

---

# TSTTC: A Large-Scale Dataset for Time-to-Contact Estimation in Driving Scenarios

---

**Yuheng Shi**  
Tianjin University  
yuheng@tju.edu.cn

**Zehao Huang**  
TuSimple  
zehoahuang18@gmail.com

**Yan Yan**  
TuSimple  
yanyan.paper@outlook.com

**Naiyan Wang**  
TuSimple  
winsty@gmail.com

**Xiaojie Guo**  
Tianjin University  
xj.max.guo@gmail.com

## Abstract

1 Time-to-Contact (TTC) estimation is a critical task for assessing collision risk and is  
2 widely used in various driver assistance and autonomous driving systems. The past  
3 few decades have witnessed development of related theories and algorithms. The  
4 prevalent learning-based methods call for a large-scale TTC dataset in real-world  
5 scenarios. In this work, we present a large-scale object oriented TTC dataset in the  
6 driving scene for promoting the TTC estimation by a monocular camera. To collect  
7 valuable samples and make data with different TTC values relatively balanced, we  
8 go through thousands of hours of driving data and select over 200K sequences with  
9 a preset data distribution. To augment the quantity of small TTC cases, we also  
10 generate clips using the latest Neural rendering methods. Additionally, we provide  
11 several simple yet effective TTC estimation baselines and evaluate them extensively  
12 on the proposed dataset to demonstrate their effectiveness. The proposed dataset  
13 and code are publicly available at dataset link and code link.

## 14 1 Introduction

15 In recent years, there has been a growing trend towards equipping vehicles with Advanced Driver  
16 Assistance System (ADAS), which consists of several subsystems such as Adaptive Cruise Control  
17 (ACC), Automated Emergency Braking (AEB) and Forward Collision Warning (FCW). ADAS  
18 aims to detect potential hazards as quickly as possible and alert the driver, or take corrective action  
19 to improve driving safety. AEB and FCW are critical features of ADAS that protect drivers and  
20 passengers and prevent traffic accidents. They both rely on the estimation of Time-to-Contact (TTC)  
21 which is defined as the time for an object to collide with the observer's plane. Although TTC can be  
22 predicted using data from various sensors, such as LiDAR, radar or camera. Vision-based methods are  
23 particularly attractive due to their low-cost and have been a popular choice among ADAS designers  
24 and manufacturers. Even in high-level (L3+) autonomous driving system, direct TTC estimation  
25 could also serve as an redundant observation when other distance measuring sensors fail.

26 Prior to the widespread adoption of deep learning, numerous vision-based theories and algorithms [25,  
27 29, 10, 6, 4] for estimating TTC had been proposed. These algorithms are not data-driven, and usually  
28 rely on hand-crafted cues. Recently, several deep learning based TTC estimation algorithms have  
29 emerged [1, 45] and demonstrate promising results in driving scenarios. The emergence of deep

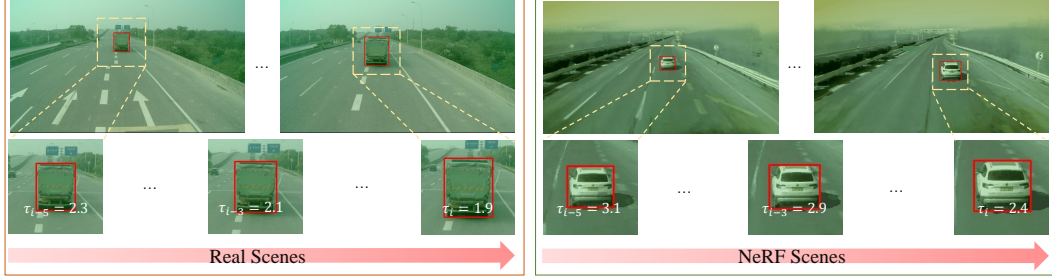


Figure 1: Example sequences and annotations from our dataset. The  $\tau$  denotes the TTC ground-truth while the subscript denotes the frame index. We could observe that, the scale of the object increases as the TTC decreases.

30 learning has brought more powerful tools to computer vision and also brought higher demands for  
 31 large-scale datasets. However, due to the lack of large-scale TTC datasets that capture real-world  
 32 driving scenarios, these methods have to pre-train their model on synthetic flow datasets [27, 12, 22].

33 In this paper, we primarily address the challenge of TTC estimation in highway scenarios. Compared to  
 34 urban scenarios, high speed driving in highway exhibits longer braking distances, thereby  
 35 necessitating a broader range of perception capabilities. In order to facilitate the development of  
 36 vision-based TTC estimation algorithms, we propose a large-scale monocular TTC dataset in this  
 37 paper, using a class 8 heavy truck as the data collection platform. From raw data collected in urban  
 38 and highway scenarios, we identify over 200K sequences covering a depth range of 400 meters.  
 39 Each sequence contains six consecutive frames captured at a rate of 10Hz, with 2D, 3D bounding  
 40 box and TTC ground-truth provided for a single object in each frame. Additionally, to address the  
 41 limited availability of samples in rare scenarios, such as sudden braking (e.g. small TTC cases), we  
 42 utilize Neural Radiance Fields (NeRF) [30] to generate additional data. These artificially generated  
 43 data can be seamlessly integrated into our dataset, thereby increasing the quantity and diversity of  
 44 data available for training. Fig.1 illustrates two typical examples from our dataset: vehicles are  
 45 gradually approaching in real scenes and NeRF scenes, respectively. Besides the proposed dataset,  
 46 we focus on object level TTC estimation rather than pixel level TTC estimation in [1]. Specifically,  
 47 we provide a sequence of images for a certain object and ask for estimating the TTC value of it in the  
 48 last frame. 2D Bounding boxes are available as optional inputs. Then we provide simple yet effective  
 49 baselines based on the relationship between the scale ratio of objects in adjacent frames and TTC. We  
 50 reformulate the problem as choosing the scale with the highest similarity in adjacent frames. Inspired  
 51 by recent studies [1, 2], we further transform scale estimation from a regression problem into a set of  
 52 binary classification tasks. A series of quantitative experiments are conducted to demonstrate the  
 53 effectiveness and feasibility of our proposed techniques. Our main contribution can be summarized  
 54 as follows:

- 55 • We propose a large-scale monocular TTC dataset for driving scenarios and will make it  
 56 publicly available along with relevant toolkits to facilitate the development of TTC estimation  
 57 algorithms for driving scenes.
- 58 • We propose two simple yet effective TTC estimation algorithms and extensively test them  
 59 on the dataset to validate their effectiveness, which could serve as baseline methods for  
 60 future study.

## 61 2 Related work

62 **Task and Methods.** In the scheme of monocular TTC estimation, TTC describes the time that an  
 63 object will cross the camera plane under concurrent relative velocity. Denote the depth of an object  
 64 in the camera coordinate as  $y$ , the time for the object under the current velocity to cross the camera  
 65 plane could be calculated by:

$$\tau = -y / \frac{dy}{dt} = -y / \dot{y}, \quad (1)$$

66 where  $j$  is the relative velocity between the object and the camera. Though estimating either velocity  
 67 or depth is an ill-posed problem, TTC can be estimated from images directly because it only depends  
 68 on the ratio of them. Researchers have proposed various approaches to accomplish TTC estimation.

69 A viable approach is to utilize hand-crafted features such as closed contours, optical flow, brightness,  
 70 or intensity from images [10, 8, 19, 29, 37, 43, 20]. Mobileye [10] adopted geometric information of  
 71 the vehicles in image to estimate TTC by establishing the relationship between TTC and the width of  
 72 vehicle. In addition to geometric-based methods, several studies have been proposed to address the  
 73 task of TTC estimation using photometric-based features, without relying on geometric features or  
 74 high-level processing. For instance, [19] adopted accumulated sums of suitable products of image  
 75 brightness derivatives from time varying images to determine the TTC value. Furthermore, [43]  
 76 elucidates the relationship between TTC and the changes in intensity in images. However, these hand-  
 77 crafted features need carefully tuned parameters or strong priors, such as constant brightness [19] or  
 78 static scene [20], which restricts their practical applicability.

79 Besides to hand-crafted methods, deep-learning approaches can also be employed for TTC estimation.  
 80 One alternative approach is to use scene flow estimation methods [28, 40, 35, 21, 45] that predict  
 81 both depth and velocity simultaneously, enabling the generation of pixel-level TTC estimation maps.  
 82 However, these methods depend on accurate optical flow information, which can result in significant  
 83 computational overhead. Recently, [1] proposed Binary TTC that bypasses optical flow computation  
 84 and directly computes TTC via estimating scale ratio between consecutive images. While these  
 85 learning-based methods may produce more promising results, they require a larger amount of data  
 86 with annotated scene flow ground-truth. Due to the expensive labeling cost, scene flow datasets are  
 87 mostly obtained through synthesis, leading to a domain gap for real-world applications.

