

# Mitigating Multi-Module Errors for Reliable Navigation in Dynamic Environments via Online Trajectory Refinement

Anonymous CVPR submission

Paper ID 0008

## Abstract

001 Perception in dynamic obstacle avoidance scenarios is inevitably  
002 affected by errors from multiple perception modules, e.g., robot  
003 localization, dynamic detection, and trajectory prediction. Existing  
004 methods mainly focus on reducing errors in a single module, overlooking  
005 their interactions and accumulation, which degrade navigation performance.  
006 We propose a systematic framework to mitigate multi-perception  
007 module errors and improve the reliability of dynamic obstacle avoidance.  
008 The proposed framework integrates two key innovations: an online trajectory  
009 refinement network called *Residual-DtACI*, and a point cloud-based  
010 dynamic target localization method. The *Residual-DtACI* adaptively  
011 refines trajectories through multi-expert residual estimation, while the  
012 point cloud-based method improves localization accuracy and robustness  
013 to occlusions and pose variations. We conducted several experiments to  
014 evaluate the proposed framework regarding its effectiveness in reducing  
015 accumulated perception errors. Results found that errors across different  
016 perception modules were reduced, alleviating error accumulation and  
017 enhancing prediction accuracy. In addition, the evaluations on public  
018 benchmarks show significant gains, with *Residual-DtACI* improving  
019 robot trajectory accuracy by up to 88.5%, and the point cloud-based  
020 method increasing localization accuracy by 59.9%.

## 026 1. Introduction

027 Perception plays a crucial role in the autonomous navigation of  
028 mobile robots in dynamic environments. It provides comprehensive  
029 information about the surrounding obstacles and the robot's own states.  
030 Such information is the foundation of path planning and decision making,  
031 determining whether the robot can reliably navigate [36]. However,  
032 perception in dynamic obstacle avoidance is unavoidably affected by  
033 errors accumulated from different components, which ultimately degrade  
034 the performance. Most  
035

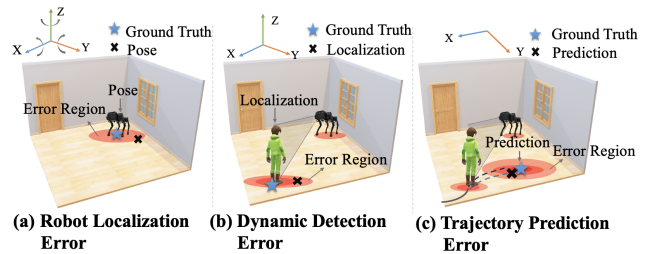


Figure 1. This figure shows perception error accumulating in dynamic obstacle avoidance. (a) Errors in the robot localization module lead to biases in the estimated pose trajectory. (b) Errors in the dynamic detection module cause inaccurate target localization. (c) Errors in the trajectory prediction module result in deviations in predicted motion. As these modules are sequentially coupled, the errors accumulate, ultimately degrading navigation performance.

existing studies have primarily focused on reducing errors within individual components such as localization and trajectory prediction. For localization, several works improved SLAM performance by enhancing feature extraction [23, 25] and tracking robustness under dynamic interference [16, 33, 39]. For trajectory prediction, destination-guided strategies have been adopted to address the indeterminacy in human motion [18, 21, 35].

However, these errors not only exist individually but also interact and accumulate over time due to the coupling among modules, directly impacting the effectiveness of dynamic obstacle avoidance. Existing works have not systematically analyzed and addressed these errors while considering their accumulated effects. Previous approaches that concentrate solely on reducing errors within a single component [7, 9, 16, 18, 21, 30, 33, 35, 39], while overlooking the error propagation from preceding and subsequent modules, cannot effectively reduce the influence of errors on the navigation system.

