

SimpleTrack: Understanding and Rethinking 3D Multi-object Tracking

Ziqi Pang^{1*} Zhichao Li² Naiyan Wang²
 UIUC¹ TuSimple²

ziqip2@illinois.edu, {leeisabug, winsty}@gmail.com

Abstract

3D multi-object tracking (MOT) has witnessed numerous novel benchmarks and approaches in recent years, especially those under the “tracking-by-detection” paradigm. Despite their progress and usefulness, an in-depth analysis of their strengths and weaknesses is not yet available. In this paper, we summarize current 3D MOT methods into a unified framework by decomposing them into four constituent parts: pre-processing of detection, association, motion model, and life cycle management. We then ascribe the failure cases of existing algorithms to each component and investigate them in detail. Based on the analyses, we propose corresponding improvements which lead to a strong yet simple baseline: SimpleTrack. Comprehensive experimental results on Waymo Open Dataset and nuScenes demonstrate that our final method could achieve new state-of-the-art results with minor modifications.

Furthermore, we take additional steps and rethink whether current benchmarks authentically reflect the ability of algorithms for real-world challenges. We delve into the details of existing benchmarks and find some intriguing facts. Finally, we analyze the distribution and causes of remaining failures in SimpleTrack and propose future directions for 3D MOT. Our code is available at <https://github.com/TuSimple/SimpleTrack>.

1. Introduction

Multi-object tracking (MOT) is a composite task in computer vision, combining both the aspects of localization and identification. Given its complex nature, MOT systems generally involve numerous interconnected parts, such as the selection of detections, the data association, the modeling of object motions, etc. Each of these modules has its special treatment and can significantly affect the system performance as a whole. Therefore, we would like to ask *which components in 3D MOT play the most important roles, and how can we improve them?*

Bearing such objectives, we revisit the current 3D MOT algorithms [3, 10, 28, 37, 43, 44]. These methods mostly adopt the “tracking by detection” paradigm, where they directly take the bounding boxes from 3D detectors and build up tracklets across frames. We first break them down into four individual modules and examine each of them: pre-processing of input detections, motion model, association, and life cycle management. Based on this modular framework, we locate and ascribe the failure cases of 3D MOT to the corresponding components and discover several overlooked issues in the previous designs.

First, we find that inaccurate input detections may contaminate the association. However, purely pruning them by a score threshold will sacrifice the recall. Second, we find that the similarity metric defined between two 3D bounding boxes need to be carefully designed. Neither distance-based nor simple IoU works well. Third, the object motion in 3D space is more predictable than that in the 2D image space. Therefore, the consensus between motion model predictions and even poor observations (low score detections) could well indicate the existence of objects. Illuminated by these observations, we propose several simple yet non-trivial solutions. The evaluation on Waymo Open Dataset [34] and nuScenes [8] suggests that our final method “SimpleTrack” is competitive among the 3D MOT algorithms (in Tab. 6 and Tab. 7).

Besides analyzing 3D MOT algorithms, we also reflect on current benchmarks. We emphasize the need for high-frequency detections and the proper handling of output tracklets in evaluation. To better understand the upper bound of our method, we further break down the remaining errors based on ID switch and MOTA metrics. We believe these observations could inspire the better design of algorithms and benchmarks.

In brief, our contributions are as follow:

- We decompose the pipeline of “tracking-by-detection” 3D MOT framework and analyze the connections between each component and failure cases.
- We propose corresponding treatments for each module and combine them into a simple baseline. The results are competitive on the Waymo Open Dataset and

*This work is complete during the first author’s internship at TuSimple.

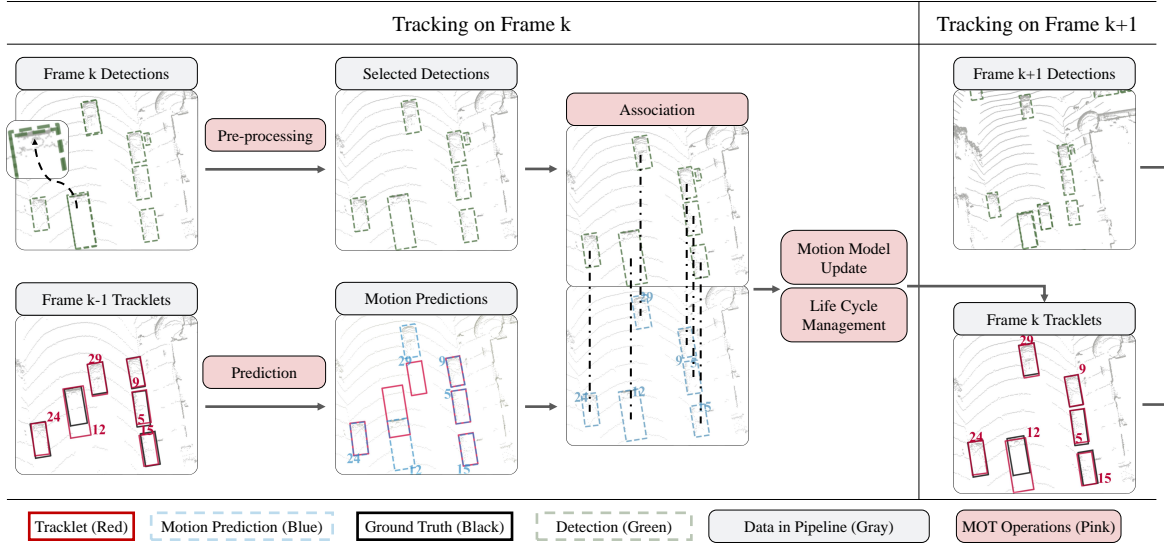


Figure 1. 3D MOT pipeline. For simplicity, we only visualize the steps between frame k and frame $k+1$. Best view in color.

nuScenes.

- We also analyze existing 3D MOT benchmarks and explain the potential influences of their designs. We hope that our analyses could shed light for future research.

2. Related Work

Most 3D MOT methods [3, 10, 28, 37, 43, 44] adopt the “tracking-by-detection” framework because of the strong power of detectors. We first summarize the representative 3D MOT work and then highlight the connections and distinctions between 3D and 2D MOT.

2.1. 3D MOT

Many 3D MOT methods are composed of hand-crafted rule-based components. AB3DMOT [37] is the common baseline of using IoU for association and a Kalman filter as the motion model. Its notable followers mainly improve on the association part: Chiu *et al.* [10] and CenterPoint [43] replace IoU with Mahalanobis and L2 distance, which performs better on nuScenes [8]. Some others notice the importance of life cycle management, where CBMOT [3] proposes a score-based method to replace the “count-based” mechanism, and Pöschmann *et al.* [28] treats 3D MOT as optimization problems on factor graphs. Despite the effectiveness of these improvements, a systematic study on 3D MOT methods is in great need, especially where these designs suffer and how to make further improvements. To this end, our paper seeks to meet the expectations.