88 In addition to the aforementioned literature, single object tracking (SOT) [3, 23, 41, 9, 46] can also  
 89 infer the scale ratio between the template and the tracked object, serving as an alternative approach to  
 90 estimating TTC. However, these methods often rely on large downsampling rates, which may lead to  
 91 the bounding box estimation not accurate enough to meet the requirements of TTC estimation. The  
 92 baseline method which utilized SOT models in our experiment validates the similar effect.

93 **Datasets.** Contrary to previous TTC esti-  
 94 mation studies, our proposed dataset facil-  
 95 itates the expansion of hand-craft features  
 96 from solely relying on RGB images to uti-  
 97 lizing the features extracted by neural net-  
 98 works. Moreover, we propose a method to  
 99 address the problem of estimating the TTC  
 100 by classifying the scale ratio. And we ext-  
 101 end the implementation on RGB images to  
 102 feature maps extracted using deep learning  
 103 models, thereby significantly enhancing the  
 104 accuracy of TTC estimation.

105 For TTC estimation, several datasets col-  
 106 lected in real scenes are available. For ex-  
 107 ample, [13] proposed a multi-person track-  
 108 ing dataset with stereo depth details, and

109 [26] presented a large-scale dataset for TTC estimation in indoor scenes. However, these datasets  
 110 mainly focus on low speed scenarios and are not suitable for direct application to TTC estimation in  
 111 driving scenarios. In addition to datasets specifically designed for TTC estimation, some datasets have  
 112 been synthesized for scene flow tasks and can provide detailed depth information, which are suitable  
 113 for TTC estimation. For example, KITTI scene flow [28] proposed an outdoor scene flow dataset  
 114 containing 400 dynamic scenes collected from KITTI [17]. These scenes are annotated using 3D  
 115 CAD models for vehicles in motion and manually mask non-rigid moving objects. The Driving [27]  
 116 dataset proposed a synthetic stereo video dataset rendering in realistic style. However, these datasets

Table 1: Comparison with several autonomous vehicle (AV) datasets. "2D-to-3D" indicates the presence of tightly bounded 2D box annotations with corresponding 3D bounding boxes, which is crucial for TTC estimation. U and H in scenes indicate Urban and Highway respectively. And we report the average velocity of the recording platform in the training set for comparison.

	KITTI	NuScenes	Waymo	Ours
Scenes	U	U	U	<b>U+H</b>
Frequency ( $hz$ )	10	2	10	10
Range ( $m$ )	[0,125]	[-80,80]	[-75,75]	<b>[-160,400]</b>
2D-to-3D	✓	✗	✗	✓
Avg Speed ( $m/s$ )	-	5.1	6.9	19.1
Ann. Frames	15K	40K	200K	<b>1M</b>
Boxes	200K	1.4M	<b>12M</b>	1M

117 are constrained by synthetic and limited scenes, which result in domain gaps with real scenes. Except  
118 to the scene flow datasets, some RGBD datasets [33, 34, 39], equipped with comprehensive depth  
119 annotations, can be utilized to train depth estimation models, which in turn can be used for TTC  
120 estimation. However, the majority depth annotations in these datasets are typically confined to a  
121 relatively small range (less than 50 meters) and the number of scenes available is limited. In addition  
122 to the aforementioned datasets, several large-scale datasets proposed for autonomous driving, such  
123 as [5, 7, 26, 38], offering comprehensive annotations like 3D LiDAR bounding boxes that can be  
124 used to generate TTC ground-truth. Nevertheless, these datasets were not specifically tailored for  
125 TTC estimation and presented issues like unbalanced TTC distribution. Furthermore, these datasets  
126 were mainly collected from urban areas, lacking data for highway scenarios where the estimation of  
127 TTC is particularly important. Please refer to Table 1 for a comparison of various datasets. Compared  
128 to the aforementioned datasets, our proposed dataset holds the distinct advantage of being large-scale  
129 and recorded in real scenarios, encompassing both urban and highway scenes.

### 130 3 Discussion

131 A common question related to TTC is that is TTC only applicable for low level assistant driving  
132 system? Why do we need TTC when we have distance and velocity measurement in advanced  
133 assistant or autonomous driving system? We argue that there are two main reasons for the essence  
134 of TTC: First, among all the three physical properties of one object (distance, velocity and time to  
135 contact), TTC is the only direct observation from monocular image. Monocular depth estimations  
136 are mostly from fully data-driven aspect, which highly relies on the recognition of the objects and  
137 scenes, which suffers from out of domain issue seriously [11]. For velocity, the situation is similar  
138 to that of depth estimation. However, TTC estimation does not require the recognition of semantic  
139 of objects (can even be achieved by pixel photometric loss), thus has better generalization in corner  
140 cases. Second, for a high level autonomous driving system, redundancy design is indispensable. Even  
141 though we have distance and velocity observation, TTC can also be fused into these observations in  
142 subsequent perception fusion modules, which serves as another independent safety guarantee.

## 143 4 TTC Dataset

### 144 4.1 Data Collection

145 The dataset consists of two parts, including a majority of data from real scenes and a minority of data  
146 from NeRF virtual rendering. After obtaining raw sensor data, we consider a sequence that contains a  
147 number of consecutive frames of a vehicle as a sample. To gather valuable data, we discard truncated  
148 objects within the camera plane and filter out 2D boxes smaller than  $15 \times 15$  pixels.

149 **Real Scenes.** For real-world scenarios, raw data was collected by our commercial trucks. We capture  
150 image data using three frontal and two backwards cameras with same  $1024 \times 576$  resolution. We  
151 adopt 2D object detection and tracking algorithms on the images to obtain 2D bounding boxes and  
152 corresponding track ID for other vehicles. And the LiDAR and Radar data are adopted to generate  
153 accurate depth and velocity of them. The sampling rate is 10Hz. Detailed sensor specification can be  
154 found in the appendix. After applying these rules to filter the raw data, we observe that the resulting  
155 data distribution is highly imbalanced across different TTC ranges. For example, vehicles with  
156 small TTC are extremely rare in driving scenes, especially for vehicles lie in the same lanes as ego  
157 vehicle. However these are the cases which FCW and AEB should focus on. In such conditions,  
158 data collection without rebalancing may result in lack of these valuable scenarios. To overcome this  
159 challenge, we pre-define a data distribution and sample the data accordingly. The sampling weights  
160 for the TTC intervals, specifically (0,3], (3,6], (6,10], (10,15], and (15,20], is set to 14%. For the  
161 TTC range of [-20,0), the sampling weight is designated at 30%.

162 **NeRF Scenes.** Despite thorough data collection efforts, we discover samples with small TTC values  
163 within the same lane are still extremely scarce, which is a crucial scenario in automated driving and

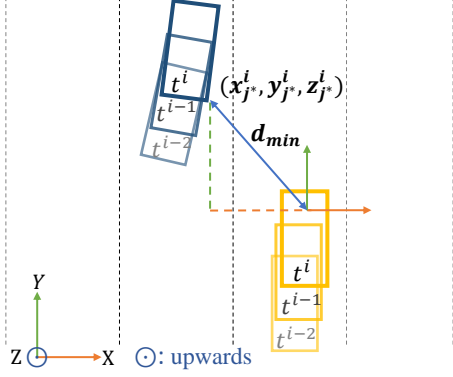


Figure 2: Relative position between ego vehicle and an object in bird-eye-view. Only three frames are plotted for brevity.

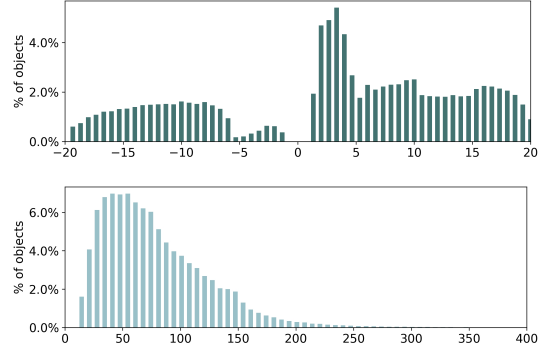


Figure 3: Histogram of TTC GT and relative depth of the training set.

164 ADAS. To supplement the absence of data in particular conditions, we adopt an internal undisclosed  
 165 project which is developed based on Instant-NGP [31] to render novel scenes. Briefly speaking, the  
 166 background models and object models are firstly extracted and trained separately. And then we can  
 167 form a new scene and render them together. Given a specific object, we pre-define some scripts in  
 168 which the TTC of the object is distributed between 0 and 6. The preset script can be found in our  
 169 appendix. After obtaining NeRF rendered images, we organize them with the same format as real  
 170 scenes data, which serves as an optional component within the training set.

