag{1}$$

The feasibility of TBSP is illustrated in the Fig. 4. Incidentally, although TBSP makes MO-YOLO appear not to be a truly end-to-end network, after thorough TALA training, this step can be omitted. Simultaneously, to better illustrate the performance of our model compared to MOTR in CO-MOT [38] and MOTRv3 [40], we have not adopted the labeling strategy used in CO-MOT and MOTRv3.

### 3.3   Decoder Embeddings and Query Positions

Differentiating itself from MOTR, MO-YOLO's decoder employs a methodology akin to RT-DETR for its embeddings. In contrast to MOTR, MO-YOLO's decoder input doesn't originate from a distinct encoder; rather, it exploits RT-DETR's IoU-aware query selection mechanism. Building upon prior works such as [16, 35, 43, 50], this process strategically selects high-quality image features as the initial object embeddings for the decoder, emphasizing IoU scores.

It's essential to clarify that in this context, the term "embeddings" specifically pertains to the initial embeddings of detection queries. The subsequent embeddings of track queries undergo generation through a TAN network, aligning with MOTR's established process.

To expedite model convergence, we implement a Query Positions generation strategy akin to CO-MOT. Drawing inspiration from DAB-DETR [19], we employ a fixed transformation to convert reference points into Query Positions. While experimentation with the RT-DETR approach involving a Multilayer Perceptron (MLP) to transform reference points shows no significant performance discrepancy, it is important to note that the use of MLP substantially increases training time. By the way, it is worth mentioning that this process differs significantly from MOTR. In MOTR, Query Positions are generated using a set of learnable parameters, whereas in MO-YOLO, they are generated using reference points.

### 3.4   Loss and Training Methods

Due to the unique training mode adopted by MOTR, its batch size in a single GPU is constrained to 1. While this training approach enables the model to effectively learn the temporal information of the targets, it concurrently requires training the network to extract features of target appearances. To expedite the model training speed, a three-stage training strategy is employed in this study. In the first stage, the emphasis is placed on training the network to extract features of target appearances; the second stage focuses on learning temporal information; and the final stage comprehensively enhances the model's performance by integrating the learning of temporal information and appearance features. This strategy is implemented to overcome the constraints posed by the unique training mode of MOTR, aiming to achieve more efficient model training.

**In the initial stage**, the network is trained as a detection network, and the loss function is formulated defined as Eq.2:

$$\mathcal{L}_o\left(\left.\widehat{Y}\right|_\omega, Y\right) = \frac{\mathcal{L}\left(\left.\widehat{Y}_{\det}^t\right|_{\omega_{\det}^t}, Y_{\det}^t\right)}{V_i} \tag{2}$$

$\mathcal{L}$ is the loss of single frame, which is similar to the detection loss in DETR, and it is shown in Eq.3.

$$\mathcal{L}\left(\left.\widehat{Y}_t\right|_{\omega_t}, Y_t\right) = \lambda_{cls}\mathcal{L}_{cls} + \lambda_{l_1}\mathcal{L}_{l_1} + \lambda_{\mathrm{giou}}\mathcal{L}_{\mathrm{giou}} \tag{3}$$

where $\hat{\omega}_{tr}^{t-1}$ represents the multi-frame predictions, $\hat{Y} = \left\{\hat{Y}_i\right\}_{i=1}^N$ represents the loss calculated over the entire video sequence, using ground truth $Y = \{Y_i\}_{i=1}^N$ and matched results $\omega = \{\omega_i\}_{i=1}^N$. $\mathcal{L}_{cls}$ denotes the focal loss [17], $\mathcal{L}_{cls}$ represents the L1 loss and $\mathcal{L}_{giou}$ represents the giou loss [26], $\lambda_{cls}, \lambda_{l1}, \lambda_{giou}$ are corresponding weight coefficients.

In this stage, the network is capable of initially learning the appearance features of objects to be tracked. During this phase, the batch size can be set larger, fully utilizing the computational power of GPUs to reduce training time. Additionally, various data augmentation techniques [13, 32, 44] can be employed to expedite the model's extraction of appearance features.

**The second stage** of training is conducted using MOTR's Collective Average Loss. During this phase, only MOTR's TAN and our TBSP are utilized. The loss functions used in this stage and the third stage are as Eq.4. Simultaneously, it is important to note that, at this juncture, the determination of the newly generated track queries relies solely on the scores of object queries and whether they have been filtered out by TBSP, as more advanced strategies like TALA have not been employed.

$$\mathcal{L}_o\left(\left.\widehat{Y}\right|_\omega, Y\right) = \frac{\sum_{n=1}^N\left(\mathcal{L}\left(\left.\widehat{Y}_{tr}^t\right|_{\omega_{tr}^t}, Y_{tr}^t\right) + \mathcal{L}\left(\left.\widehat{Y}_{\det}^t\right|_{\omega_{\det}^t}, Y_{\det}^t\right)\right)}{\sum_{n=1}^N (V_t)}, N = 5 \tag{4}$$

where $V_i = V_{tr}^i + V_{det}^i$ denotes the total number of ground-truths objects at frame $i$. $V_{tr}^i$ and $V_{det}^i$ are the numbers of tracked objects and newborn objects at frame i, respectively.