Different from the methods mentioned above, many others attempt to solve 3D MOT with fewer manual designs. [2, 9, 15, 38] leverage rich features from RGB images for association and life cycle control, and Chiu *et al.* [9]

specially uses neural networks to handle the feature fusion, association metrics, and tracklet initialization. Recently, OGR3MOT [44] follows Guillem *et al.* [7] and solves 3D MOT with graph neural networks (GNN) in an end-to-end manner, focusing on the data association and the classification of active tracklets, especially.

2.2. 2D MOT

2D MOT shares the common goal of data association with 3D MOT. Some notable attempts include probabilistic approaches [1, 16, 30, 32], dynamic programming [11], bipartite matching [6], min-cost flow [4, 46], convex optimization [27, 35, 36, 45], and conditional random fields [42]. With the rapid progress of deep learning, many methods [7, 12–14, 19, 40] learn the matching mechanisms and others [17, 20, 21, 24, 26] learn the association metrics.

Similar to 3D MOT, many 2D trackers [5, 22, 33, 48] also benefit from the enhanced detection quality and adopt the “tracking-by-detection” paradigm. However, the objects on RGB images have varied sizes because of scale variation; thus, they are harder for association and motion models. But 2D MOT can easily take advantage of rich RGB information and use appearance models [18, 19, 33, 39], which is not available in LiDAR based 3D MOT. In summary, the design of MOT methods should fit the traits of each modality.

3. 3D MOT Pipeline

In this section, we decompose 3D MOT methods into the following four parts. An illustration is in Fig. 1.

Pre-processing of Input Detections. It pre-processes the bounding boxes from detectors and selects the ones to be

used for tracking. Some exemplar operations include selecting the bounding boxes with scores higher than a certain threshold. (In “Pre-processing” Fig. 1, some redundant bounding boxes are removed.)

Motion Model. It predicts and updates the states of objects. Most 3D MOT methods [3, 10, 37] directly use the Kalman filter, and CenterPoint [43] uses the velocities predicted by detectors from multi-frame data. (In “Prediction” and “Motion Model Update” Fig. 1.)

Association. It associates the detections with tracklets. The association module involves two steps: similarity computation and matching. The similarity measures the distance between a pair of detection and tracklet, while the matching step solves the correspondences based on the pre-computed similarities. AB3DMOT [37] proposes the baseline of using IoU with Hungarian algorithm, while Chiu *et al.* [10] uses Mahalanobis distance and greedy algorithm, and CenterPoint [43] adopts the L2 distance. (In “Association” Fig. 1.)

Life Cycle Management. It controls the “birth”, “death” and “output” policies. “Birth” determines whether a detection bounding box will be initialized as a new tracklet; “Death” removes a tracklet when it is believed to have moved out of the attention area; “Output” decides whether a tracklet will output its state. Most of the MOT algorithm adopts a simple count-based rule [10, 37, 43], and CB-MOT [3] improves birth and death by amending the logic of tracklet confidences. (In “Life Cycle Management” Fig. 1.)

4. Analyzing and Improving 3D MOT

In this section, we analyze and improve each module in the 3D MOT pipeline. For better clarification, we ablate the effects of every modification by removing it from the final variant of SimpleTrack. By default, the ablations are all on the validation split with the CenterPoint [43] detection. We also provide additive ablation analyses and the comparison with other methods in Sec. 4.5.

4.1. Pre-processing

To fulfill the recall requirements, current detectors usually output a large number of bounding boxes with scores roughly indicating their quality. However, if these boxes are treated equally in the association step of 3D MOT, the bounding boxes with low quality or severe overlapping may deviate the trackers to select the inaccurate detection for extending or forming tracklets (as in the “raw detection” of Fig. 2). Such a gap between the detection and MOT task needs careful treatment.

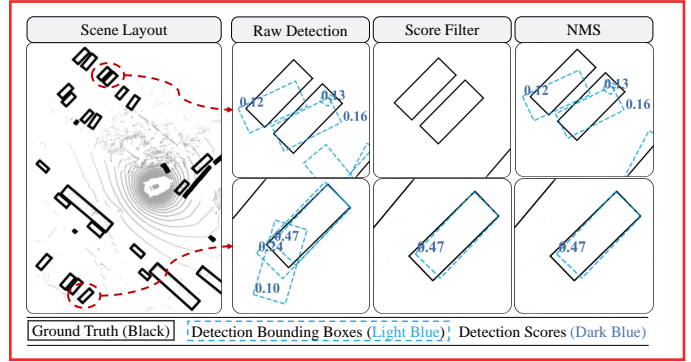


Figure 2. Comparison between score filtering and NMS. To remove the redundant bounding boxes on row 2, score filtering needs at least a 0.24 threshold, but this will eliminate the detections on row 1. However, NMS can well satisfy both by removing the overlapping on row 2 and maintaining the recall on row 1.

NMS	AMOTA↑	AMOTP↓	MOTA↑	IDS↓
×	0.673	0.574	0.581	557
✓	0.687	0.573	0.592	519

Table 1. Ablation for NMS on nuScenes.

NMS	Vehicle			Pedestrian		
	MOTA↑	MOTP↓	IDS(%)↓	MOTA↑	MOTP↓	IDS(%)↓
×	0.5609	0.1681	0.09	0.4962	0.3090	5.00
✓	0.5612	0.1681	0.08	0.5776	0.3125	0.42

Table 2. Ablation for NMS on WOD.

3D MOT methods commonly use confidence scores to filter out the low-quality detections and improve the precision of input bounding boxes for MOT. However, such an approach may be detrimental to the recall as they directly abandon the objects with poor observations (top row in Fig. 2). It is also especially harmful to metrics like AMOTA, which needs the tracker to use low score bounding boxes to fulfill the recall requirements.

To improve the precision without significantly decreasing the recall, our solution is simple and direct: we apply stricter **non-maximum suppression (NMS)** to the input detections. As shown in the right of Fig. 2, the NMS operation alone can effectively eliminate the overlapped low-quality bounding boxes while keeping the diverse low-quality observations, even on regions like sparse points or occlusion. *Therefore, by adding NMS to the pre-processing module, we could roughly keep the recall, but greatly improves the precision and benefits MOT.*

On WOD, our stricter NMS operation removes 51% and 52% bounding boxes for vehicles and pedestrians and nearly doubles the precision: 10.8% to 21.1% for vehicles, 5.1% to 9.9% for pedestrians. At the same time, the recall drops relatively little from 78% to 74% for vehicles and 83% to 79% for pedestrians. According to Tab. 1 and Tab. 2, this largely benefits the performance, especially on

Method	Vehicle			Pedestrian		
	MOTA↑	MOTP↓	IDS(%)↓	MOTA↑	MOTP↓	IDS(%)↓
KF	0.5612	0.1681	0.08	0.5776	0.3125	0.42
CV	0.5515	0.1691	0.14	0.5661	0.3159	0.58
KF PD	0.5516	0.1691	0.14	0.5654	0.3158	0.63

Table 3. Comparison of motion models on Waymo Open Dataset. “KF” denotes Kalman filters; “CV” denotes constant velocity model; “KF-PD” denotes the variant using Kalman filter only for motion prediction. Details in Sec. 4.2.