## 171 4.2 Annotation

172 In each sequence, we provide the TTC ground-truth for every frame as the annotation. In the  
 173 following, we will describe how we generate the TTC annotation. Given a frame, we first run 2D  
 174 detection on the image and 3D detection on LiDAR. The long range LiDAR detection algorithm  
 175 which reliably covers [-160, 400] meter range. Then, we could obtain its corresponding 2D detection  
 176 box by projecting the 3D box to the image plane then picking the 2D detection box which has highest  
 177 IoU between the projected box. In the vehicle coordinate system, one corner of the 3D bounding  
 178 box could be denoted as  $x_j, y_j, z_j$  where  $j \in \{1, 2, \dots, 8\}$  is the corner index. Here, we take the  
 179 y-coordinate of the corner point which is the closest to ego as the depth of this object:

$$j^* = \arg \min_{j \in \{1, 2, \dots, 8\}} (\sqrt{(x_j^i - 0)^2 + (y_j^i - 0)^2 + (z_j^i - 0)^2}), \text{ where } y^i = y_{j^*}^i, \quad (2)$$

180 where  $y^i$  is the depth of the vehicle in  $i$ -th frame. Here, we assume the relative velocity between the  
 181 vehicle and ego is constant in a short time interval (e.g.  $q$  frames) to acquire a stable estimation of the  
 182 velocity. Given the depth of the vehicle in the past  $q$  frames, we fit the depth to obtain the relative  
 183 velocity  $v^i$  by RANSAC [15] algorithm. We set the value of  $q$  to 10 by default. It is worth noting  
 184 that the velocity is acquired prior to sequence splitting, thereby  $q$  may be greater than the sequence  
 185 length. After obtaining the depth and relative velocity of current frame, the TTC ground-truth could  
 186 be obtained by  $\tau^i = y^i / v^i$ .

187 Since the camera planes are almost parallel to the XZ plane in the vehicle coordinate and the origins  
 188 of these coordinate systems are highly aligned, we regard the TTC value calculated in the vehicle  
 189 coordinate system as the TTC ground-truth for camera planes. The annotation process is illustrated in  
 190 Fig 2. One may challenge that using past 10 frames to fit the velocity may result in latency in velocity  
 191 estimation especially for sudden break. To remedy this issue, we first generate TTC ground-truth  
 192 with different  $q$  (e.g. 3, 5 and 10) and consider the vehicle with varied TTC values as an accelerating  
 193 or decelerating object. For object with varied TTC values, we select the one closest to the TTC  
 194 computed by the velocity of radar sensor as the ground truth. Note that we do not directly utilize  
 195 radar sensors for all objects since they could not cover too distant objects. **The data annotations are**  
 196 **then manually checked by our annotation team to ensure the quality**

### 197 4.3 Dataset Statistics

198 We construct the dataset following the pre-defined rules and annotation pipeline, resulting in 206K  
199 sequences comprising over 1M frames from real scenes, as well as 1K sequences from 6.0K NeRF  
200 rendered images. We split the sequences from the real scenes based on their recorded date to training,  
201 validation and test sets, yielding 149.1K, 28.8K, and 28.6K sequences respectively. For better  
202 understanding the data distribution, we plot the histogram of the TTC ground-truth and depth in the  
203 training set in Fig. 3. The distribution in validation and test set is similar. Note that the far away  
204 samples are rare because we set a minimum 2D bounding box size. More detailed information about  
205 the dataset statistics is provided in our appendix.

### 206 4.4 Task Definition

207 As the tracklet of a vehicle is collected in the format of fixed length sequences, we formulate the TTC  
208 estimation task in sequence level. For a sequence of a specific vehicle, we provide six consecutive  
209 frames and their corresponding 2D bounding boxes as input. The last frame in the sequence is  
210 considered as the target frame while the rest of the frames serve as references. With all frames and  
211 boxes available in the sequence, the objective is to predict the TTC value of the object in the target  
212 frame or equivalently relative scale ratio between the target box and the referenced one.

## 213 5 Metrics & Method

214 In this section, we first review the relationship between TTC estimation and scale ratio briefly. Then,  
215 we explicate the evaluation metrics for our TTC estimation task. Subsequently, we introduce our  
216 approach, which comprises two variants: the pixel MSE approach and its deep learning counterpart.

### 217 5.1 Estimate TTC via Scale Ratio

218 As shown in Sec. 4, TTC could be obtained by depth and its rate of change. However, estimating the  
219 depth and relative velocity of an object directly with only a monocular camera is very challenging. To  
220 address this issue, researchers proposed to estimate the TTC of frontal-parallel, planar non-deformable  
221 objects according to the change of object scales. As described in [19], we can obtain the image size  
222 of a frontal-parallel, non-deformable object of size  $S$  at distance  $y$  as:

$$s = fS/y, \tag{3}$$

223 where  $f$  is the focal length of the camera. For an object without rotation, TTC can be formulated as a  
224 function of object size in image space by combining Eq. (1) and (3):

$$\tau = \frac{t_1 - t_0}{1 - \frac{s(t_0)}{s(t_1)}} = \frac{t_1 - t_0}{1 - \alpha}, \tag{4}$$

225 where  $s(t_0)$  and  $s(t_1)$  are the sizes of image of an object in frame  $t_0$  and  $t_1$  correspondingly. Thus,  
226 the estimation of TTC can be simplified as a scale ratio estimation problem which can be done with  
227 only observations in image space. With the development of deep learning, modern object detectors or  
228 trackers could produce relatively accurate 2D bounding boxes for vehicles. Given the detection or  
229 tracking bounding box in consecutive frames of a vehicle, one intuitive idea is that we can use the  
230 ratio of the box or mask area to accomplish the scale ratio estimation. However, are these bounding  
231 boxes accurate enough for the TTC estimation task and does there exist more accurate scale ratio  
232 estimation algorithms under such conditions? We try to answer these questions via experimental  
233 validation on the proposed dataset.

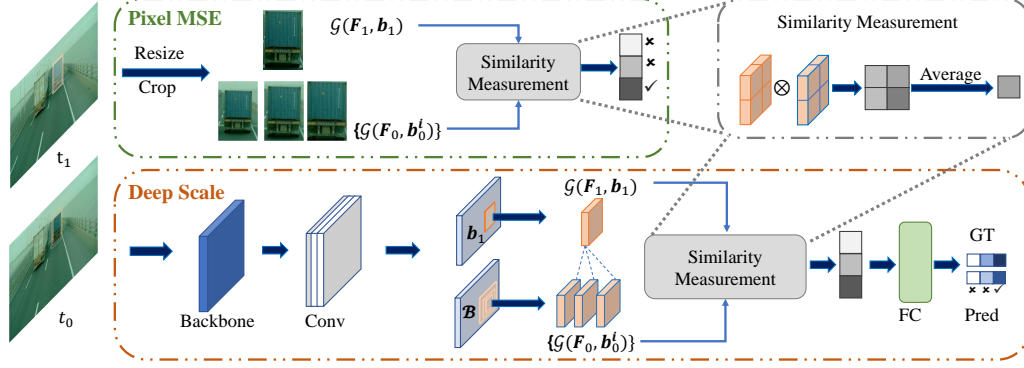


Figure 4: Framework of our proposed methods. After aligning the size between contents in  $\mathbf{b}_1$  and  $\mathcal{B}$ , an operator  $\otimes$  is applied to measure their similarity. For simplicity, we only illustrate three scale rates of  $\mathcal{B}$ . Some operations such as center shift are omitted for brevity. The green dashed box and the orange dashed box represent **Pixel MSE** and **Deep Scale**, respectively.

## 234 5.2 Evaluation Metrics

235 Before introducing the detailed design, we need to design the metrics for evaluation. Here, we adopt  
 236 Motion-in-Depth (MiD) error and Relative TTC Error (RTE) as performance indicators, which could  
 237 be denoted as:

$$\begin{aligned} \text{MiD} &= \|\log(\alpha) - \log(\hat{\alpha})\|_1 \times 10^4, \\ \text{RTE} &= \left\| \frac{\tau - \hat{\tau}}{\tau} \right\|_1 \times 100\%, \end{aligned} \quad (5)$$

238 where  $\alpha$  and  $\hat{\alpha}$  mean the ground-truth and predicted scale ratios while the  $\hat{\tau}$  denotes predicted TTC  
 239 value. The scale ratio ground-truth  $\alpha$  is obtained by Eq. (4). MiD is utilized in previous works [45, 1]  
 240 to describe the TTC error indirectly from the perspective of scale ratio. Due to the instability of  
 241 TTC at larger values, **we prioritize the MiD as the primary metric**. With the purpose of revealing  
 242 more information from the evaluation metrics, we partition several TTC intervals, namely crucial(c) /  
 243 small(s) / large(l) / negative(n)<sup>1</sup>, which correspond to TTC values of 0~3, 3~6, 6~20, and -20~0,  
 244 respectively. We mark them on the indices of the RTE and MiD. The threshold for crucial cases is  
 245 determined by rounding the typical TTC threshold of 2.7 seconds used in FCW systems [36, 48].  
 246 This division allows a more detailed analysis of the performance for the TTC estimation algorithms  
 247 in different TTC intervals, providing a better understanding of their limitations and strengths. During  
 248 the prediction, TTC values that exceed the predefined range will be truncated to the boundary value.