To overcome the above challenges, we proposed a framework that systematically addresses cumulative perception errors throughout the navigation process and improves the accuracy of dynamic target prediction. Specifically, our

059 framework focuses on three key components in the per-  
 060 ception pipeline, including robot localization, dynamic ob-  
 061 stacle detection, and trajectory prediction. It integrates  
 062 trajectory refinement and target localization optimization  
 063 to reduce accumulated perception errors and enhance the  
 064 reliability of dynamic obstacle avoidance. In particular,  
 065 we achieved the framework through two main innovations.  
 066 Firstly, we designed the Residual-DtACI network for tra-  
 067 jectory refinement. It builds upon the Conformal Infer-  
 068 ence method DtACI [10], extending with multi-dimensional  
 069 residual prediction. It does not require prior knowledge,  
 070 leverages multiple expert estimations, and demonstrates  
 071 strong adaptability to different types of trajectory errors. In  
 072 addition, we proposed a point cloud-based target localiza-  
 073 tion method. This method transforms RGB-D inputs into  
 074 3D point clouds of the target through depth back-projection.  
 075 The resulting point cloud is then denoised through a filter-  
 076 ing function, from which the target’s position is extracted.  
 077 In the context of dynamic navigation, and in contrast to ex-  
 078 isting 2D localization [20], our method demonstrates im-  
 079 proved localization accuracy in datasets with heavy occlu-  
 080 sions and pose variations.

081 We integrated the above two innovations into a unified  
 082 framework. Specifically, the Residual-DtACI is applied to  
 083 refine the robot pose trajectory in the robot localization  
 084 module as well as the predicted human trajectory in the tra-  
 085 jectory prediction module. In parallel, the proposed point  
 086 cloud-based method is employed in the dynamic detection  
 087 module to perform target localization. By systematically  
 088 reducing errors across perception modules and mitigating  
 089 their accumulation, the framework ultimately enables reli-  
 090 able, collision-free obstacle avoidance. The main contribu-  
 091 tions are as follows:

- 092 • We proposed a framework that systematically analyzes  
 093 the sources of perception errors and addresses their ac-  
 094 cumulated effects, effectively improving the accuracy of  
 095 dynamic target prediction and thereby enhancing the reli-  
 096 ability of dynamic obstacle avoidance.
- 097 • We proposed an online learning trajectory refinement net-  
 098 work called Residual-DtACI. This network adopts trajec-  
 099 tory residuals as the loss function to update the weights  
 100 of multiple residual estimation experts. Using a unified  
 101 scale, it refines both the robot pose trajectory in the local-  
 102 ization module and the person prediction trajectory in the  
 103 trajectory prediction module.
- 104 • We introduced a 3D point cloud-based dynamic target  
 105 localization method for the dynamic detection module.  
 106 Quantitative experiments demonstrate that this method  
 107 achieves higher localization accuracy than conventional  
 108 2D approaches, while also providing robustness to occlu-  
 109 sions and variations in target pose.

## 2. Related Works 110

### 2.1. Perception Errors in Dynamic Obstacle Avoid- 111 ance 112

In visual SLAM, both tracking and mapping are affected by 113  
 multiple factors, including illumination changes, dynamic 114  
 objects, and scene texture characteristics, which inevitably 115  
 introduce localization errors. To address these issues, sev- 116  
 eral works have sought to reduce photometric error to im- 117  
 prove the accuracy of camera pose tracking [7, 9, 30]. Other 118  
 approaches attempt to remove dynamic feature points to 119  
 mitigate the interference of moving objects [16, 33, 39]. 120  
 In addition, alternative feature representations, such as 3D 121  
 Gaussians [23] and neural point clouds [25], have been ex- 122  
 plored to improve reconstruction quality in complex and 123  
 highly textured 3D environments. Moreover, in the detec- 124  
 tion module, YOLOv8 [13] balances speed and accuracy, 125  
 enabling real-time detection on resource-constrained hard- 126  
 ware. However, existing localization methods often neglect 127  
 the geometric structure of objects, relying on 2D informa- 128  
 tion for position estimation, which leads to errors [20]. 129