Method	AMOTA↑	AMOTP↓	MOTA↑	IDS↓
KF	0.687	0.573	0.592	519
CV	0.690	0.564	0.592	516

Table 4. Comparison of motion models on nuScenes. Abbreviations are identical to Tab. 3. Details in Sec. 4.2.

the pedestrian (right part of Tab. 2), where the object detection task is harder.

4.2. Motion Model

Motion models depict the motion status of tracklets. They are mainly used to predict the candidate states of objects in the next frame, which are the proposals for the following association step. Furthermore, the motion models like the Kalman filter can also potentially refine the states of objects. In general, there are two commonly adopted motion models for 3D MOT: Kalman filter (KF), *e.g.* AB3DMOT [37], and constant velocity model (CV) with **predicted speeds from detectors, *e.g.* CenterPoint [43]**. The advantage of KF is that it could utilize the information from multiple frames and provide smoother results when facing low-quality detection. Meanwhile, CV deals better with abrupt and unpredictable motions with its explicit speed predictions, but its effectiveness on motion smoothing is limited. In Tab. 3 and Tab. 4, we compare the two of them on WOD and nuScenes, which provides clear evidence for our claims.

In general, these two motion models demonstrate similar performance. On nuScenes, CV marginally outperforms KF, while it is the opposite on WOD. The advantages of KF on WOD mainly come from the refinement for the bounding boxes. To verify this, we implement the “KF-PD” variant, which uses KF only for providing motion predictions prior to association, and the outputs are all original detections. Eventually, the marginal gap between “CV” and “KF-PD” in Tab. 3 supports our claim. On nuScenes, the CV motion model is slightly better due to the lower frame rates on nuScenes (2Hz). To prove our conjecture, we apply KF and CV both under a higher frequency 10Hz setting on nuScenes¹, and KF marginally outperforms CV by 0.696 versus 0.693 in AMOTA this time.

To summarize, *the Kalman Filter fits better for high-*

¹Please check Sec. 5.1 for how we build 10Hz settings on nuScenes.

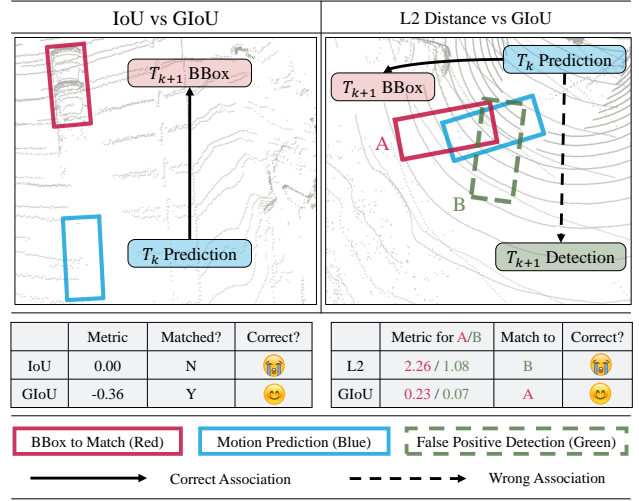


Figure 3. Illustration of association metrics. Left: IoU versus GIoU. Right: L2 Distance versus GIoU. Details are in Sec. 4.3.1.

frequency cases because of more predictable motions, and the constant velocity model is more robust for low-frequency scenarios with explicit speed prediction. Since inferring the speed is not yet common for detectors, we adopt the Kalman filter for SimpleTrack without loss of generality.

4.3. Association

4.3.1 Association Metrics: 3D GIoU

IoU based [37] and distance based [10, 43] association metrics are the two prevalent choices in 3D MOT. As in Fig. 3, they have typical but different failure modes. IoU computes the overlapping ratios between bounding boxes, so it cannot connect the detections and motion predictions if the IoU between them are all zeros, which are common at the beginnings of tracklets or on objects with abrupt motions (the left of Fig. 3). The representatives for distance-based metrics are Mahalanobis [10] and L2 [43] distances. With larger distance thresholds, they can handle the failure cases of IoU based metrics, but they may not be sensitive enough for nearby detection with low quality. We explain such scenarios on the right of Fig. 3. On frame k , the blue motion prediction has smaller L2 distances to the green false positive detection, thus it is wrongly associated. Illuminating by such example, we conclude that the distance-based metrics lack discrimination on orientations, which is just the advantage of IOU based metrics.

To get the best of two worlds, we propose to generalize “Generalized IoU” (GIoU) [31] to 3D for association. Briefly speaking, for any pair of 3D bounding boxes B_1, B_2 , their 3D GIoU is as Eq. 1, where I, U are the intersection and union of B_1 and B_2 . C is the enclosing convex hull of U . V represents the volume of a polygon. We set $\text{GIoU} > -0.5$ as the threshold for every category of objects

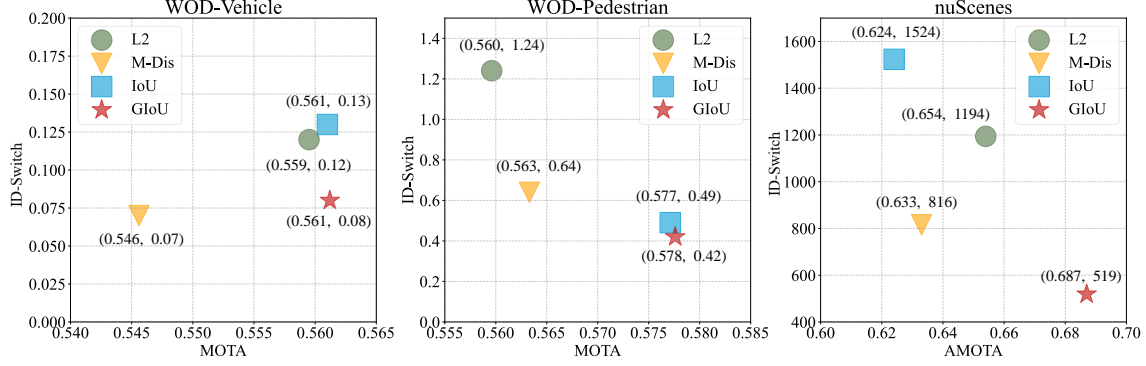


Figure 4. Comparison of association metrics on WOD (left & middle) and nuScenes (right). “M-Dis” is the short for Mahalanobis distance. The best method is closest to the bottom-right corner, having the lowest ID-Switches and highest MOTA/AMOTA.

on both WOD and nuScenes for this pair of associations to enter the subsequent matching step.

$$V_U = V_{B_1} + V_{B_2} - V_I, \quad (1)$$

$$\text{GIoU}(B_1, B_2) = V_I/V_U - (V_C - V_U)/V_C.$$

As in Fig 3, the GIoU metric can handle both patterns of failures. The quantitative results in Fig. 4 also show the ability of GIoU for improving the association on both WOD and nuScenes.