## 249 5.3 Our Design

250 **Formulation.** We denote the 2D bounding box as  $\mathbf{b} = [x, y, w, h]$  where  $x, y$  represent the center  
 251 coordinate and  $w, h$  denote the width and height. Given a reference frame  $\mathbf{F}_0$  and a target frame  
 252  $\mathbf{F}_1$ ,  $\mathbf{b}_0 = [x_0, y_0, w_0, h_0]$  and  $\mathbf{b}_1 = [x_1, y_1, w_1, h_1]$  are bounding boxes of a specific object in these  
 253 two frames. The core idea of our methods is to estimate the relative scale ratio of this vehicle  
 254 between the reference frame and the target frame. A straightforward approach is simply adopting  
 255 the scale ratio between  $\mathbf{b}_0$  and  $\mathbf{b}_1$  as the result. However, this strategy is largely influenced by the  
 256 precision of detection algorithms. In our methods, we estimate the scale ratio change in these two  
 257 frames in pixel space or feature space, yielding two kinds of implementation: Pixel MSE and Deep  
 258 Scale. For the reference frame, we enumerate  $n$  different scale ratios to obtain a series of scaled  
 259 boxes  $\mathcal{B} = [\mathbf{b}_0^{\alpha_1}, \mathbf{b}_0^{\alpha_2}, \dots, \mathbf{b}_0^{\alpha_i}, \dots, \mathbf{b}_0^{\alpha_n}]$  where  $\mathbf{b}_0^{\alpha_i} = [x_0, y_0, \alpha_i w_1, \alpha_i h_1]$ . Then, we crop  $\mathbf{F}_0$  via  $\mathcal{B}$   
 260 and resize them to a target size of  $W, H$ , which could be denoted as  $\mathcal{G}(\mathbf{F}_0, \mathbf{b}_0^{\alpha_i})$  for the  $i$ -th scale  
 261 ratio, where  $\mathcal{G}(\mathbf{F}, \mathbf{b})$  denotes the crop  $\mathbf{b}$  on  $\mathbf{F}$  and resize it. Similarly, we could get  $\mathcal{G}(\mathbf{F}_1, \mathbf{b}_1)$   
 262 for the target frame. Finally, we use an operator to measure the similarity between  $\mathcal{G}(\mathbf{F}_0, \mathbf{b}_0^{\alpha_i})$  and

<sup>1</sup>Negative TTC indicates away from the observer.

263  $\mathcal{G}(\mathbf{F}_1, \mathbf{b}_1)$ , yielding  $n$  similarity scores. With five frames free to reference, taking different frames  
264 as the reference will produce different scale ratios. To address this issue, we convert scale ratio to  
265 corresponding TTC value via Eq. (4) and convert it to the scale ratio under the setting of 10Hz when  
266 computing MiD. We list the relationship between different scale ratios of same  $\tau$  in the appendix.

267 **Pixel MSE.** We can measure the similarity between  $\mathcal{G}(\mathbf{F}_0, \mathbf{b}_0^{\alpha_i})$  and  $\mathcal{G}(\mathbf{F}_1, \mathbf{b}_1)$  in image space with  
268 Mean Squared Error (MSE). The weighted sum of the top  $k$  similar scale ratios is adopted as the final  
269 estimation and the weight is normalized by the reciprocal of MSE. The top of Fig. 4 illustrates the  
270 pipeline of Pixel MSE.

271 **Deep Scale.** For the deep version, we first input two images into a backbone network for feature  
272 extraction. Afterwards, the grid sampling operation is used to align the features of different box sizes  
273 into one fixed size. Then, we calculate the similarity of each position in the two feature maps via  
274 cosine similarity, yielding a similarity map  $\mathbf{S}_i$ . Afterwards, the similarity score of scale  $\alpha_i$  is obtained  
275 by adopting a Global Average Pooling (GAP) operation to  $\mathbf{S}_i$ . Then we concatenate the similarity  
276 scores of different scale ratios and use a Fully-Connected (FC) layer to obtain the final prediction.  
277 During training, binary cross-entropy (BCE) loss is used and we adopt the strategy proposed in  
278 [24, 47] to convert the scale ratio ground-truth to a size  $n$  vector as soft label. Similar to Pixel  
279 MSE, we apply a top  $k$  weighted sum operation to get the final results and the weight is defined  
280 as the sigmoid of the FC output. For fast inference, we adopt a convolutional layer followed by a  
281 stage of modified CSPNet [42] used in [16] as our backbone. After obtaining backbone features,  
282 we fed them into one transposed convolutional layer and two convolutional layers before similarity  
283 measurement. To capture subtle details, all the stride in the network is set to 1 except the first and the  
284 transposed convolutional layer which are set to 2 yielding a feature map with the same size of input.  
285 The framework is illustrated at the bottom of Fig. 4.

286 **Center Shift.** The boxes predicted by the object detection model may be inaccurate, which will  
287 bring misalignment between the centers of reference and target boxes. To address this problem, we  
288 introduce a center shift operation. Specifically, we enumerate an offset of  $[-c, c]$  along height and  
289 width direction, respectively, which yields a total of  $(2c + 1) \times (2c + 1)$  candidates. After obtaining  
290  $(2c + 1) \times (2c + 1)$  similarity scores for a single scale, we adopt the highest score as the final score  
291 for both Pixel MSE and Deep Scale. Our experiments show that this operation will bring significant  
292 improvement in terms of MiD and RTE with little time cost.

## 293 6 Experimental Validation

### 294 6.1 Implementation Details

295 **Pixel MSE.** For Pixel MSE, we validate its performance on validation, and test set. The target size  
296  $W, H$  after interpolation is set to the size of  $\mathbf{b}_1$ . The scale ratio is set to the range of  $[0.65, 1.5]$  to  
297 cover samples with different scale ratios in the training set. The number of scale bins  $n$  is 125, the  
298 top  $k$  for the weighted sum and the  $c$  for center shift are 3. Besides, the detection boxes are manually  
299 expanded with a maximum ratio of 1.1 (if the expanded boxes do not exceed the image boundary) to  
300 reduce the influence of inaccurate detection results.

301 **Deep Scale.** In terms of Deep Scale, we train it for 36 epochs on the train/train+val set dataset for  
302 evaluation on val/test respectively with a batch size of 16 using SGD [18]. The image is resized  
303 to  $1024 \times 576$  for both training and test phases. We adopt random color on HSV space as data  
304 augmentation during training. The weight decay and SGD momentum parameters are set to 0.0005  
305 and 0.9, respectively. We start from a learning rate of  $10^{-4} \times \text{BatchSize}$  and adopt cosine learning  
306 rate schedule. The target size is set to  $50 \times 50$  for grid sample as larger size does not bring more  
307 benefits. The input channel for the backbone is set to 12, while the channel for the three followed  
308 convolutional layers is set to 24. The kernel size of the convolutional layers is 7 while the transposed  
309 one is 3. For the scale range and box expansion, we keep them the same as Pixel MSE. Besides, the



number of scale bins  $n$ , the top  $k$  for the weighted sum and the  $c$  for center shift are set to 20, 4 and 1 respectively. We test the latency of all models with FP16 and batch size of 1 on a 3090 GPU.

Besides the aforementioned methods, we further propose two baselines termed as Detection and SOT in Table 4. For Detection, the scale ratio is obtained by simply computing the ratio between the area of the target box and the reference box. As for SOT, given a target box, we adopt a state-of-the-art (SOTA) SOT tracker [46] to obtain a reference box and then estimate the scale ratio as the same as in Detection. Additionally, we include results from an internal monocular depth algorithm and a LiDAR detection + tracking algorithm for a comprehensive evaluation. Details of these algorithms are provided in our appendix. Although there are other methods available, such as [45, 1], the models released by the authors achieve poor results in our dataset due to the domain gap, and no training codes are provided.

Table 2: Main results of different methods on the validation set. The † means the result is obtained under padding NeRF data. The % after RTE is omitted for brevity.

Methods	MiD	MiD <sub>c</sub>	MiD <sub>s</sub>	MiD <sub>l</sub>	MiD <sub>n</sub>	RTE
Detection	213.9	675.4	305.1	112.4	115.3	58.1
SOT [46]	200.8	641.1	261.1	77.4	158.8	57.1
Pixel MSE	41.0	57.4	36.5	32.5	48.4	29.9
Depth	62.3	111.9	74.6	36.1	68.4	47.3
LiDAR	6.4	12.5	6.0	5.3	3.3	-
Deep Scale	14.4	27.1	16.4	10.9	13.5	12.1
Deep Scale <sup>†</sup>	14.3	26.5	15.2	10.8	13.5	12.0

## 6.2 Main Results

The detailed comparison between various methods is presented in Tab. 4. Due to page constraints, we report only their performance on the MiD metric and overall RTE. As observed, the RTE predicted by box detection or tracking methods is approximately 50%, which is inadequate for practical applications. In contrast, our Pixel MSE method demonstrates significantly lower MiD errors, indicating more accurate scale ratio estimations and consequently, lower RTEs. Among learning-based methods, our Deep Scale significantly outperforms the depth estimation method, although it remains inferior to the LiDAR detection + tracking algorithm. Nevertheless, it achieves the best performance among methods that utilize images. **More detailed comparison, ablations and visualizations could be found in the appendix.**

## 7 Limitations

Our work is a large-scale benchmark for Time-To-Contact estimation, but there are still some limitations that need to be addressed in future works. In regard to our dataset, the collected data primarily focuses on trucks and cars in highway and urban scenarios, lacking more diverse categories that are commonly found in autonomous driving datasets, such as cyclists and pedestrians. Besides, we have to acknowledge that due to the inherent differences in distribution between our dataset and real-world scenarios, there may be potential challenges when directly applying the model trained on our dataset to real-world contexts.