Human trajectory prediction is particularly challenging 130  
 due to the dynamic nature of human motion, such as the 131  
 unknown long-term goals and the influence of social in- 132  
 teractions on individual behavior. Recent studies have fo- 133  
 cused on forecasting long-term destinations as an interme- 134  
 diate step, first estimating the destination and then generat- 135  
 ing the trajectory [18, 21, 35]. Recently, PPT [18] adopts a 136  
 progressive approach that enables the model to capture both 137  
 short-term dynamics and long-term dependencies. Other 138  
 works have emphasized the role of group dynamics by ex- 139  
 plicitly modeling pair-wise and higher-order interactions 140  
 [1, 15, 34]. MART [15], for instance, employs a hyper- 141  
 graph transformer network to capture both individual be- 142  
 haviors and group-level influences. Overall, existing works 143  
 primarily focus on addressing errors within individual mod- 144  
 ules, while overlooking cross-module interactions and ac- 145  
 cumulation of errors. Therefore, in this work, we aim to 146  
 systematically address accumulated perception errors. 147

### 2.2. Trajectory Refinement 148

Trajectory refinement has emerged as an effective strategy 149  
 to reduce prediction errors and has been widely applied in 150  
 motion prediction for autonomous driving [6, 11, 19, 27, 40, 151  
 41]. Shi et al. [27] incorporated local map information as 152  
 fine-grained trajectory features to refine predictions. Other 153  
 studies employed residual learning by adopting a simple 154  
 MLP to predict residuals [11] or introducing angular con- 155  
 straints to enhance refinement [19]. In addition, Zhou et al. 156  
 [40] proposed a scenario-adaptive framework that dynami- 157  
 cally adjusts the refinement process to improve prediction 158  
 efficiency. Trajectory refinement has also been applied in 159  
 human trajectory prediction. Yoon et al. [38] employed 160

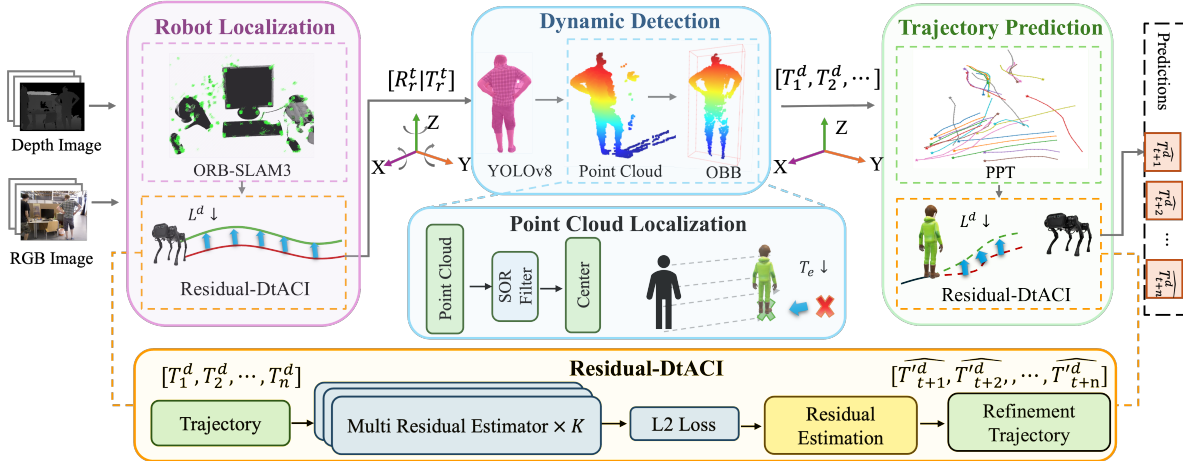


Figure 2. Overview of the proposed framework for systematically addressing perception errors in dynamic obstacle avoidance. The framework consists of three key modules: the robot localization module, which adopts ORB-SLAM3 as the baseline and applies the proposed Residual-DtACI network to refine the robot’s pose estimation trajectory; the dynamic detection module, which uses YOLOv8 for person detection and incorporates our point cloud–based method to improve dynamic target localization accuracy; and the trajectory prediction module, which employs PPT as the baseline predictor and integrates the Residual-DtACI with a unified scale to reduce trajectory prediction errors. Overall, the proposed framework systematically reduces accumulated perception errors.