4.3.2 Matching Strategies

Generally speaking, there are two approaches for the matching between detections and tracklets: 1) Formulating the problem as a bipartite matching problem, and then solving it using Hungarian algorithm [37]. 2) Iteratively associating the nearest pairs by greedy algorithm [10, 43].

We find that these two methods heavily couples with the association metrics: IoU based metrics are fine with both, while distance-based metrics prefer greedy algorithms. We hypothesize that the reason is that the range of distance-based metrics are large, thus methods optimizing global optimal solution, like the Hungarian algorithm, may be adversely affected by outliers. In Fig. 5, we experiment with all the combinations between matching strategies and association metrics on WOD. As demonstrated, IoU and GIoU function well for both strategies, while Mahalanobis and L2 distance demand greedy algorithm, which is also consistent with the conclusions from previous work [10].

4.4. Life Cycle Management

We analyze all the ID-Switches on WOD², and categorize them into two groups as in Fig. 6: wrong association and early termination. Different from the major focus of many work, which is association, we find that the early termination is actually the dominating cause of ID-Switches:

²We use py-motmetrics [23] for the analysis.

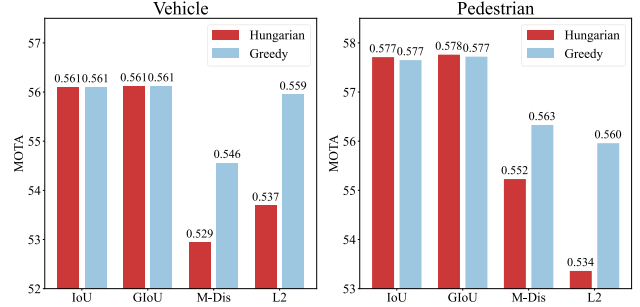


Figure 5. Comparison of matching strategies on WOD.

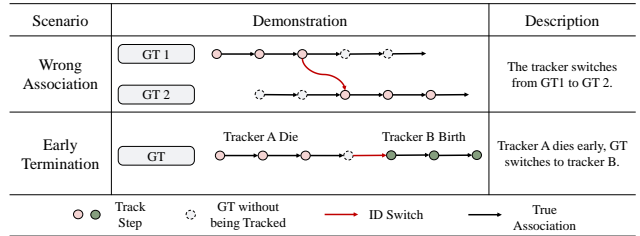


Figure 6. Illustration for two major types of ID-Switches.

95% for vehicle and 91% for pedestrian. Among the early terminations, many of them are caused by point cloud sparsity and spatial occlusion. To alleviate this issue, we utilize the free yet effective information: consensus between motion models and detections with low scores. *These bounding boxes are usually of low localization quality, however they are strong indication of the existence of objects if they agree with the motion predictions.* Then we use these to extend the lives of tracklets.

Bearing such motivation, we propose “Two-stage Association.” Specifically, we apply two rounds of association with different score thresholds: a low one T_l and a high one T_h (e.g. 0.1 and 0.5 for pedestrian on WOD). In stage one, we use the identical procedure as most current algorithms [10, 37, 43]: only the bounding boxes with scores higher than T_h are used for association. In stage two, we focus on the tracklets unmatched to detections in stage one

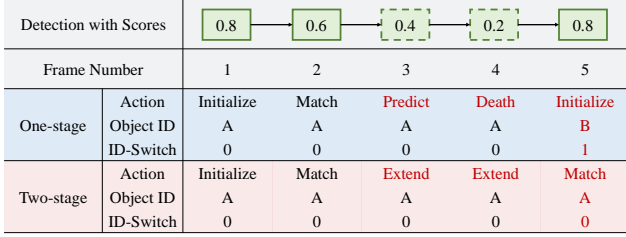


Figure 7. Comparison for “One-stage” and “Two-stage” association with a hypothetical example. “Extend” means “extending the life cycles,” and “Predict” means “using motion predictions due to no association.” Suppose $T_h = 0.5$ and $T_l = 0.1$ are the score thresholds, the “one-stage” method early terminates the tracklet because of consecutively lacking associations. Details in Sec. 4.4.

Method	Vehicle			Pedestrian		
	MOTA↑	MOTP↓	IDS(%)↓	MOTA↑	MOTP↓	IDS(%)↓
One	0.5567	0.1682	0.46	0.5718	0.3125	0.96
Two	0.5612	0.1681	0.08	0.5776	0.3125	0.42

Table 5. Ablation for “Two-stage Association” on WOD. “One” and “Two” denotes the previous one-stage association and our two-stage association methods. Details in Sec. 4.4.

and relax the conditions on their matches: detections having scores larger than T_l will be sufficient for a match. If the tracklet is successfully associated with one bounding box in stage two, it will still keep being alive. However, as the low score detections are generally in poor quality, we don’t output them to avoid false positives, and they are also not used for updating motion models. Instead, we use motion predictions as the latest tracklet states, replacing the low quality detections.

We intuitively explain the differences between our “Two-stage Association” and traditional “One-stage Association” in Fig. 7. Suppose $T = 0.5$ is the original score threshold for filtering detection bounding boxes, the trackers will then neglect the boxes with scores 0.4 and 0.2 on frames 3 and 4, which will die because of lacking matches in continuous frames and this eventually causes the final ID-Switch. In comparison, our two-stage association can maintain the active state of the tracklet.

In Tab. 5, our approach greatly decreases the ID-Switches without hurting the MOTA. This proves that SimpleTrack is effective in extending the life cycles by using detections more flexibly. Parallel to our work, a similar approach is also proven to be useful for 2D MOT [47].

4.5. Integration of SimpleTrack

In this section, we integrate the aforementioned techniques into the unified SimpleTrack and demonstrate how they improve the performance step by step.