**The third stage** of training involves the utilization of TALA policy during the training phase, while simultaneously deactivating the TBSP. Experimental evidence indicates that if the training epochs during this stage are sufficient, TBSP can also be disabled during the inference phase without significantly compromising performance. This finding underscores the robustness of the model, suggesting that once adequately trained in the third stage, the TBSP policy can be turned off during inference with negligible impact on overall performance.

In addition, it is important to note that while the first and second stages of MO-YOLO share similarities in generating a new track query from the object query, discrepancies may arise in the third stage. Tracking objects with the track query does not involve additional techniques, and the process across the three stages of MO-YOLO closely resembles that of MOTR.

## 4 Experiments

### 4.1 Datasets and Metrics

This paper introduces MO-YOLO, an innovative MOT method proficient in addressing intricate scenarios involving occlusions, interactions, and diverse motions. To validate its performance, we evaluate the algorithm on three demanding datasets: DanceTrack [29], MOT17 [24], and KITTI [12].

DanceTrack is a vast dataset designed for human tracking, showcasing scenarios with occlusion, frequent crossovers, uniform appearances, and varied body gestures. Comprising 100 videos featuring diverse dance styles, it underscores the significance of motion analysis in multi-object tracking.

MOT17 focuses on multiple object tracking in public spaces, primarily pedestrians. The dataset includes seven scenes, divided into training and testing clips, providing object detection annotations suitable for both online and offline tracking approaches.

KITTI serves as a research dataset for autonomous driving, encompassing stereo, optical flow, odometry, 3D object detection, and tracking tasks. Captured across diverse environments using advanced sensors, it provides raw data and task-specific benchmarks.

**Evaluation Metrics**: We assessed our method using widely recognized MOT evaluation metrics. The primary metrics employed for evaluation included HOTA, AssA, DetA ,IDF1 and MOTA [20,27]. These metrics collectively provided a comprehensive and robust evaluation of our model's tracking performance. HOTA assessed the tracking accuracy based on the spatial and temporal overlap between predicted and ground truth tracks. AssA measured the quality of associations between objects, while IDF1 quantified the accuracy of identity tracking. MOTA offered an overall evaluation of tracking accuracy, considering false positives, false negatives, and identity switches. These metrics formed the basis for a thorough evaluation of our model's tracking capabilities.

### 4.2 Implementation Details

The rationale for adopting MO-YOLO originates from our observation of the high training costs and suboptimal inference speed associated with MOTR. Thus, our aim is to propose a new baseline model surpassing MOTR's performance. To ensure a fair comparison, we exclusively employ methods enhancing tracking performance, while omitting improvement techniques utilized in MeM-OTR [10], MOTRv2 [47], MOTRv3 [40], and CO-MOT [38]. To achieve a comparable parameter count with MOTR, we simultaneously utilize YOLOv8-L [14]

and RT-DETR-L [21] in constructing MO-YOLO (henceforth referred to as MO-YOLO-L concerning model size). The initial weights are obtained from a model pre-trained on the COCO dataset [18].

The majority of experiments were conducted using PyTorch on a single NVIDIA GeForce RTX 4090 GPU. For equitable comparison with MOTR, training on the MOT17 dataset was performed on a single NVIDIA GeForce RTX 2080ti. Additionally, to better contrast with the MOTR series in terms of speed, inference speed tests were carried out on a single NVIDIA V100 GPU. One more thing, he image size for the input network is $640 \times 640$.

During the initial training, batch sizes were set to 32 on the 4090 and 16 on the 2080ti. Subsequent stages utilized a uniform batch size of 1. **For MOT17**, like MOTR, we integrated extra data from the CrowdedHuman dataset [28], employing random displacement to generate video segments with pseudo-trajectories. Training comprised 120 epochs in the first stage, followed by 30 epochs in the second. The third stage extended to 20 epochs with TBSP, or 55 without. **For Dancetrack**, ensuring parity without extra data like MOTRv2, the first stage was 15 epochs, the second was 6, and the third, if TBSP was used, was 6 epochs, or 15 without TBSP. **For KITTI**, initial training involved 80 epochs, followed by 25 in the second stage. The third stage consisted of 30 epochs with TBSP or 70 without. Additionally, approximately 5k images from BDD100k [41] were selected for additional training.

**Note:** With the exception of Table 5b, all data concerning MO-YOLO in this paper pertains to MO-YOLO-L. Additionally, unless stated otherwise, the provided data pertains to the final end-to-end model.