161 convolutional neural networks to refine predicted trajectories, while Wang et al. [29] focused on modeling future  
 162 interactions among multiple agents to optimize refinement.  
 163 However, these refinement methods lack a unified scale ap-  
 164 plicable across different trajectory tasks in robotic systems.  
 165

### 166 3. Method

167 As illustrated in Fig. 2, we proposed a framework to reduce  
 168 accumulated perception errors in dynamic obstacle avoid-  
 169 ance scenarios. This framework focused on reducing errors  
 170 in three components, which are robot localization, dynamic  
 171 detection, and trajectory prediction. As a whole, these three  
 172 modules provide the robot with comprehensive perception  
 173 information for dynamic navigation, including surrounding  
 174 obstacles, self-localization, and target trajectory estimation.  
 175 We developed an online learning network of trajectory re-  
 176 finement called Residual-DtACI, which performed refine-  
 177 ment on both robot pose trajectories and dynamic target  
 178 prediction trajectories within a unified structure. In addi-  
 179 tion, we introduced a point cloud–based method to im-  
 180 prove the accuracy of object localization in the dynamic  
 181 detection module. First, we adopted ORB-SLAM3 [5] to estimate the  
 182 robot’s pose trajectory from RGB-D image inputs and ap-  
 183 plied our proposed Residual-DtACI network for trajectory  
 184 refinement. Next, YOLOv8 [13] was employed to detect  
 185 pedestrians in RGB images, and depth information was in-  
 186 corporated with our method to generate 3D point clouds of  
 187 the detected persons. A filtering function was then applied  
 188 to denoise the point cloud, and oriented bounding boxes  
 189 were used to extract the person’s localization. Finally, we

employed Progressive Pretext Task learning (PPT) [18] pre-  
 190 dictor to predict future trajectories of the persons, which  
 191 were further refined using the Residual-DtACI. The refined  
 192 trajectories were subsequently integrated into the navigation  
 193 path planner to enable reliable dynamic obstacle avoidance.  
 194

### 195 3.1. Accumulation of Perception Errors in Dynamic 196 Obstacle Avoidance Scenarios

This section describes how perception errors accumulate in  
 the process of dynamic obstacle avoidance scenarios. First,  
 in the SLAM process, dynamic objects, sensor noise, and  
 long-term drift, which affect the capture of feature points,  
 may cause the estimated robot pose  $\hat{P}_r^t$  to deviate from its  
 ground truth  $P_r^t$ , which inevitably introduces a pose error  
 $\varepsilon_r^t$ . This can be expressed as follows:

$$\hat{P}_r^t = P_r^t + \varepsilon_r^t \quad (1)$$

Next, YOLOv8 provides high detection accuracy with  
 relatively fast inference speed to obtain the person mask.  
 By combining the mask of the person with the depth image,  
 the depth values  $Z(u, v)$  corresponding to the masked pix-  
 els are retrieved. Through the intrinsic matrix  $K$  and the  
 extrinsic transformation  $\hat{P}_r^t$ , the estimated person localiza-  
 tion  $\hat{T}_p^t$  is obtained. However, due to noise in the point cloud  
 and errors in the extrinsic parameters (the robot pose esti-  
 mation error), an additional error term  $\varepsilon_p^t(\varepsilon_r^t)$  is introduced,  
 which can be formulated as:

$$\hat{T}_p^t = \hat{P}_r^t \cdot K \cdot Z(u, v) = T_p^t + \varepsilon_p^t(\varepsilon_r^t) \quad (2)$$