In Fig. 8, we illustrate how the performance of 3D MOT trackers improve from the baselines. On WOD, although

Method	Vehicle			Pedestrian		
	MOTA↑	MOTP↓	IDS(%)↓	MOTA↑	MOTP↓	IDS(%)↓
AB3DMOT [37]	0.5773	0.1614	0.26	0.5380	0.3163	0.73
Chiu <i>et al.</i> [10]	0.4932	0.1689	0.62	0.4438	0.3227	1.83
CenterPoint [43]	0.5938	0.1637	0.32	0.5664	0.3116	1.07
SimpleTrack	0.6030	0.1623	0.08	0.6013	0.3114	0.40

Table 6. Comparison on WOD test split (L2). CenterPoint [43] detections are used. We mark the best in **red** and the second in **blue**. We list the methods using public detection. For AB3DMOT [37] and Chiu *et al.* [10], we report their best leaderboard entries.

Methods	AMOTA↑	AMOTP↓	MOTA↑	IDS ↓
AB3DMOT [37]	0.151	1.501	0.154	9027
Chiu <i>et al.</i> [10]	0.550	0.798	0.459	776
CenterPoint [43]	0.638	0.555	0.537	760
CBMOT [3]	0.649	0.592	0.545	557
OGR3MOT [44]	0.656	0.620	0.554	288
SimpleTrack (2Hz)	0.658	0.568	0.557	609
SimpleTrack (10Hz)	0.668	0.550	0.566	575

Table 7. Comparison on nuScenes test split. CenterPoint [43] detections are used. We list the methods using public detection. We mark the best in **red** and the second in **blue**. For CBMOT [3] and OGR3MOT [44], we report their numbers with CenterPoint [43] detection. Our numbers using both 2Hz and 10Hz frame rate detections are reported (details of our 10Hz setting are in Sec. 5).

the properties of vehicles and pedestrian are much different, each technique is applicable to both. On nuScenes, every proposed improvement is also effective for both the AMOTA and ID-Switch.

We also report the test set performance and compare with other 3D MOT methods. Combining our techniques leads to new state-of-the-art results (in Tab. 6, Tab. 7).³

5. Rethinking nuScenes

Besides the techniques mentioned above, we delve into the design of benchmarks. The benchmarks greatly facilitate the development of research and guide the designs of algorithms. Contrasting WOD and nuScenes, we have the following findings: 1) The frame rate of nuScenes is 2Hz, while WOD is 10Hz. Such low frequency adds unnecessary difficulties to 3D MOT (Sec. 5.1). 2) The evaluation of nuScenes requires high recalls with low score thresholds. And it also pre-processes the tracklets with interpolation, which encourages trackers to output the confidence scores reflecting the entire tracklet quality, but not the frame quality (Sec. 5.2). We hope these two findings could inspire the community to rethink the benchmarks and evaluation protocols of 3D tracking.

5.1. Detection Frequencies

Tracking generally benefits from higher frame rates, because motion is more predictable in short intervals. We

³Validation split comparisons are in the supplementary.

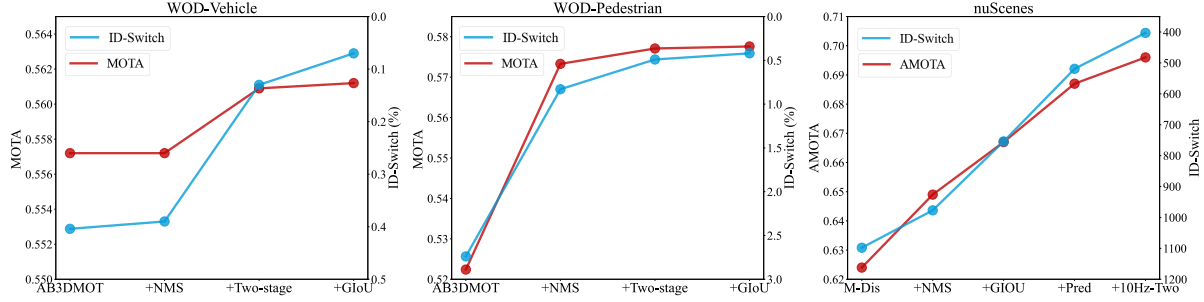


Figure 8. Improvements from SimpleTrack on WOD (left & middle) and nuScenes (right). We use the common baselines of AB3DMOT [37] on WOD and Chiu *et al.* [10] on nuScenes. For nuScenes, the improvements of “10Hz-Two” (using 10Hz detection and two-stage association) is in Sec. 5.1, and “Pred” (outputting motion model predictions) is in Sec. 5.2. The names for modifications are on the x-axis. Better MOTA and ID-Switch values are higher on the y-axis for clearer visualization.

Benchmark	Data	Annotation	Model
Waymo Open Dataset	10Hz	10Hz	10Hz
nuScenes	20Hz	2Hz	2Hz

Table 8. Frequency comparison of benchmarks.

compare the frequencies of point clouds, annotations, and common MOT frame rates on the two benchmarks in Tab. 8. On nuScenes, it has 20Hz point clouds but only 2Hz annotations. This leads to most common detectors and 3D MOT algorithms work under 2Hz, even they actually utilize all the 20Hz LiDAR data and operate faster than 2Hz. Therefore, we investigate the effect of high-frequency data as follows.

Although the information is more abundant with high frequency (HF) frames, it is non-trivial to incorporate them because nuScenes only evaluates on the low-frequency frames, which we refer to as “evaluation frames.” In Tab. 9, simply using all the 10Hz frames does not improve the performance. This is because the low-quality detection on the HF frames may deviate the trackers and hurt the performance on the sampled evaluation frames. To overcome this issue, we explore by first applying the “One-stage Association” on HF frames, where only the bounding boxes with scores larger than $T_h = 0.5$ are considered and used for motion model updating. We then adopt the “Two-stage Association” (described in Sec. 4.4) by using the boxes with scores larger than $T_l = 0.1$ to extend the tracklets. As in Tab. 9, our approach significantly improves both the AMOTA and ID-Switches. We also try to even increase the frame rate to 20Hz, but this barely leads to further improvements due to the deviation issue. So SimpleTrack uses the 10Hz setting in our final submission to the test set.⁴

5.2. Tracklet Interpolation

The AMOTA metric used in nuScenes calculates the average MOTAR [37] at different recall thresholds, which re-

Setting	AMOTA↑	AMOTP↓	MOTA↑	IDS↓
2Hz	0.687	0.573	0.592	519
10Hz	0.687	0.548	0.599	512
10Hz - One	0.696	0.564	0.603	450
10Hz - Two	0.696	0.547	0.602	403
20Hz - Two	0.690	0.547	0.598	416

Table 9. MOT with higher frame rates on nuScenes. “10Hz” is the vanilla baseline of using all the detections on high frequency (HF) frames. “-One” denotes “One-stage,” and “-Two” denotes “Two-stage.” Details in Sec. 5.1.