Furthermore, the baseline methods proposed in this study operate under the assumption that objects exhibit frontal-parallel characteristics and are non-deformable. However, it is essential to acknowledge that real-world conditions are considerably more intricate, and these methods may yield suboptimal performance when the underlying assumption is not met. Additionally, as we mentioned earlier, our methods are sensitive to the alignment between the box center of the target and reference frame, which poses another limitation.

## 8 Conclusion

In this work, we built a large-scale TTC dataset and provided a simple yet effective TTC estimation algorithm as baselines for the community. Our dataset is characterized by its focus on objects in driving scenes, which contains both urban and highway scenarios and covers a wider range of depth. We hope that our proposed dataset could facilitate the development of TTC estimation algorithms.

## References

- 356
- 357 [1] Abhishek Badki, Orazio Gallo, Jan Kautz, and Pradeep Sen. *Binary TTC: A temporal geofence*  
358 *for autonomous navigation*. In *CVPR*, 2021.
- 359 [2] Abhishek Badki, Alejandro Troccoli, Kihwan Kim, Jan Kautz, Pradeep Sen, and Orazio Gallo.  
360 *Bi3D: Stereo depth estimation via binary classifications*. In *CVPR*, 2020.
- 361 [3] Luca Bertinetto, Jack Valmadre, Joao F Henriques, Andrea Vedaldi, and Philip HS Torr. *Fully-*  
362 *convolutional siamese networks for object tracking*. In *ECCV*, 2016.
- 363 [4] Jeffrey Byrne and Camillo J Taylor. *Expansion segmentation for visual collision detection and*  
364 *estimation*. In *ICRA*, 2009.
- 365 [5] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu,  
366 Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. *nuScenes: A multimodal*  
367 *dataset for autonomous driving*. In *CVPR*, 2020.
- 368 [6] TA Camus. *Calculating Time-to-Contact using real-time quantized optical flow*. 1995.
- 369 [7] Ming-Fang Chang, John Lambert, Patsorn Sangkloy, Jagjeet Singh, Slawomir Bak, Andrew  
370 Hartnett, De Wang, Peter Carr, Simon Lucey, Deva Ramanan, et al. *Argoverse: 3d tracking and*  
371 *forecasting with rich maps*. In *CVPR*, 2019.
- 372 [8] Roberto Cipolla and Andrew Blake. *Surface orientation and time to contact from image*  
373 *divergence and deformation*. In *ECCV*, 1992.
- 374 [9] Yutao Cui, Cheng Jiang, Limin Wang, and Gangshan Wu.  
375 *MixFormer: End-to-End tracking with iterative mixed attention*. In *CVPR*, 2022.
- 376 [10] Erez Dagan, Ofer Mano, Gideon P Stein, and Amnon Shashua. *Forward collision warning with*  
377 *a single camera*. In *IV*, 2004.
- 378 [11] Tom van Dijk and Guido de Croon. *How do neural networks see depth in single images?* In  
379 *CVPR*, 2019.
- 380 [12] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov,  
381 Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. *FlowNet: Learning optical flow*  
382 *with convolutional networks*. In *ICCV*, 2015.
- 383 [13] Andreas Ess, Bastian Leibe, Konrad Schindler, and Luc Van Gool. *Robust multiperson tracking*  
384 *from a mobile platform*. *PAMI*, 2009.
- 385 [14] Lue Fan, Feng Wang, Naiyan Wang, and Zhaoxiang Zhang. *Fsd v2: Improving fully sparse 3d*  
386 *object detection with virtual voxels*. *arXiv preprint arXiv:2308.03755*, 2023.
- 387 [15] Martin A Fischler and Robert C Bolles. *Random sample consensus: a paradigm for model*  
388 *fitting with applications to image analysis and automated cartography*. *Communications of the*  
389 *ACM*, 1981.
- 390 [16] Zheng Ge, Songtao Liu, Feng Wang, Zeming Li, and Jian Sun.  
391 *YOLOX: Exceeding yolo series in 2021*. *arXiv preprint arXiv:2107.08430*, 2021.
- 392 [17] Andreas Geiger, Philip Lenz, and Raquel Urtasun. *Are we ready for Autonomous Driving? the*  
393 *KITTI vision benchmark suite*. In *CVPR*, 2012.
- 394 [18] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz  
395 Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming  
396 He. *Accurate, large minibatch SGD: Training ImageNet in 1 hour*. *arXiv preprint*  
397 *arXiv:1706.02677*, 2017.

- 398 [19] Berthold KP Horn, Yajun Fang, and Ichiro Masaki. Time to Contact relative to a planar surface.  
399 In *IV*, 2007.
- 400 [20] Berthold KP Horn, Yajun Fang, and Ichiro Masaki. Hierarchical framework for direct gradient-  
401 based Time-to-Contact estimation. In *IV*, 2009.
- 402 [21] Junhwa Hur and Stefan Roth. Self-supervised monocular scene flow estimation. In *CVPR*,  
403 2020.
- 404 [22] Eddy Ilg, Tonmoy Saikia, Margret Keuper, and Thomas Brox. Occlusions, motion and depth  
405 boundaries with a generic network for disparity, optical flow or scene flow estimation. In *ECCV*,  
406 2018.
- 407 [23] Bo Li, Junjie Yan, Wei Wu, Zheng Zhu, and Xiaolin Hu. High performance visual tracking with  
408 siamese region proposal network. In *CVPR*, 2018.
- 409 [24] Zhichao Li and Naiyan Wang. DML0: Deep matching lidar odometry. In *IROS*, 2020.
- 410 [25] MIA LOURAKIS. Using planar parallax to estimate the Time-to-Contact. In *CVPR*, 1999.
- 411 [26] Aashi Manglik, Xinshuo Weng, Eshed Ohn-Bar, and Kris Kitani. Future near-collision prediction  
412 from monocular video: Feasibility, dataset, and challenges. In *IROS*, 2019.
- 413 [27] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy,  
414 and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow,  
415 and scene flow estimation. In *CVPR*, 2016.
- 416 [28] Moritz Menze and Andreas Geiger. Object scene flow for Autonomous vehicles. In *CVPR*,  
417 2015.
- 418 [29] F. Meyer and P. Bouthemy. Estimation of Time-to-Collision maps from first order motion  
419 models and normal flows. In *ICPR*, 1992.
- 420 [30] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi,  
421 and Ren Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. In *ECCV*,  
422 2020.
- 423 [31] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics  
424 primitives with a multiresolution hash encoding. *ToG*, 2022.
- 425 [32] Ziqi Pang, Zhichao Li, and Naiyan Wang. Simpletrack: Understanding and rethinking 3d  
426 multi-object tracking. *arXiv preprint arXiv:2111.09621*, 2021.
- 427 [33] Ashutosh Saxena, Min Sun, and Andrew Y Ng. Make3d: Depth perception from a single still  
428 image. In *AAAI*, 2008.
- 429 [34] Thomas Schops, Johannes L Schonberger, Silvano Galliani, Torsten Sattler, Konrad Schindler,  
430 Marc Pollefeys, and Andreas Geiger. A multi-view stereo benchmark with high-resolution  
431 images and multi-camera videos. In *CVPR*, 2017.
- 432 [35] René Schuster, Oliver Wasenmüller, and Didier Stricker. Dense scene flow from stereo disparity  
433 and optical flow. *arXiv preprint arXiv:1808.10146*, 2018.
- 434 [36] MohammadReza Seyedi, MohammadReza Koloushani, Sungmoon Jung, and Arda Vanli. Safety  
435 assessment and a parametric study of forward collision-avoidance assist based on real-world  
436 crash simulations. *Journal of Advanced Transportation*, 2021.
- 437 [37] Muralidhara Subbarao. Bounds on Time-to-Collision and rotational component from first-order  
438 derivatives of image flow. *Computer Vision, Graphics, and Image Processing*, 1990.