Finally, a sequence of person localizations is used as the  
 observed trajectory input to the trajectory predictor in order

to forecast the person's future trajectory. However, due to the indeterminate nature of human motion, existing models can hardly achieve highly accurate predictions, which introduces a trajectory prediction error  $\varepsilon^{t+1:t+H}$ . This error is further influenced by the preceding errors, as formulated below:

$$\hat{\tau}^{t+1:t+H} = g\left(\hat{T}_p^{1:t}\right) = \tau^{t+1:t+H} + \varepsilon_{\tau}^{t+1:t+H}(\varepsilon_p^t, \varepsilon_r^t) \quad (3)$$

Therefore, accumulated perception errors need to be addressed systematically. In our framework, we explicitly introduce refinement and correction mechanisms at multiple stages to reduce error propagation across the pipeline.

### 3.2. Trajectory Refinement

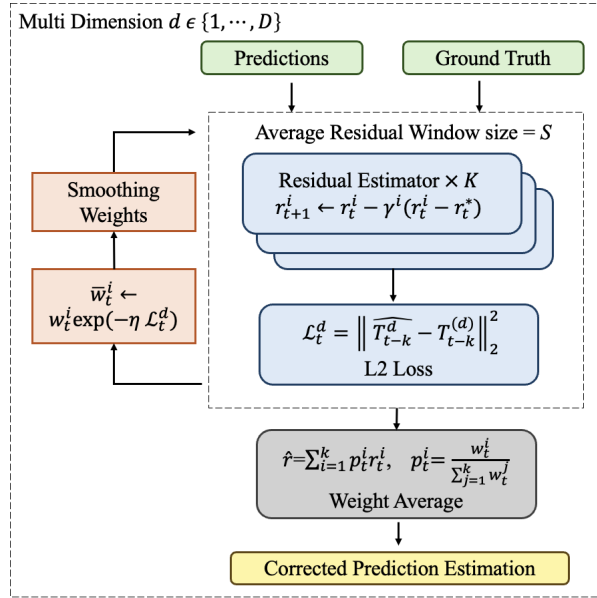


Figure 3. Architecture of the proposed Residual-DtACI network. The network consists of multiple expert models that perform online residual prediction. Their weights are updated using an  $L_2$  loss, with a smoothing mechanism applied to prevent expert collapse and maintain expert diversity.

Trajectory refinement aims to minimize the deviation between the ground-truth trajectory  $T^{1:t+H}$  and the predicted trajectory  $\hat{T}^{1:t+H}$ . In this stage, we leverage both the past ground truth and predicted trajectories  $\{T^{1:t}, \hat{T}^{1:t}\}$  as inputs to compute the offset  $\Delta T = T - \hat{T}$ . This offset is then used to refine the future prediction. We adopt an  $L_2$  loss to optimize the process:

$$\mathcal{L}^d = \left\| \Delta \hat{T}_t^d - \Delta T_t^d \right\|_2^2, \quad \forall 1 \leq d \leq D \quad (4)$$

where  $D$  denotes the dimensionality of the trajectory, and the loss is computed independently for each dimension.

### 3.3. Residual Dynamically-tuned Adaptive Conformal Inference

Building upon the trajectory refinement loss function introduced above, we further modified the Dynamically-Tuned Adaptive Conformal Inference (DtACI)[10] framework with our model and designed the proposed Residual-DtACI network, as shown in Fig 3. This network enabled online correction of predicted trajectories, allowing the system to rapidly adapt to distributional shifts without relying on prior knowledge of the underlying dynamics.