Predictions	AMOTA↑	AMOTP↓	MOTA↑	IDS↓	RECALL↑
×	0.667	0.612	0.572	754	0.696
✓	0.687	0.573	0.592	519	0.725

Table 10. Improvement from “outputting motion model predictions” on nuScenes (2Hz setting).

quires the trackers output the boxes of all score segments. In order to further improve the recall, we output the motion model predictions for frames and tracklets without associated detection bounding boxes, and empirically assign them lower scores than any other detection. In our case, their scores are $0.01 \times S_P$, where S_P is the confidence score of the tracklet in the previous fram. As shown in Tab. 10, this simple trick improves the overall recall and AMOTA.

However, we discover that enhancing the recall is not the only reason for such improvement. Besides the bounding boxes, the scores of the motion model predictions also make a significant contribution. This starts with the evaluation protocol on nuScenes, where they interpolate the input tracklets to fill in the missing frames and change all the scores with their tracklet-average scores as illustrated in Fig. 9. Under this context, our approach can explicitly penalize the low-quality tracklets, which generally contain more missing boxes replaced by motion model predictions.

In summary, such interpolation on nuScenes encourages the trackers to treat tracklet quality holistically and output calibrated quality-aware scores. However the quality of

⁴Because of the submission time limits to nuScenes test set, we are only able to report the “10Hz-One” variant in Tab. 7. It will be updated to “10Hz-Two” once we had the chance.

Detection with Scores		0.5	None	0.5	None	0.5
Frame Number		1	2	3	4	5
without Simple-Track Prediction	Tracker Output	0.5	None	0.5	None	0.5
	nuScenes Interpolate	0.5	0.5	0.5	0.5	0.5
with Simple-Track Prediction	Tracker Output	0.5	0.05	0.5	0.05	0.5
	nuScenes Interpolate	0.32	0.32	0.32	0.32	0.32

Figure 9. How the motion predictions and nuScenes interpolation changes tracklet scores. Dashed arrows are the directions for interpolation. On Frame 2 and 4 the boxes with score 0.05 are our motion predictions. The “0.5” and “0.32” are the tracklet-average scores with or without motion predictions. Details in Sec. 5.2.

Method	Vehicle				Pedestrian			
	MOTA↑	IDS(%)↓	FP↓	FN↓	MOTA↑	IDS(%)↓	FP↓	FN↓
SimpleTrack	0.561	0.078	0.104	0.334	0.578	0.425	0.109	0.309
GT Output	0.741	0.104	0.000	0.258	0.778	0.504	0.003	0.214
GT All	0.785	0.000	0.000	0.215	0.829	0.000	0.000	0.171

Table 11. Oracle Experiments on WOD.

boxes may vary a lot across frames even for the same tracklet, thus we suggest depicting the quality of a tracklet by only one score is imperfect. Moreover, future information is also introduced in this interpolation step and it changes the tracklet results. This could also bring the concern on whether the evaluation setting is still a fully online one.

6. Error Analyses

In this section, we conduct analyses on the remaining failure cases of SimpleTrack and propose potential future directions for improving “tracking by detection” paradigm. Without loss of generality, we use WOD as an example.

6.1. Upper Bound Experiment Settings

To quantitatively evaluate the causes of failure cases, we contrast SimpleTrack with two different oracle variants. The results are summarized in Tab. 11.

GT Output erases the errors caused by “output” policy. We compute the IoU between the bounding boxes from SimpleTrack with the GT boxes at the “output” stage, then use the IoU to decide if a box should be output instead of the detection score.⁵

GT All is the upper bound of tracking performance with CenterPoint boxes. We greedily match the detections from CenterPoint to GT boxes, keep the true positive and assign them ground-truth ID.

⁵The ID-Switch increases because we output more bounding boxes and IDs. The 0.003 false positives in pedestrians are caused by some boxes matching with the same GT box in crowded scenes.

6.2. Analyses for “Tracking By Detection”

ID-Switches. We break down the causes of ID-Switches as in Fig. 6. Although early termination has been greatly decreased by the scale of 86% for vehicle and 70% for pedestrian with “Two-stage Association,” it still takes up 88% and 72% failure cases in the remaining ID-Switches in SimpleTrack for vehicle and pedestrian, respectively. We inspect these cases and discover that most of them result from long-term occlusion or the returning of objects from being temporarily out of sight. Therefore, in addition to improving the association, potential future work can develop appearance models like in 2D MOT [18, 19, 33, 39] or silently maintain their states to re-identify these objects after they are back.

FP and FN. The “GT All” in Tab. 11 shows the upper bound for MOT with CenterPoint [43] detection, and we analyze the class of vehicle for example. Even with “GT All” the false negatives are still 0.215, which are the detection FN and can hardly be fixed under the “tracking by detection” framework. Comparing “GT All” and SimpleTrack, we find that the tracking algorithm itself introduces 0.119 false negatives. We further break them down as follows. Specifically, the difference between “GT Output” and “GT ALL” indicates that the 0.043 false negatives are caused by the uninitialized tracklets resulting from NMS and score threshold in pre-processing. The others come from life-cycle management. The “Initialization” requires two frames of accumulation before outputting a tracklet, which is same as AB3DMOT [37]. This yields a marginal 0.005 false negatives. Our “Output” logic uses detection score to decide output or not, taking up the false negatives number 0.076. Based on these analyses, we can conclude that the gap is mainly caused by the inconsistency between the scores and detection quality. By using historical information, 3D MOT can potentially provide better scores compared to single frame detectors, and this has already drawn some recent attention [3, 44].

7. Conclusions and Future Work

In this paper, we decouple the “tracking by detection” 3D MOT algorithms into several components and analyze their typical failures. With such insights, we propose corresponding enhancements of using *NMS*, *GloU*, and *Two-stage Association*, which lead to our SimpleTrack. In addition, we also rethink the frame rates and interpolation pre-processing in nuScenes. We eventually point out several possible future directions for “tracking by detection” 3D MOT.

However, beyond the “tracking by detection” paradigm, there are also branches of great potential. For better bounding box qualities, 3D MOT can refine them using long term information [25, 29, 41], which are proven to outperform the detections based only on local frames. The future work can also transfer the current manual rule-based methods into

learning-based counterparts, *e.g.* using learning based intra-frame mechanisms to replace the NMS, using inter-frame reasoning to replace the 3D GIoU and life cycle management, etc.

Acknowledgment. We would like to thank Tianwei Yin for kindly helping us during our applying the CenterPoint detection to 3D multi-object tracking.