### 4.3   Performance Evaluation and Comparison

In Table 1, a comprehensive comparison of training time and speed was conducted between MO-YOLO and the MOTR series.It is crucial to note that the

| **(a)** Comparison of Training Time on MOT17 | | | **(b)** Comparison of Speed on Dancetrack | |
|---|---|---|---|---|
| Method | Device | Training Time | Method | FPS |
| MOTR [42] | **8 2080ti** | **$\sim$ 96.0 hours** | MOTR [42] | 9.5 |
| MO-YOLO stage1 | 1 2080ti | 12.1 hours | MOTRv2 [47] | 6.9 |
| MO-YOLO stage2 | 1 2080ti | 30.8 hours | MOTRv3 [40] | 10.6 |
| MO-YOLO(TBSP) | 1 2080ti | 43.4 hours | MO-YOLO(TBSP) | 18.5 |
| MO-YOLO | **1 2080ti** | **68.7 hours** | MO-YOLO | **19.6** |

**Table 1:** Comparison of training time (including all previous stages) and speed between MO-YOLO and MOTR series on a single V100 GPU. The data for the MOTR series is sourced from [42], [40] or their open source address.

datasets used for comparing training time and inference speed exhibit variations.

This disparity primarily arises from the act that the publicly disclosed training times and inference speeds of the MOTR series are based on different datasets.

MO-YOLO demonstrates a substantial advantage over the MOTR series, particularly in terms of training time and speed. The primary reason behind this efficiency gain lies in the architectural simplicity of MO-YOLO. Unlike MOTR, MO-YOLO adopts a streamlined structure that eliminates the need for an encoder, resulting in a more computationally frugal design. This reduction in complexity enables MO-YOLO to achieve training completion within 68.7 hours, utilizing only 1 Nvidia GeForce 2080ti GPU. In contrast, MOTR relies on 8 Nvidia GeForce 2080ti GPUs and takes about 96 hours for training. The absence of an encoder in MO-YOLO contributes to its swift training process, emphasizing the model's efficiency in resource utilization.

On the Dancetrack test dataset, MO-YOLO(TBSP) achieves a speed of 18.5 frames per second (FPS) using a single V100, while the MOTR series attains speeds of 9.5 FPS (MOTR), 6.9 FPS (MOTRv2), and 10.6 FPS (MOTRv3). Particularly noteworthy is the remarkable 19.6 FPS achieved by MO-YOLO, emphasizing its outstanding performance in inference speed.

**Table 2:** Performance comparison between MO-YOLO and existing methods on the Dancetrack [29] dataset under the private detection protocols. Numbers are highlighted in bold when they correspond to either MO-YOLO or the superior metric from MOTR.

| Methods | HOTA | AssA | DetA | IDF1 | MOTA |
|---|---|---|---|---|---|
| *Non-End-to-end:* | | | | | |
| ByteTrack [45] | 47.3 | 31.3 | 71.6 | 52.5 | 89.5 |
| GTR [49] | 48.0 | 31.9 | 72.5 | 50.3 | 84.7 |
| DST-track [6] | 51.9 | 34.6 | 82.3 | 51.0 | 84.9 |
| OC-SORT [5] | 54.6 | 38.0 | 80.4 | 54.2 | 89.4 |
| *End-to-end:* | | | | | |
| MeMOTR [10] | 63.4 | 52.3 | 77.0 | 65.5 | 85.4 |
| **MOTR** [42] | **54.2** | **40.2** | 73.5 | **51.5** | 79.7 |
| MO-YOLO(Ours) | 52.0 | 35.7 | **75.9** | 51.4 | **83.5** |

Table 2, Table 3, and Table 4 exhibit the final results of end-to-end MO-YOLO on the Dancetrack, MOT17, and KITTI datasets, respectively. In contrast to MOTR, MO-YOLO demonstrates competitive performance across multiple metrics on both datasets. It is important to highlight that the performance of MeMOTR is provided in Tables 2 and 3. The inclusion of MeMOTR's performance aims to underscore the considerable potential for advancement and enhancement in both MOTR and our MO-YOLO. MeMOTR [10] integrates additional modules and performance-enhancing techniques, along with supplementary datasets during training. Consequently, comparing its tracking performance directly with that of MO-YOLO and MOTR may not be entirely fair.

**Table 3:** Performance comparison between MO-YOLO and existing methods on the MOT17 [24] dataset.

| method | HOTA | AssA | DetA | IDF1 | MOTA |
|---|---|---|---|---|---|
| *Non-End-to-end:* | | | | | |
| FairMOT [46] | 59.3 | 58.0 | 60.9 | 72.3 | 73.7 |
| ByteTrack [45] | 62.8 | 62.2 | 63.8 | 77.2 | 78.9 |
| StrongSORT [49] | 63.5 | 63.7 | 63.6 | 78.3 | 78.5 |
| UTM [39] | 64.0 | 62.5 | 65.9 | 78.7 | 81.8 |
| OC-SORT [5] | 63.2 | 63.2 | 63.2 | 77.5 | 78.0 |
| Deep OC-SORT [22] | 64.9 | 65.9 | 64.1 | 80.6 | 79.4 |
| *End-to-end:* | | | | | |
| MeMOTR [10] | 58.8 | 58.4 | 59.6 | 69.0 | 72.5 |
| **MOTR** [42] | **57.2** | 55.8 | **58.9** | 59.4 | **71.9** |
| MO-YOLO(Ours) | 55.5 | **55.9** | 55.8 | **70.1** | 66.8 |

For Dancetrack (Table 2), MO-YOLO achieves comparable results to MOTR, with slightly lower scores in HOTA and AssA but outperforming MOTR in DetA. Specifically, MO-YOLO achieves an impressive MOTA, indicating its overall superior performance in multi-object tracking tasks.