The network consists of multiple experts that perform error estimation  $\hat{r}_{t+1}^{i,d}$  in parallel. Each expert  $i$  is updated based on the gradient of the residual loss  $\nabla \mathcal{L}_t^d$ , while adopting a distinct learning rate  $\gamma^i$  to control the update speed. Larger values of  $\gamma^i$  enable faster adaptation to abrupt distributional shifts, whereas smaller values lead to more stable but slower updates. The update formula is as follows:

$$\hat{r}_{t+1}^{i,d} \leftarrow \hat{r}_t^{i,d} - \gamma^i \nabla \mathcal{L}_t^d = \hat{r}_t^{i,d} - \gamma^i (\hat{r}_t^{i,d} - r_t^{i,d}) \quad (5)$$

Each expert's weight  $\bar{w}_t^i$  is adaptively adjusted according to its loss  $\mathcal{L}_t^d$ , where the exponential weighting scheme places greater emphasis on recent performance. The parameter  $\eta$  is a hyperparameter that adjusts the speed at which the weights change and determines the sensitivity of the update:

$$\bar{w}_t^i \leftarrow w_t^i \exp(-\eta \mathcal{L}_t^d), \quad \forall 1 \leq i \leq k \quad (6)$$

After the exponential weighting, a smoothing step is applied to preserve weight diversity by mixing the normalized weights with a uniform distribution, ensuring that no expert loses all of its weight and that all experts retain a minimum probability of being selected. The  $\sigma$  is a smoothing coefficient:

$$w_{t+1}^i \leftarrow (1 - \sigma) \frac{\bar{w}_t^i}{\sum_{i=1}^k \bar{w}_t^i} + \frac{\sigma}{k}, \quad \forall 1 \leq i \leq k \quad (7)$$

Based on the probability distribution of each expert  $p_{t+1}^i$ , the final prediction correction  $\hat{r}_{t+1}^{d}$  is obtained as the weighted average of all expert predicted residuals  $r_{t+1}^{i,d}$ :

$$\hat{r}_{t+1}^d = \sum_{i=1}^k p_{t+1}^i \hat{r}_{t+1}^{i,d}, \quad p_{t+1}^i = \frac{w_{t+1}^i}{\sum_{i=1}^k w_{t+1}^i} \quad (8)$$

Finally, the correction values for each dimension are added to the corresponding predicted trajectory, thereby completing the refinement of the trajectory:

$$\hat{T}_{t+1}^d = \hat{T}_{t+1}^d + \hat{r}_{t+1}^d, \quad \forall 1 \leq d \leq D \quad (9)$$

### 3.4. Point Cloud-based Dynamic Target Localization

This section introduces a point cloud-based dynamic target localization method, which is designed to reduce the localization errors that occur during the dynamic detection process. Traditionally, localization is performed by applying

197  
198  
199  
200  
201

204  
205  
206  
207  
208  
209  
210  
211  
212  
213

214  
215  
216  
217  
218  
219

220 the YOLOv8 detection model to obtain bounding boxes of  
 221 persons from RGB images [20]. The center of each bound-  
 222 ing box, together with the corresponding depth value, is  
 223 then transformed into the world coordinate system and re-  
 224 garded as the center position of the person, as illustrated in  
 225 Fig. 4(a). However, this approach overlooks the inherent  
 226 three-dimensional characteristics of a person, which often  
 227 leads to errors in center localization.