A. Appendix for SimpleTrack

A.1. Validation Split Comparison

We compare our SimpleTrack with other 3D MOT methods on the validation splits as in Tab. A and Tab. B. In the experiments, our SimpleTrack also demonstrates strong performance. On both Tab. A and Tab. B, our SimpleTrack can outperform the methods without learning based modules, which is consistent with the test set performance in the main paper (Tab. 6 and Tab. 7). In addition, we find it interesting in Tab. B that the learning based method OGR3MOT [44] can achieve better performance than our 2Hz SimpleTrack, which demonstrate the potential of applying learning techniques for 3D MOT. However, such advantage of OGB3MOT vanishes for AMOTA on the test set, as in the Tab. 7 of the main paper. This suggests that our learning-free modifications may have the ability to adapt to the domain gaps in the data.

A.2. Experimental Setup

Due to the space constraints, we discuss the detailed hyper-parameters and settings for our SimpleTrack here.

Waymo Open Dataset

1. **Pre-process.** We use CenterPoint detection [43], and then apply NMS with the IoU threshold equals to 1/4 onto the detection bounding boxes.
2. **Association.** We use GIoU as the association metric and Hungarian algorithm to solve the matchings. The threshold for GIoU is -0.5 across all types of objects.
3. **Motion Model.** We use the default Kalman filter parameters as AB3DMOT [37], and pair the usages of Hungarian algorithm.
4. **Life Cycle Management.** The life cycle management is the same as AB3DMOT [37], 3 hits to start outputting a tracker and consecutive 2 misses terminates a tracklet. We set the threshold for outputting detection bounding boxes as 0.7 for vehicle and cyclist, and 0.5 for pedestrian. In our “Two-stage association,” we adopt the low score threshold as 0.1.

Method	Vehicle			Pedestrian		
	MOTA↑	MOTP↓	IDS(%)↓	MOTA↑	MOTP↓	IDS(%)↓
AB3DMOT* [37]	0.5572	0.1679	0.40	0.5224	0.3098	2.74
Chiu <i>et al.</i> * [10]	0.5406	0.1665	0.37	0.4810	0.3086	3.34
CenterPoint [43]	0.5505	0.1691	0.26	0.5493	0.3137	1.13
SimpleTrack	0.5612	0.1681	0.08	0.5776	0.3125	0.42

Table A. Results on WOD validation split (L2). We mark the best in **red** and the second in **blue**. For fair comparison, we list the methods using the public CenterPoint [43] detection, * means the numbers from our own implementations.

Methods	AMOTA↑	AMOTP↓	MOTA↑	IDS ↓
AB3DMOT* [37]	0.598	0.771	0.537	1570
AB3DMOT [37] [44]	0.578	0.807	0.514	1275
Chiu <i>et al.</i> * [10]	0.624	0.655	0.542	1098
Chiu <i>et al.</i> [10] [44]	0.617	0.984	0.533	680
CenterPoint [43]	0.665	0.567	0.562	562
CBMOT [3]	0.675	0.591	0.583	494
MPN-Baseline [44]	0.593	0.832	0.514	1079
OGR3MOT [44]	0.693	0.627	0.602	262
SimpleTrack (2Hz)	0.687	0.573	0.592	519
SimpleTrack (10Hz)	0.696	0.547	0.602	405

Table B. Results on nuScenes validation set. We mark the best in **red** and the second in **blue**. For fair comparison, we list the methods using public CenterPoint [43] detection, the numbers marked with * are our own implementations, the numbers marked with [44] are from OGR3MOT [44].

nuScenes

1. **Pre-process.** We apply NMS according to IoU threshold equal to 1/10. After NMS, all the remaining detections are kept as the input to 3D MOT algorithms.
2. **Association.** We adopt the same settings as on Waymo Open Dataset.
3. **Motion Model.** The settings for our motion model is identical to that on Waymo Open Dataset.
4. **Life Cycle Management.** We adopt the similar strategy as the tracking algorithm in Center Point [43], where the trackers start outputting upon the first association, and are terminated after two continuous misses.

References

- [1] Yaakov Bar-Shalom, Thomas E Fortmann, and Peter G Cable. Tracking and data association, 1990. **2**
- [2] Erkan Baser, Venkateshwaran Balasubramanian, Prarthana Bhattacharyya, and Krzysztof Czarnecki. FANTrack: 3D multi-object tracking with feature association network. In *IV*, 2019. **2**
- [3] Nuri Benbarka, Jona Schröder, and Andreas Zell. Score refinement for confidence-based 3D multi-object tracking. In *IROS*, 2021. **1, 2, 3, 6, 8, 9**
- [4] Jerome Berclaz, Francois Fleuret, Engin Turetken, and Pascal Fua. Multiple object tracking using k-shortest paths op-