- 439 [38] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul  
440 Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for  
441 autonomous driving: Waymo open dataset. In *CVPR*, 2020.
- 442 [39] Igor Vasiljevic, Nick Kolkin, Shanyi Zhang, Ruotian Luo, Haochen Wang, Falcon Z Dai,  
443 Andrea F Daniele, Mohammadreza Mostajabi, Steven Basart, Matthew R Walter, et al. Diode:  
444 A dense indoor and outdoor depth dataset. *arXiv preprint arXiv:1908.00463*, 2019.
- 445 [40] Sundar Vedula, Simon Baker, Peter Rander, Robert Collins, and Takeo Kanade.  
446 Three-Dimensional scene flow. In *ICCV*, 1999.
- 447 [41] Paul Voigtlaender, Jonathon Luiten, Philip HS Torr, and Bastian Leibe.  
448 Siam R-CNN: Visual tracking by re-detection. In *CVPR*, 2020.
- 449 [42] Chien-Yao Wang, Hong-Yuan Mark Liao, Yueh-Hua Wu, Ping-Yang Chen, Jun-Wei Hsieh, and  
450 I-Hau Yeh. CSPNet: A new backbone that can enhance learning capability of cnn. In *CVPR*  
451 *workshops*, 2020.
- 452 [43] Yukitoshi Watanabe, Fumihiko Sakaue, and Jun Sato. Time-to-Contact from image intensity.  
453 In *CVPR*, 2015.
- 454 [44] Yan Yan, Yuxing Mao, and Bo Li. Second: Sparsely embedded convolutional detection. *Sensors*,  
455 2018.
- 456 [45] Gengshan Yang and Deva Ramanan. Upgrading optical flow to 3D scene flow through optical  
457 expansion. In *CVPR*, 2020.
- 458 [46] Botao Ye, Hong Chang, Bingpeng Ma, Shiguang Shan, and Xilin Chen. Joint feature learning  
459 and relation modeling for tracking: A one-stream framework. In *ECCV*, 2022.
- 460 [47] Zhichao Yin, Trevor Darrell, and Fisher Yu. Hierarchical discrete distribution decomposition  
461 for match density estimation. In *CVPR*, 2019.
- 462 [48] Meixin Zhu, Xuesong Wang, and Jingyun Hu. Impact on car following behavior of a forward  
463 collision warning system with headway monitoring. *Transportation research part C: emerging*  
464 *technologies*, 2020.

465 **Appendix**

466 **Sensor Specification**

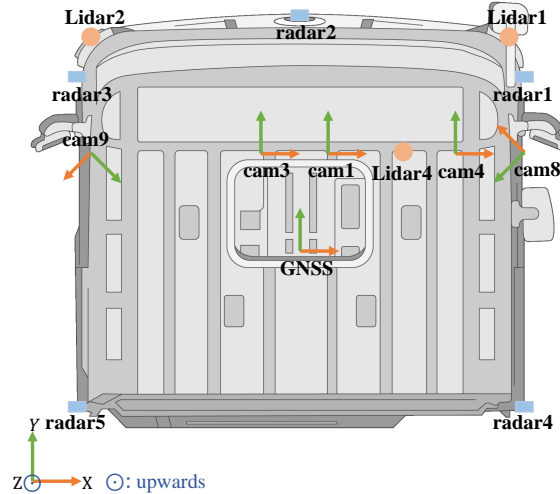


Figure 5: Sensor layout and coordinate system. The coordinate system of radars and Lidars are omitted for brevity. And GNSS means the Global Navigation Satellite System.

467 In our research, we conducted a comprehensive data collection process utilizing five cameras with  
 468 varying focal lengths, five radar sensors, and three Lidar sensors. The spatial arrangement of these  
 469 sensors is visually depicted in Fig. 5. Furthermore, we provide detailed specifications of the cameras  
 470 in Tab. 3.

471 **More details about data statistic**

472 Real scenes comprise 10.8% ur-  
 473 ban/suburban and 89.2% highway  
 474 data. However, rare cases with small  
 475 TTC on the same lane ( $[0,6]$ ) make up  
 476 only 0.02% of real scenes. To address  
 477 this, we have supplemented these rare  
 478 cases using data synthesized by NeRF.  
 479 The synthesized data, constituting 5K  
 480 sequences, represents highway scenes and  
 481 is exclusively used for the training set. The  
 482 data distribution of training set for object  
 483 sizes and time-of-day are shown in Fig 6.  
 484 As data in training, validation, and test sets  
 485 are randomly split, their distributions are

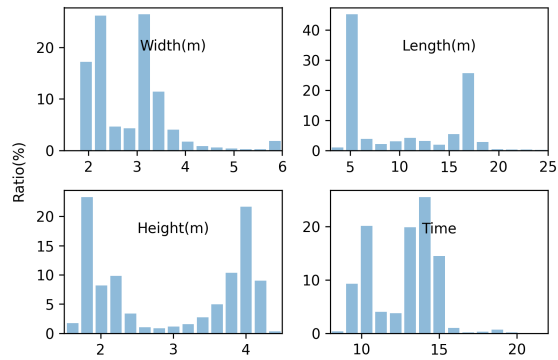


Figure 6: Distribution of object sizes and time-of-day.

Table 3: Camera specification. all images captured by the cameras are subjected to downsampling and cropping, resulting in a size of 1024x576 pixels. The camera’s horizontal field of view (HFOV) refers to the angular range covered by the camera along the y-axis in the x-y plane of the camera sensor frame.

Camera	1	3	4	8	9
HFOV	$\pm 63.2^\circ$	$\pm 40.4^\circ$	$\pm 18.4^\circ$	$\pm 63.2^\circ$	$\pm 63.2^\circ$

486 consistent, so we’ve omitted the latter two. Besides, the geographical distribution of the dataset spans  
 487 two cities: Shanghai and Hebei in China.

### 488 More explanation on MiD computing

489 In the Sec 5.3 of our main paper, we first convert the predicted scale ratio under different FPS settings  
 490 to corresponding TTC value and then convert the TTC value to the scale ratio under the setting of  
 491 10Hz. Actually, the scale ratios computed under different FPS settings could convert to each other.  
 492 Since the TTC ground-truth for the target frame is specific, we could get the relationship between the  
 493 scale ratios (*e.g.*  $\alpha_m, \alpha_n$ ) under different FPS settings (*e.g.*  $FPS_m, FPS_n$ ) by Eq.(5) in our paper:

$$\alpha_m = \frac{1}{\frac{FPS_n}{FPS_m}(1/\alpha_n - 1) + 1}. \quad (6)$$

494 In practice, the scale ratios obtained under different FPS settings are converted to the 10Hz setting  
 495 via Eq (6) directly.

### 496 Detailed Results

Table 4: Detailed results of different methods. The † means the result is obtained under padding NeRF data. The % after RTE is omitted.

Methods	Validation Set									
	MiD	MiD <sub>c</sub>	MiD <sub>s</sub>	MiD <sub>l</sub>	MiD <sub>n</sub>	RTE	RTE <sub>c</sub>	RTE <sub>s</sub>	RTE <sub>l</sub>	RTE <sub>n</sub>
Detection	213.9	675.4	305.1	112.4	115.3	54.2	58.1	56.3	55.0	50.2
SOT [46]	200.8	641.1	261.1	77.4	158.8	52.8	57.1	50.9	48.1	58.4
Pixel MSE	41.0	57.4	36.5	32.5	48.4	29.9	11.3	13.0	31.0	44.3
Depth	62.3	111.9	74.6	36.1	68.4	47.3	-	-	-	-
LiDAR	6.4	12.5	6.0	5.3	3.3	-	-	-	-	-
Deep Scale	14.4	27.1	16.4	10.9	13.5	12.1	6.3	8.7	13.0	14.8
Deep Scale†	14.3	26.5	15.2	10.8	13.5	12.0	6.2	8.4	12.9	14.6

497 We report detailed results of different methods on  
 498 validation set in Tab. 4 of both MiD and RTE. For  
 499 the Depth method, we first estimate the depth using  
 500 a mono depth estimation model and then utilize  
 501 RANSAC to fit the TTC value. Our experiments on  
 502 the validation set demonstrate that the relative error  
 503 of the depth estimation is only 10.7%. However, the  
 504 relative TTC error and MiD error are significantly  
 505 worse than our proposed method, reaching 47.3%  
 506 and 62.3 respectively. The primary reason for such  
 507 large errors is the inherent noise present in the depth  
 508 estimation. For the LiDAR model, it is a internal  
 509 sparse detector like SECOND [44] and FSDv2 [14].  
 510 The 3D bboxes in the training set were generated by  
 511 an internal algorithm. The tracker we used is Simple  
 512 Track [32]. The average BEV IoU between the  
 513 3D bboxes generated by internal algorithm and the  
 514 manual annotations on the validation and test sets is  
 515 83.1%.