In the MOT17 (Table 3), MO-YOLO demonstrates competitive performance, particularly in handling complex tracking scenarios, with competitive scores in HOTA and AssA. Despite trailing slightly behind MOTR in these metrics, MO-YOLO compensates with a comparable DetA and excels in IDF1 by almost 10 points. This advantage can be attributed to MO-YOLO's simplified network structure, which enables competitive tracking accuracy with reduced training data demand. Additionally, MO-YOLO maintains robust overall performance, as evidenced by its high MOTA score, highlighting its suitability for real-world tracking applications.

**Table 4:** Performance comparison between MO-YOLO and existing methods on the KITTI [12] dataset.

| Tracker | Car | | | Pedestrian | | |
|---|---|---|---|---|---|---|
| | HOTA | MOTA | AssA | HOTA | MOTA | AssA |
| *Non-End-to-end:* | | | | | | |
| CenterTr [48] | 73.0 | 88.8 | 71.1 | 40.3 | 53.8 | 36.9 |
| PermaTr [30] | 77.4 | 90.8 | 77.6 | 47.4 | 65.0 | 43.6 |
| PolarMOT [15] | 75.1 | 80.0 | 76.9 | 43.6 | 46.9 | 48.1 |
| TripletTrack [23] | 73.6 | 84.3 | 74.7 | 50.0 | 46.5 | 57.8 |
| *End-to-end:* | | | | | | |
| MO-YOLO | 72.2 | 83.2 | 73.8 | 51.5 | 56.8 | 58.4 |

Table 4 provides an overview of MO-YOLO's performance on the KITTI dataset. MO-YOLO achieves competitive results, with an HOTA of 70.2 and 82.9 for Car tracking and Pedestrian tracking, respectively. These results position MO-YOLO as a strong contender in multi-object tracking on the KITTI.

In conclusion, MO-YOLO stands out as a highly advantageous model in the realm of multi-object tracking, demonstrating superior object detection accuracy and overall tracking performance compared to MOTR, as evidenced by competitive results on diverse datasets, including Dancetrack, MOT17, and KITTI. The notable improvements observed in MeMOTR relative to MOTR suggest that MO-YOLO has significant untapped potential for performance enhancement. Positioned as a robust alternative, MO-YOLO can serve as an excellent new baseline, driving advancements in end-to-end models within the field of multi-object tracking.

### 4.4    Ablation Study

To assess the proposed training strategy and the efficacy of TBSP, ablation experiments were conducted on the Dancetrack dataset. The training set includes the entire Dancetrack training data, with testing involving 25 sequences from the validation set. Additionally, the impact of two methods for generating Query Positions—utilizing the formula from DAB-DETR [19] or employing an MLP—on the model's performance was investigated. Notably, MO-YOLO-L, MO-YOLO-M, and MO-YOLO-S differ based on their foundations in YOLOv8-L,YOLOv8-M, and YOLOv8-S, respectively, providing insights into model variations and their implications.

(a) Performance of MO-YOLO at different training stages

| Training Stage | HOTA | AssA | MOTA |
|---|---|---|---|
| Initial | 38.5 | 21.6 | 76.8 |
| Intermediate | 44.5 | 25.5 | 84.2 |
| Final (TBSP) | 51.9 | 34.6 | 84.9 |
| Final | 53.6 | 35.9 | 84.7 |

(b) Performance across different model sizes

| Model | HOTA | MOTA | FPS |
|---|---|---|---|
| MO-YOLO-L | 53.6 | 84.7 | 19.6 |
| MO-YOLO-M | 49.7 | 82.9 | 22.7 |
| MO-YOLO-S | 42.9 | 79.3 | 26.2 |

(c) The impact of $\delta_{iou}$ on model performance when using TBSP

| $\delta_{iou}$ | HOTA | AssA | MOTA |
|---|---|---|---|
| 0.3 | 49.7 | 32.3 | 79.1 |
| 0.5 | 51.9 | 34.6 | 84.9 |
| 0.7 | 51.3 | 34.8 | 83.2 |
| 0.9 | 50.0 | 32.8 | 81.3 |

(d) Effect of query position generation strategy

| Strategy | HOTA | MOTA | FPS |
|---|---|---|---|
| MLP | 53.3 | 85.1 | 18.5 |
| Formula | 53.6 | 84.7 | 19.6 |