- timization. *IEEE transactions on pattern analysis and machine intelligence*, 33(9):1806–1819, 2011. 2
- [5] Philipp Bergmann, Tim Meinhardt, and Laura Leal-Taixe. Tracking without bells and whistles. In *ICCV*, 2019. 2
- [6] Alex Bewley, Zongyuan Ge, Lionel Ott, Fabio Ramos, and Ben Uprocroft. Simple online and realtime tracking. In *ICIP*, 2016. 2
- [7] Guillem Braso and Laura Leal-Taixe. Learning a neural solver for multiple object tracking. In *CVPR*, 2020. 2
- [8] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuScenes: A multi-modal dataset for autonomous driving. In *CVPR*, 2020. 1, 2
- [9] Hsu-kuang Chiu, Jie Li, Rares Ambrus, and Jeannette Bohg. Probabilistic 3D multi-modal, multi-object tracking for autonomous driving. *ICRA*, 2021. 2
- [10] Hsu-kuang Chiu, Antonio Prioletti, Jie Li, and Jeannette Bohg. Probabilistic 3D multi-object tracking for autonomous driving. *arXiv:2001.05673*, 2020. 1, 2, 3, 4, 5, 6, 7, 9
- [11] Francois Fleuret, Jerome Berclaz, Richard Lengagne, and Pascal Fua. Multicamera people tracking with a probabilistic occupancy map. *IEEE transactions on pattern analysis and machine intelligence*, 30(2):267–282, 2007. 2
- [12] Jiawei He, Zehao Huang, Naiyan Wang, and Zhaoxiang Zhang. Learnable graph matching: Incorporating graph partitioning with deep feature learning for multiple object tracking. In *CVPR*, 2021. 2
- [13] Andrea Hornakova, Roberto Henschel, Bodo Rosenhahn, and Paul Swoboda. Lifted disjoint paths with application in multiple object tracking. In *ICML*, 2020. 2
- [14] Xiaolong Jiang, Peizhao Li, Yanjing Li, and Xiantong Zhen. Graph neural ased end-to-end data association framework for online multiple-object tracking. *arXiv preprint arXiv:1907.05315*, 2019. 2
- [15] Aleksandr Kim, Aljosa Osep, and Laura Leal-Taixé. EagerMOT: 3D multi-object tracking via sensor fusion. *arxiv:2104.14682*, 2021. 2
- [16] Chanh Kim, Fuxin Li, Arridhana Ciptadi, and James M Rehg. Multiple hypothesis tracking revisited. In *ICCV*, 2015. 2
- [17] Long Lan, Dacheng Tao, Chen Gong, Naiyang Guan, and Zhigang Luo. Online multi-object tracking by quadratic pseudo-boolean optimization. In *IJCAI*, 2016. 2
- [18] Laura Leal-Taixé, Cristian Canton-Ferrer, and Konrad Schindler. Learning by tracking: Siamese CNN for robust target association. In *CVPR Workshops*, 2016. 2, 8
- [19] Jiahe Li, Xu Gao, and Tingting Jiang. Graph networks for multiple object tracking. In *WACV*, 2020. 2, 8
- [20] Tianyi Liang, Long Lan, and Zhigang Luo. Enhancing the association in multi-object tracking via neighbor graph. *arXiv preprint arXiv:2007.00265*, 2020. 2
- [21] Qiankun Liu, Qi Chu, Bin Liu, and Nenghai Yu. GSM: Graph similarity model for multi-object trackin. In *IJCAI*, 2020. 2
- [22] Zhichao Lu, Vivek Rathod, Ronny Votel, and Jonathan Huang. RetinaTrack: Online single stage joint detection and tracking. In *CVPR*, 2020. 2
- [23] Py motmetrics Contributors. py-motmetrics. <https://github.com/cheind/py-motmetrics>. 5
- [24] Jiangmiao Pang, Linlu Qiu, Xia Li, Haofeng Chen, Qi Li, Trevor Darrell, and Fisher Yu. Quasi-dense similarity learning for multiple object tracking. In *CVPR*, 2021. 2
- [25] Ziqi Pang, Zhichao Li, and Naiyan Wang. Model-free vehicle tracking and state estimation in point cloud sequences. *IROS*, 2021. 8
- [26] 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
- [27] Hamed Pirsiavash, Deva Ramanan, and Charless C Fowlkes. Globally-optimal greedy algorithms for tracking a variable number of objects. In *CVPR*, 2011. 2
- [28] Johannes Pöschmann, Tim Pfeifer, and Peter Protzel. Factor graph based 3D multi-object tracking in point clouds. In *IROS*, 2020. 1, 2
- [29] Charles R Qi, Yin Zhou, Mahyar Najibi, Pei Sun, Khoa Vo, Boyang Deng, and Dragomir Anguelov. Offboard 3D object detection from point cloud sequences. *CVPR*, 2021. 8
- [30] Donald Reid. An algorithm for tracking multiple targets. *IEEE transactions on Automatic Control*, 24(6):843–854, 1979. 2
- [31] Hamid Rezaatofghi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union. *CVPR*, 2019. 4
- [32] Seyed Hamid Rezaatofghi, Anton Milan, Zhen Zhang, Qinfeng Shi, Anthony Dick, and Ian Reid. Joint probabilistic data association revisited. In *ICCV*, 2015. 2
- [33] Amir Sadeghian, Alexandre Alahi, and Silvio Savarese. Tracking the untrackable: Learning to track multiple cues with long-term dependencies. In *ICCV*, 2017. 2, 8
- [34] Pei Sun, Henrik Kretschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott Ettinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. Scalability in perception for autonomous driving: Waymo Open Dataset. *arxiv:1912.04838*, 2019. 1
- [35] Siyu Tang, Bjoern Andres, Miykhaylo Andriluka, and Bernt Schiele. Subgraph decomposition for multi-target tracking. In *CVPR*, 2015. 2
- [36] Siyu Tang, Bjoern Andres, Mykhaylo Andriluka, and Bernt Schiele. Multi-person tracking by multicut and deep matching. In *ECCV*, 2016. 2
- [37] Xinshuo Weng, Jianren Wang, David Held, and Kris Kitani. 3D multi-object tracking: A baseline and new evaluation metrics. *IROS*, 2020. 1, 2, 3, 4, 5, 6, 7, 8, 9
- [38] Xinshuo Weng, Yongxin Wang, Yunze Man, and Kris Kitani. GNN3DMOT: Graph neural network for 3D multi-object tracking with 2D-3D multi-feature learning. *CVPR*, 2020. 2
- [39] Nicolai Wojke, Alex Bewley, and Dietrich Paulus. Simple online and realtime tracking with a deep association metric. In *ICIP*, 2017. 2, 8

- [40] Yihong Xu, Aljosa Osep, Yutong Ban, Radu Horaud, Laura Leal-Taixé, and Xavier Alameda-Pineda. How to train your deep multi-object tracker. In *CVPR*, 2020. 2
- [41] Bin Yang, Min Bai, Ming Liang, Wenyuan Zeng, and Raquel Urtasun. Auto4D: Learning to label 4D objects from sequential point clouds. *arXiv:2101.06586*, 2021. 8
- [42] Bo Yang, Chang Huang, and Ram Nevatia. Learning affinities and dependencies for multi-target tracking using a CRF model. In *CVPR*, 2011. 2
- [43] Tianwei Yin, Xingyi Zhou, and Philipp Krähenbühl. Center-based 3D object detection and tracking. *CVPR*, 2021. 1, 2, 3, 4, 5, 6, 8, 9
- [44] Jan-Nico Zaech, Dengxin Dai, Alexander Liniger, Martin Danelljan, and Luc Van Gool. Learnable online graph representations for 3D multi-object tracking. *arXiv:2104.11747*, 2021. 1, 2, 6, 8, 9
- [45] Amir Roshan Zamir, Afshin Dehghan, and Mubarak Shah. GMCP-Tracker: Global multi-object tracking using generalized minimum clique graphs. In *ECCV*, 2012. 2
- [46] Li Zhang, Yuan Li, and Ramakant Nevatia. Global data association for multi-object tracking using network flows. In *CVPR*, 2008. 2
- [47] Yifu Zhang, Peize Sun, Yi Jiang, Dongdong Yu, Zehuan Yuan, Ping Luo, Wenyu Liu, and Xinggang Wang. ByteTrack: Multi-object tracking by associating every detection box. *arXiv preprint arXiv:2110.06864*, 2021. 6
- [48] Yifu Zhang, Chunyu Wang, Xinggang Wang, Wenjun Zeng, and Wenyu Liu. FairMOT: On the fairness of detection and re-identification in multiple object tracking. *International Journal of Computer Vision*, pages 1–19, 2021. 2