516 In addition to quantitative results, we also illustrate  
 517 several cases in Fig. 7 for intuitive understanding. In  
 518 the last column, we draw the scaled target box in the  
 519 reference frame which is obtained by utilizing the center of the reference box and the scale ratio

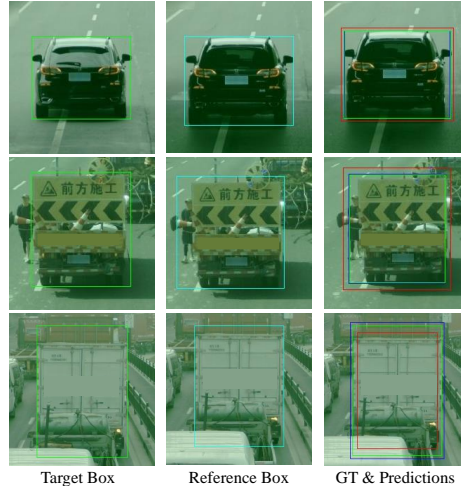


Figure 7: Case study for our Pixel MSE and Deep Scale, best viewed in color. For the last column, box with green, red and blue color denote the scaled box obtained by GT, Pixel MSE and Deep Scale.

Table 5: Influence of scale bins  $n$  in Pixel MSE

$n$	50	75	100	<b>125</b>	150
MiD	42.5	41.8	41.3	41.0	41.0
RTE(%)	31.9	31.9	30.8	29.9	30.0
Time(ms)	10.0	13.6	15.5	18.5	21.2

Table 6: Influence of scale bins  $n$  in Deep Scale.

$n$	10	15	<b>20</b>	25	30
MiD	16.8	14.9	14.4	14.4	14.4
RTE(%)	13.5	12.7	12.1	12.4	12.2
Time(ms)	11.7	11.9	12.0	12.2	12.4

Table 7: Ablation on the center shift operation for Pixel MSE and Deep Scale

	Pixel MSE		Deep Scale	
	w/ S	w/o S	w/ S	w/o S
MiD	41.0	73.3	14.4	17.2
RTE(%)	29.9	37.6	12.1	14.6
Time(ms)	18.5	17.4	12.0	10.8

Table 8: Ablation on padding different amounts of NeRF rendered sequences.

Seqs	MiD	MiD <sub>c</sub>	MiD <sub>s</sub>	MiD <sub>l</sub>	MiD <sub>n</sub>
0 K	14.4	27.1	15.5	10.9	13.5
3 K	14.5	27.1	15.4	10.8	13.6
<b>5 K</b>	14.3	26.5	15.2	10.8	13.7
7 K	14.6	26.7	15.4	11.2	13.8

520 obtained from various methods and ground truth. From the first and second cases, we can observe that  
 521 the Deep Scale is more robust to the illumination changes and inaccurate detection boxes compared  
 522 to Pixel MSE. In the last row, we showcase a failure case caused by a severely inaccurate detection  
 523 box. As described before, our methods are sensitive to the alignment between the box center of  
 524 the target and reference frame, which is also a limitation of our methods. These cases also reveal  
 525 the key difference between TTC estimation and tracking task: *TTC estimation requires far more*  
 526 *accurate estimation than tracking, while tracking usually focuses on complicated appearance change.*  
 527 Furthermore, we have included visualizations in GIF format in the supplementary material located  
 528 at `./Vis/monoDepth_visualization` to compare our proposed method and TTC estimation via  
 529 depth estimation. These visualizations provide a more intuitive understanding of the noise in the  
 530 depth estimation.

## 531 Ablation Study

532 To validate the contribution of different components and hyper-parameters, we perform extensive  
 533 experiments and report the results on the validation set unless otherwise specified. Default hyper-  
 534 parameters are denoted in **bold**.

535 **Number of Scale Bin.** To probe the suitable scale bins for our Pixel MSE, we conduct ablation  
 536 experiments, as shown in Table 5. The performance continues to boost until scale bin number reaches  
 537 125 and then becomes stable. Thus, we take the 125 bins as the default setting. To determine the  
 538 optimal scale bin number  $n$  for Deep Scale, we vary the value of  $n$  from 10 to 30. As shown in  
 539 Table 6, the MiD and RTE continuously decrease until  $n = 20$ , after which they tend to be stable.  
 540 However, an increase in scale bins results in a larger computation cost. As a consequence, we set  $n$  to  
 541 20 by default.

542 **Center Shift.** In this experiment, we present a comparison between the results obtained with and  
 543 without center shift operation. We list the results for both pixel MSE and Deep scale, as shown in  
 544 Table 7. The center shift operation is effective in reducing the estimation error.

545 **NeRF Augmentation.** To investigate the impact of supplementing NeRF data on the training process  
 546 for the Deep Scale, we conducted ablation experiments by adding different numbers of NeRF  
 547 sequences. The added NeRF sequences were randomly selected from the NeRF data we rendered.  
 548 To avoid the potential impact of data randomness on the experimental results, we report the average  
 549 results of five runs with different random seeds in Table 8. The results show that padding NeRF  
 550 sequences can effectively reduce the MiD error in crucial scenes. However, too many NeRF sequences  
 551 lead to a slight decrease in overall performance. Therefore, we use 5K rendered NeRF sequences in  
 552 our experiments.

553 **On Target Size.** In this experiment, we ablate the target size  $W, H$  for grid sample in Deep Scale  
 554 and we keep  $W = H$  when conducting ablation for convenience. Table 9 lists the result for different  
 555 settings and we can observe that the performance continuously to boost until 50 and then becomes  
 556 stable. As a consequence, we set the default target size to 50.

Table 9: Influence of the target size for grid sampling.

Size	10	25	50	75	100
MiD	20.0	14.9	14.4	14.8	14.9
RTE(%)	16.3	12.4	12.1	12.3	12.7

Table 10: Influence of the kernel size.

K	1	3	5	7	9
MiD	18.5	15.1	14.8	14.4	14.0
RTE(%)	15.2	12.3	12.3	12.1	11.8

557 **On Kernel Size.** Without large downsampling rate in our Deep Scale, the receptive field is mainly  
 558 decided by the kernel size of the convolutional layers. To find the optimal kernel size, we train our  
 559 model with the kernel size ranging from 1 to 9 and report the results in Table 10. Although increasing  
 560 the kernel size can improve the performance, it comes at the cost of longer inference latency. As a  
 561 trade off, we set the default kernel size to 7.

562 **On Plate Blur.** To verify whether the plate blur will influence the scale ratio estimation, we conduct  
 563 ablation experiments. We test the MiD and RTE on the validation, and test set in w/ and w/o blur  
 564 settings for the Pixel MSE. For the Deep Scale, we train the model on the blurred train/train+val  
 565 set and report the result of val/test set in w/ and w/o blur settings. As shown in Table 11, for both  
 566 methods, the plate blur operation brings negligible differences. As a conclusion, we take the dataset  
 567 with blurred license plates as the release version.

568 **Number of Frame Gap.** In this experiment, we validate the influence of frame gap for both Pixel  
 569 MSE and Deep Scale. The minimum and maximum scale ratio for different frame gaps are adjusted  
 570 according to Eq. (6). For the Deep Scale, we train our model in different frame gaps. As for testing,  
 571 we maintain the frame gap consistent with the training settings to ensure optimal results. We list  
 572 detailed results in Table 12. Larger frame gap brings more obvious scale changes and thus benefits  
 573 the classification process. As a result, we set the default frame gap to 5.

## 574 More Visualization

575 In Fig. 8, we present more samples from our dataset, covering different ranges of TTC, meteorological  
 576 fluctuations, and models rendered using NeRF. To get more intuitive understanding for different  
 577 methods, we present more cases for visualization in Fig. 9. In the first column of Fig. 9, we plot  
 578 the detection boxes of the target frames. In the second column, we show the boxes generated from  
 579 detection and tracking model. For the last column, we show the scaled boxes obtained by **GT**,  
 580 **Pixel MSE** and **Deep Scale**. As we can observe, the Pixel MSE produces unsatisfactory outcomes  
 581 when object images encounter significant illumination changes or low quality, as exemplified in the  
 582 first, second, and fourth cases. In contrast, the Deep Scale metric continues to perform robustly.  
 583 Nonetheless, severe occlusion remains a challenge that adversely affects the performance of both  
 584 Pixel MSE and Deep Scale, as demonstrated in the third case.

585 Additionally, we have included enhanced visualizations in GIF format within our supplementary  
 586 material. These visualizations, which can be found in the `./Vis/prediction_visualization`

Table 11: Ablation on plate blur for our methods

		Pixel MSE		Deep Scale	
		w/ blur	w/o blur	w/ blur	w/o blur
Val	MiD	41.0	41.0	14.4	14.4
	RTE(%)	29.9	30.0	12.1	12.1
Test	MiD	40.3	40.3	14.8	14.8
	RTE(%)	28.7	28.7	12.3	12.3



Table 12: Ablation on the frame gap for Pixel MSE and Deep Scale.