**Table 5:** Ablation Studies on Our Proposed MO-YOLO on the DanceTrack validation set.

The performance of MO-YOLO improves consistently across different training stages, peaking in the final stage. Different model sizes (MO-YOLO-L, MO-YOLO-M, MO-YOLO-S) achieve a balance between accuracy and computational efficiency. The similarity in their speeds is likely attributable to the consistent scale of the decoder. The choice of $\delta_{iou}$ impacts precision and recall trade-offs, with higher values leading to a slight decrease in performance. Position generation strategies, MLP and Formula, yield similar performance in terms of detection accuracy, with slight differences in FPS. Overall, the ablation studies demonstrate the effectiveness of MO-YOLO across different training conditions and model configurations, showcasing its adaptability and versatility in MOT field.



MOTR                                MO-YOLO

**Fig. 5:** MO-YOLO and MOTR failed in the same scenario. The probable cause may lie in the insufficient scale of the MOT17 dataset.

## 5    Conclusion and Limitation

This paper presents MO-YOLO, an efficient and lightweight end-to-end multi-object tracking model that optimizes computational resources. By integrating YOLO and RT-DETR principles, the study establishes a novel tracking paradigm with accelerated training strategies. While MO-YOLO outperforms MOTR in resource efficiency, it marginally lags behind in performance, attributed to factors like YOLO's detection capabilities and reduced input image dimensions ($640 \times 640$). Despite enhancing computational efficiency, the dimension reduction raises concerns about potential information loss, impacting overall effectiveness. Additionally, latent issues from MOTR may surface in MO-YOLO, as depicted in Fig.5. This research contributes to advancing real-time computer vision applications.

# References

1. Aharon, N., Orfaig, R., Bobrovsky, B.Z.: BoT-SORT: Robust Associations Multi-Pedestrian Tracking (Jul 2022). https://doi.org/10.48550/arXiv.2206.14651 1

2. Bergmann, P., Meinhardt, T., Leal-Taixe, L.: Tracking without bells and whistles. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 941–951 (2019) 1

3. Bewley, A., Ge, Z., Ott, L., Ramos, F., Upcroft, B.: Simple online and realtime tracking. In: 2016 IEEE International Conference on Image Processing (ICIP). pp. 3464–3468 (Sep 2016). https://doi.org/10.1109/ICIP.2016.7533003 1, 4

4. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. Advances in neural information processing systems **33**, 1877–1901 (2020) 2

5. Cao, J., Pang, J., Weng, X., Khirodkar, R., Kitani, K.: Observation-centric sort: Rethinking sort for robust multi-object tracking. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9686–9696 (2023) 11, 12

6. Cao, J., Wu, H., Kitani, K.: Track targets by dense spatio-temporal position encoding. arXiv preprint arXiv:2210.09455 (2022) 11

7. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: European conference on computer vision. pp. 213–229. Springer (2020) 2

8. Chen, S., Sun, P., Song, Y., Luo, P.: Diffusiondet: Diffusion model for object detection. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 19830–19843 (2023) 4

9. Du, Y., Zhao, Z., Song, Y., Zhao, Y., Su, F., Gong, T., Meng, H.: StrongSORT: Make DeepSORT Great Again. IEEE Transactions on Multimedia pp. 1–14 (2023). https://doi.org/10.1109/TMM.2023.3240881 1

10. Gao, R., Wang, L.: Memotr: Long-term memory-augmented transformer for multi-object tracking. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9901–9910 (2023) 3, 9, 11, 12

11. Ge, Z., Liu, S., Wang, F., Li, Z., Sun, J.: Yolox: Exceeding yolo series in 2021. arXiv preprint arXiv:2107.08430 (2021) 3

12. Geiger, A., Lenz, P., Urtasun, R.: Are we ready for autonomous driving? The KITTI vision benchmark suite. In: 2012 IEEE Conference on Computer Vision and Pattern Recognition. pp. 3354–3361 (Jun 2012). https://doi.org/10.1109/CVPR.2012.6248074 9, 12

13. Ghiasi, G., Cui, Y., Srinivas, A., Qian, R., Lin, T.Y., Cubuk, E.D., Le, Q.V., Zoph, B.: Simple copy-paste is a strong data augmentation method for instance segmentation. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 2918–2928 (2021) 8

14. Jocher, G., Chaurasia, A., Qiu, J.: Ultralytics yolov8 (2023), https://github.com/ultralytics/ultralytics 2, 9

15. Kim, A., Brasó, G., Ošep, A., Leal-Taixé, L.: Polarmot: How far can geometric relations take us in 3d multi-object tracking? In: European Conference on Computer Vision. pp. 41–58. Springer (2022) 12

16. Li, F., Zhang, H., Liu, S., Guo, J., Ni, L.M., Zhang, L.: Dn-detr: Accelerate detr training by introducing query denoising. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13619–13627 (2022) 7

17. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: Proceedings of the IEEE international conference on computer vision. pp. 2980–2988 (2017) 8

18. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13. pp. 740–755. Springer (2014) 10, 19

19. Liu, S., Li, F., Zhang, H., Yang, X., Qi, X., Su, H., Zhu, J., Zhang, L.: Dab-detr:dynamic anchor boxes are better queries for detr (2022) 7, 13

20. Luiten, J., Osep, A., Dendorfer, P., Torr, P., Geiger, A., Leal-Taixé, L., Leibe, B.: Hota: A higher order metric for evaluating multi-object tracking. International journal of computer vision 129, 548–578 (2021) 9

21. Lv, W., Zhao, Y., Xu, S., Wei, J., Wang, G., Cui, C., Du, Y., Dang, Q., Liu, Y.: DETRs Beat YOLOs on Real-time Object Detection (Jul 2023). https://doi.org/10.48550/arXiv.2304.08069 2, 4, 10

22. Maggiolino, G., Ahmad, A., Cao, J., Kitani, K.: Deep OC-SORT: Multi-Pedestrian Tracking by Adaptive Re-Identification (Feb 2023) 1, 12

23. Marinello, N., Proesmans, M., Van Gool, L.: Triplettrack: 3d object tracking using triplet embeddings and lstm. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4500–4510 (2022) 12

24. Milan, A., Leal-Taixe, L., Reid, I., Roth, S., Schindler, K.: MOT16: A Benchmark for Multi-Object Tracking (May 2016). https://doi.org/10.48550/arXiv.1603.00831 9, 12, 19

25. Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 779–788 (2016) 2, 3

26. Rezatofighi, H., Tsoi, N., Gwak, J., Sadeghian, A., Reid, I., Savarese, S.: Generalized intersection over union: A metric and a loss for bounding box regression. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 658–666 (2019) 8

27. Ristani, E., Solera, F., Zou, R., Cucchiara, R., Tomasi, C.: Performance measures and a data set for multi-target, multi-camera tracking. In: European conference on computer vision. pp. 17–35. Springer (2016) 9

28. Shao, S., Zhao, Z., Li, B., Xiao, T., Yu, G., Zhang, X., Sun, J.: Crowdhuman: A benchmark for detecting human in a crowd. arXiv preprint arXiv:1805.00123 (2018) 10

29. Sun, P., Cao, J., Jiang, Y., Yuan, Z., Bai, S., Kitani, K., Luo, P.: DanceTrack: Multi-Object Tracking in Uniform Appearance and Diverse Motion. In: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 20961–20970. IEEE, New Orleans, LA, USA (Jun 2022). https://doi.org/10.1109/CVPR52688.2022.02032 1, 9, 11

30. Tokmakov, P., Li, J., Burgard, W., Gaidon, A.: Learning to track with object permanence. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10860–10869 (2021) 12

31. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. Advances in neural information processing systems 30 (2017) 2

32. Wang, C.Y., Bochkovskiy, A., Liao, H.Y.M.: Scaled-YOLOv4: Scaling cross stage partial network. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 13029–13038 (June 2021) 8

33. Wang, C.Y., Bochkovskiy, A., Liao, H.Y.M.: Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7464–7475 (2023) 22
34. Wang, C.Y., Yeh, I.H., Liao, H.Y.M.: Yolov9: Learning what you want to learn using programmable gradient information. arXiv preprint arXiv:2402.13616 (2024) 22
35. Wang, Y., Zhang, X., Yang, T., Sun, J.: Anchor detr: Query design for transformer-based detector. In: Proceedings of the AAAI conference on artificial intelligence. vol. 36, pp. 2567–2575 (2022) 7
36. Xia, Z., Pan, X., Song, S., Li, L.E., Huang, G.: Vision transformer with deformable attention. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 4794–4803 (2022) 4
37. Yan, F., Luo, W., Zhong, Y., Gan, Y., Ma, L.: Bridging the Gap Between End-to-end and Non-End-to-end Multi-Object Tracking 3
38. Yan, F., Luo, W., Zhong, Y., Gan, Y., Ma, L.: Bridging the Gap Between End-to-end and Non-End-to-end Multi-Object Tracking (May 2023). https://doi.org/10.48550/arXiv.2305.12724 1, 7, 9
39. You, S., Yao, H., Bao, B.K., Xu, C.: Utm: A unified multiple object tracking model with identity-aware feature enhancement. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 21876–21886 (2023) 12
40. Yu, E., Wang, T., Li, Z., Zhang, Y., Zhang, X., Tao, W.: MOTRv3: Release-Fetch Supervision for End-to-End Multi-Object Tracking 1, 3, 7, 9, 10
41. Yu, F., Chen, H., Wang, X., Xian, W., Chen, Y., Liu, F., Madhavan, V., Darrell, T.: Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 2636–2645 (2020) 10
42. Zeng, F., Dong, B., Zhang, Y., Wang, T., Zhang, X., Wei, Y.: Motr: End-to-end multiple-object tracking with transformer. In: European Conference on Computer Vision. pp. 659–675. Springer (2022) 1, 2, 5, 10, 11, 12, 19
43. Zhang, H., Li, F., Liu, S., Zhang, L., Su, H., Zhu, J., Ni, L.M., Shum, H.Y.: Dino:detr with improved denoising anchor boxes for end-to-end object detection (2022) 4, 7
44. Zhang, H., Cisse, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk minimization. arXiv preprint arXiv:1710.09412 (2017) 8
45. Zhang, Y., Sun, P., Jiang, Y., Yu, D., Weng, F., Yuan, Z., Luo, P., Liu, W., Wang, X.: ByteTrack: Multi-object Tracking by Associating Every Detection Box. In: Avidan, S., Brostow, G., Cissé, M., Farinella, G.M., Hassner, T. (eds.) Computer Vision – ECCV 2022. pp. 1–21. Lecture Notes in Computer Science, Springer Nature Switzerland, Cham (2022) 11, 12
46. Zhang, Y., Wang, C., Wang, X., Zeng, W., Liu, W.: Fairmot: On the fairness of detection and re-identification in multiple object tracking. International Journal of Computer Vision 129, 3069–3087 (2021) 12
47. Zhang, Y., Wang, T., Zhang, X.: Motrv2: Bootstrapping end-to-end multi-object tracking by pretrained object detectors. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 22056–22065 (2023) 3, 9, 10
48. Zhou, X., Koltun, V., Krähenbühl, P.: Tracking objects as points. In: European conference on computer vision. pp. 474–490. Springer (2020) 12