	Gap	1	2	3	4	5
PixelMSE	MiD	102.1	61.2	46.4	41.2	41.0
	RTE(%)	95.7	59.0	42.8	35.1	29.9
DeepScale	MiD	31.5	22.7	18.1	15.8	14.4
	RTE(%)	29.5	20.1	15.9	13.5	12.1

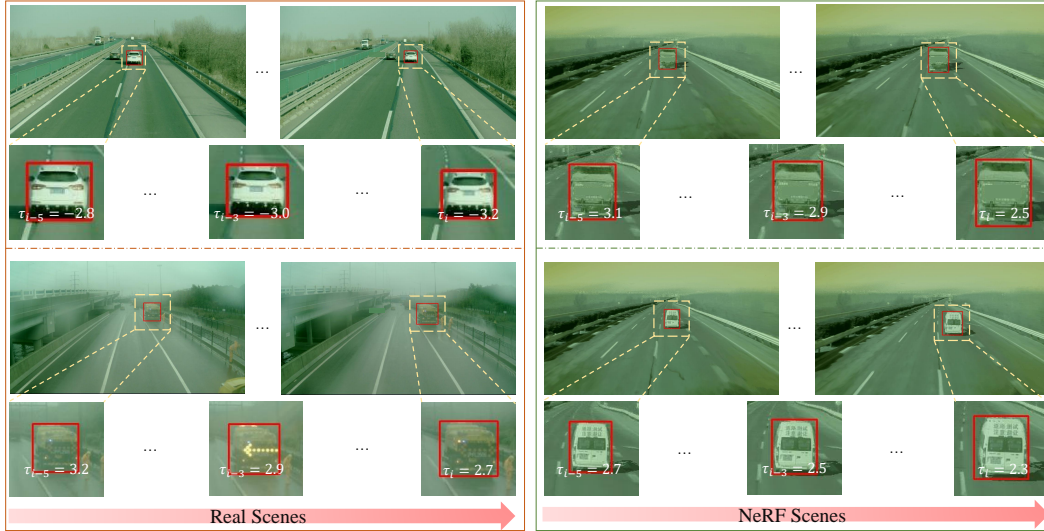


Figure 8: More visualization for samples in real and NeRF scenes.

587 and `./Vis/scene_visualization` directories, separately showcase qualitative results and various  
 588 scene representations.

### 589 NeRF Script

590 For the scripts used for NeRF rendering, we list them in Table 13.  $v_{ego}$  and  $v_{vehicle}$  denote the initial  
 591 speed of ego and the target vehicle respectively. The  $y$  denotes the relative distance in depth direction  
 592 and we only consider the straight lane when setting the scripts. We permute and combine the speed of  
 593 ego and vehicle to get more scenarios. For the rendered images, we also adopt the detection model to  
 594 generate 2D bounding boxes and the truncated objects will be discarded to ensure the completeness.

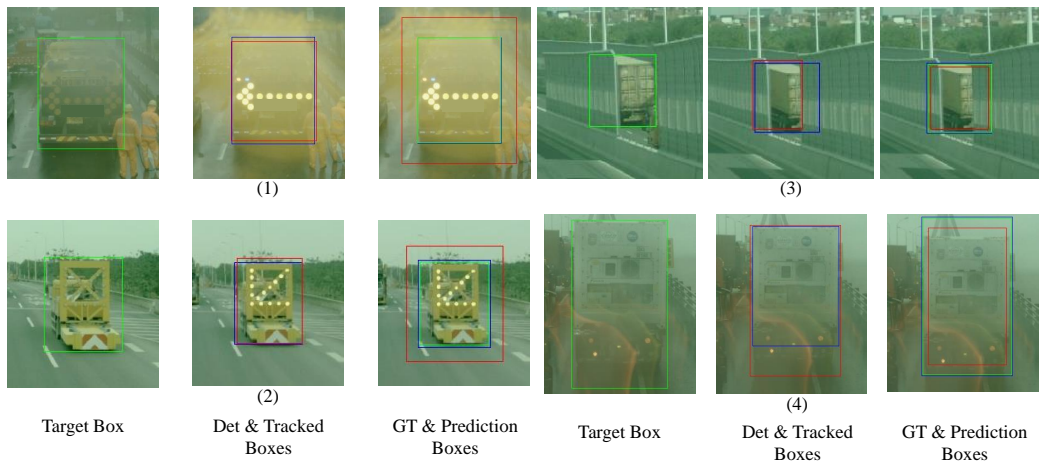


Figure 9: More visual comparison of the detection, tracking, and proposed methods, best viewed in color. In the second column, we use the blue and red color to distinguish the box from detection and tracking. For the last column, box with green, red and blue color denote the scaled box obtained by GT, Pixel MSE and Deep Scale.

Table 13: Script for NeRF rendering.

No.	Initial Status			Script
	$v_{ego}(km/h)$	$v_{vehicle}(km/h)$	$y(m)$	
1	40 60 80	$v_{ego}$	50	<ol style="list-style-type: none"> <li>Ego drives at <math>v_{ego}</math> while the target vehicle decelerates at a speed of <math>-3m/s^2</math>.</li> <li>After 3 seconds, ego decelerates with an acceleration of <math>-4m/s^2</math> until its velocity matches that of the target vehicle.</li> </ol>
2	40 60 80	$v_{ego}-20$	65	<ol style="list-style-type: none"> <li>Ego drives at <math>v_{ego}</math>.</li> <li>At a distance range of <math>(10, 50, 5)m</math> to the target vehicle, ego performs a lane change at a constant lateral relative speed of <math>2m/s</math>, until it is completely in a different lane from the target vehicle.</li> </ol>
3	60 80	20 40	65	<ol style="list-style-type: none"> <li>Ego gradually accelerates towards the target vehicle with an acceleration of range <math>(0.5, 3, 0.5)m/s^2</math>.</li> <li>When the distance between ego and the target vehicle is within the range of <math>(10, 50, 5)m</math>, the target vehicle changes lanes with a constant lateral relative speed of <math>2m/s</math> from an adjacent lane, until ego and the target vehicle are completely in the same lane. At the same time, ego decelerates at <math>-4m/s^2</math> within the same distance range of <math>(10, 50, 5)m</math>.</li> </ol>
4	60 80	60	65	<ol style="list-style-type: none"> <li>Ego drives at <math>v_{ego}</math> while the target vehicle decelerates at a speed of <math>-3m/s^2</math>.</li> <li>After 3 seconds, ego decelerates with an acceleration of <math>-4m/s^2</math> until its velocity matches that of the target vehicle.</li> </ol>
5	40 60	20 30	65	<ol style="list-style-type: none"> <li>Ego drives at <math>v_{ego}</math> while the target vehicle gradually accelerates with an acceleration range of <math>(0.5, 3, 0.5)m/s^2</math>.</li> <li>At a distance range of <math>(20, 60, 4)m</math> to the target vehicle, ego smoothly changes lanes with a lateral velocity range of <math>(0.5, 1.5, 0.2)m/s</math>, until ego and the target vehicle are completely in different lanes or the distance between them is less than <math>5m</math>.</li> </ol>
6	40 60 80	$v_{ego}-30$	65	<ol style="list-style-type: none"> <li>Ego drives at <math>v_{ego}</math> while the target vehicle gradually accelerates with an acceleration range of <math>(0.5, 3, 0.5)m/s^2</math>.</li> <li>When the distance to the target vehicle is in the range of <math>(10, 50, 4)m</math>, ego gradually decelerates with an acceleration of <math>-(1, 4, 0.5)m/s^2</math> until ego's velocity matches the target vehicle's velocity or the distance between them is less than <math>5m</math>.</li> </ol>

595 **Checklist**

- 596 1. For all authors...
- 597 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
- 598 contributions and scope? [Yes]
- 599 (b) Did you describe the limitations of your work? [Yes] See Sec. 7.
- 600 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See
- 601 Supplementary materials.
- 602 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
- 603 them? [Yes]
- 604 2. If you are including theoretical results...
- 605 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Sec. 5.
- 606 (b) Did you include complete proofs of all theoretical results? [Yes]
- 607 3. If you ran experiments (e.g. for benchmarks)...
- 608 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 609 mental results (either in the supplemental material or as a URL)? [Yes] See abstract.
- 610 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 611 were chosen)? [Yes]
- 612 (c) Did you report error bars (e.g., with respect to the random seed after running ex-
- 613 periments multiple times)? [No] The paper does not report error bars or statistical
- 614 significance information because the results of multiple experiments are consistent and
- 615 show minimal variation.
- 616 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 617 of GPUs, internal cluster, or cloud provider)? [Yes]
- 618 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 619 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 620 (b) Did you mention the license of the assets? [Yes]
- 621 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 622 See abstract.
- 623 (d) Did you discuss whether and how consent was obtained from people whose data you’re
- 624 using/curating? [Yes]
- 625 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 626 information or offensive content? [Yes]
- 627 5. If you used crowdsourcing or conducted research with human subjects...
- 628 (a) Did you include the full text of instructions given to participants and screenshots, if
- 629 applicable? [N/A]
- 630 (b) Did you describe any potential participant risks, with links to Institutional Review
- 631 Board (IRB) approvals, if applicable? [N/A]
- 632 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 633 spent on participant compensation? [N/A]