49. Zhou, X., Yin, T., Koltun, V., Krähenbühl, P.: Global tracking transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8771–8780 (2022) 11, 12
50. Zhu, X., Su, W., Lu, L., Li, B., Wang, X., Dai, J.: Deformable detr: Deformable transformers for end-to-end object detection. arXiv preprint arXiv:2010.04159 (2020) 4, 7, 19
51. Zong, Z., Song, G., Liu, Y.: Detrs with collaborative hybrid assignments training. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 6748–6758 (2023) 4

# A   Appendix

## A.1   Disclaimer

The data used in this paper is sourced from publicly available datasets related to human subjects. While every effort has been made to ensure the accuracy and reliability of the data, the authors make no representations or warranties of any kind, express or implied, about the completeness, accuracy, reliability, suitability, or availability with respect to the data. The authors shall not be liable for any loss or damage, including but not limited to indirect or consequential loss or damage, arising from the use of the data in this research paper. Users of the data are encouraged to independently verify the accuracy and reliability of the data before making any decisions or taking any actions based on it.

## A.2   More Comprehensive Analysis of Experimental Results and Model Parameter

The initial motivation for building MO-YOLO was not because of its slow inference speed, but rather because of its high training resource requirements. With 8 V100 GPUs using pre-trained weights, it took approximately 84 hours to train 200 epochs on an extended MOT17 [24] with only about 10,000 images. Each epoch took around 0.42 hours, and on average, each GPU could train on approximately 3,000 images per hour.

It is worth noting that MOTR [42] is based on Deformable-DETR [50], which was trained on the MSCOCO dataset [18] with approximately 110,000 images. In single-scale mode (which can be considered similar to MOTR's training), using eight V100 GPUs, each epoch took approximately 0.4 hours, and on average, each GPU could train on around 30,000 images per hour.

Upon careful consideration of the reasons, the first aspect to be contemplated is the batch size. The collective average loss (CAL) mechanism and the tracklet-aware label assignment (TALA) strategy employed in MOTR enable the network to learn some latent temporal information, consequently restricting the batch size on a single card to 1. Setting the batch size to 1 may result in underutilization of the corresponding computational resources (e.g., GPU, TPU, or other possible devices). Additionally, MOTR is trained in a one-shot manner, requiring the model not only to grasp motion information but also to learn object appearance features. Hence, we proposed a multi-step training strategy to first allow the model to learn appearance features before proceeding to learn object motion characteristics. Experimental results confirm the effectiveness of this strategy.

As for the TBSP, its initial purpose was simply to provide the detection network with a basic tracking capability. However, it also pleasantly surprised us by further reducing training time. Here, we also provide the required rounds and results if TALA were to be used directly without the TBSP strategy, as shown in Table 6. The training here is conducted on a single NVIDIA 4090 GPU. And this training process is not included in the first stage

**Table 6:** Comparison of MO-YOLO with and without using TBSP on the Dancetrack val dataset.

| Initial | Intermediate | Final | HOTA | MOTA | Epochs | Total Time |
|---------|--------------|-------|------|------|--------|------------|
| ✓ | ✓ | ✓ | 53.6 | 84.7 | 27 | 56.7 hours |
| ✓ | | ✓ | 52.9 | 83.8 | 40 | 85.2 hours |

In order to provide a more comprehensive illustration of the characteristics of our model, we list the sizes of different models, namely MO-YOLO-L (45.8M), MO-YOLO-M (29.9M), and MO-YOLO-S (19.7M), as depicted in Fig. 6. The
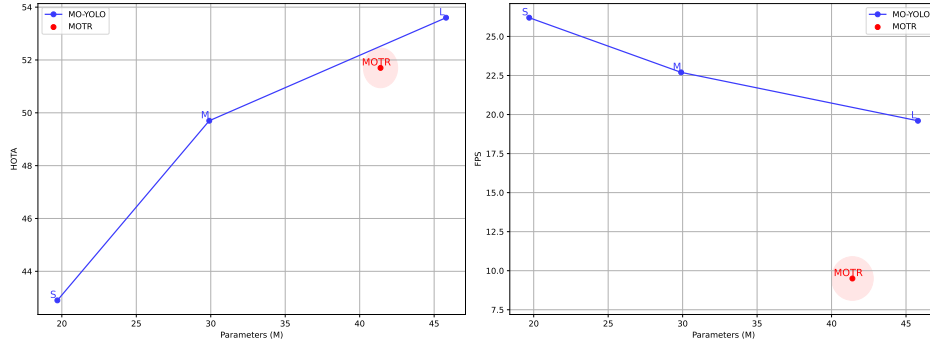


**Fig. 6:** The model sizes of MOTR and MO-YOLO, as well as their FPS and HOTA on the Dancetrack.

table compares the performance of MO-YOLO with and without using the TBSP strategy on the Dancetrack validation dataset across different training phases: Initial, Intermediate, and Final. Comparison of the two rows reveals that employing TBSP throughout all phases leads to slightly better performance in terms of HOTA and MOTA, with a significantly lower total training time. This underscores the importance of the TBSP strategy in reducing the overall training duration while maintaining or even improving tracking accuracy.

We hypothesize that the efficacy of the TBSP strategy may stem from its endowment of the model with a degree of autonomy, facilitating the preliminary acquisition of object motion information that would otherwise necessitate numerous training iterations. Analogously, the strategy resembles the use of training wheels in bicycle riding, which assist novices in achieving smooth operation while providing an opportunity to experience the sensation of riding. Consequently, this approach aids in the comprehensive mastery of the skill. We posit that a comparable scenario may manifest within this context.

### A.3    RT-MOTR

At the same time, to demonstrate the advantages of our MO-YOLO models constructed using YOLO and RT-DETR relative to end-to-end tracking models built directly on RT-DETR, we have also constructed such an end-to-end tracking model based on RT-DETR (referred to as RT-MOTR), and conducted performance testing on it.

**NOTE**: in Table 7, both the train time and epochs refer to the second and third stages like Table 6. For the first stage, both models were trained for 15 epochs, and the training device is a single 4090, while the testing device is a single V100.

The table compares the performance of MO-YOLO and RT-MOTR on the Dancetrack validation dataset. It presents various metrics including HOTA (Higher Order Tracking Accuracy), training epochs, training speed (hours per epoch), total training time, and frames per second (FPS).

MO-YOLO achieves a HOTA of 53.6 with 27 training epochs, an average training speed of 2.08 hours per epoch, a total training time of 56.7 hours, and an FPS of 19.6.

In comparison, RT-MOTR achieves a HOTA of 53.4 with 35 training epochs, an average training speed of 2.27 hours per epoch, a total training time of 90.8 hours, and an FPS of 20.3.

Comparing the performance metrics of the two methods, we observe that MO-YOLO slightly outperforms RT-MOTR in terms of HOTA, while also having a shorter training time albeit with a slightly lower FPS. This indicates that MO-YOLO exhibits advantages in both tracking accuracy and training efficiency.

**Table 7:** Comparison of MO-YOLO and RT-MOTR on the Dancetrack val dataset.

| Methods | HOTA | Train Epochs | Speed (hours/epoch) | Train Times | FPS |
|---------|------|--------------|---------------------|-------------|-----|
| MO-YOLO | 53.6 | 27 | 2.08 | 56.7 h | 19.6 |
| RT-MOTR | 53.4 | 35 | 2.27 | 79.5 h | 20.3 |

Regarding why the speeds of these two models are comparable, our experimental findings suggest that the primary time-consuming component is the decoder responsible for tracking. Therefore, optimizing the decoder represents a crucial area for future investigation.

This outcome underscores the rationale behind integrating elements from YOLOv8 and RT-DETR to construct a rapid end-to-end tracking network. Furthermore, the observed disparity in training efficiency between MO-YOLO and RT-DETR can be attributed to the former's notably simpler architecture.

Moreover, our decision to adopt YOLO as the backbone of the network is motivated by its propensity for swifter updates compared to DETR-based models.

Recent developments such as YOLOv9 [34] exemplify this trend. While YOLOv9 may trade off training time for enhanced precision and speed (resulting in at least a twofold increase in training duration compared to YOLOv7 [33]), we anticipate future iterations will offer refined improvements. Such advancements will, in turn, facilitate further enhancements to MO-YOLO.