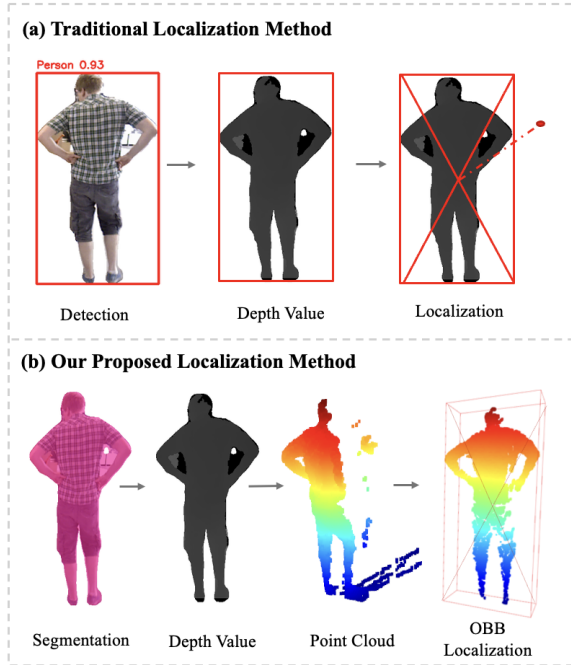


Figure 4. Illustration of dynamic target localization methods. (a) Traditional 2D approach using detection and depth value at the center of the box. (b) Our proposed point cloud-based method integrates segmentation, depth values, and point cloud cleaning with OBB to improve localization accuracy.

228 Thus, to better capture the 3D characteristics and spatial  
 229 geometric structure of a person, we propose a point cloud-  
 230 based dynamic target localization method. The proposed  
 231 method allows for a more accurate estimation of the per-  
 232 son’s center position in environments, as illustrated in Fig.  
 233 4(b).

We used the YOLOv8 segmentation model to obtain the person mask  $(u, v)$  and extracted the corresponding depth values  $Z(u, v)$ . These pixel depth pairs were then transformed into the world coordinate system using the camera intrinsic matrix  $K$  and extrinsic parameters  $[R | T]$ , resulting in the person point cloud  $T_p$ :

$$T_p(u, v) = R \cdot (Z(u, v) K^{-1}[u, v, 1]^T) + T \quad (10)$$

The extrinsic matrix  $[R | T] = P_r \cdot P_c$  is the composition of the robot’s pose in the world coordinate system  $P_r$  and

the camera’s pose relative to the robot  $P_c$ . Since the camera is rigidly mounted on the robot,  $P_c$  is typically fixed. Next, to eliminate the influence of noise in the person mask, we applied the Statistical Outlier Removal (SOR) filter function to clean the raw point cloud. For each point, the mean distance to its  $k$  nearest neighbors was computed. Any point with a mean distance greater than threshold  $D_{\max}$  was considered an outlier and removed:

$$D_{\max} = \mu + \lambda \times \sigma \quad (11)$$

Following cleaning the point cloud, an Oriented Bound-  
 ing Box (OBB) was fitted to the target person’s point cloud,  
 and its center was taken as the person’s localization.

## 4. Experiment

In this section, we elevated the effectiveness of our pro-  
 posed framework in systematically reducing accumulated  
 perception errors in dynamic obstacle avoidance scenarios.  
 First, we validate the quantitative benefits of our methods  
 in reducing errors on the three key perception modules,  
 which are robot localization, dynamic detection, and trajec-  
 tory prediction. Subsequently, we evaluate the performance  
 of the complete framework on the THUD++ [17] dataset  
 to demonstrate its ability to systematically reduce accumu-  
 lated perception errors. In particular, our experiments are  
 designed to answer the following four research questions:

- What is the quantitative benefit of trajectory refinement in reducing errors of robot pose trajectory estimation?
- What is the quantitative benefit of the point cloud-based method in reducing errors of dynamic target localization?
- What is the quantitative benefit of trajectory refinement in reducing errors of person trajectory prediction?
- What are the quantitative benefits of the complete framework in reducing accumulated perception errors across the perception pipeline?

### 4.1. Trajectory Refinement for Robot Pose Trajectory

We adopted ORB-SLAM3 as the baseline SLAM system and incorporated our Residual-DtACI network for robot pose trajectory refinement. Experiments were conducted on the TUM RGB-D [28] dataset, including 5 static sequences and 4 dynamic sequences, with the dynamic sequences containing highly dynamic objects interference. In addition to the baseline, we compared our approach with several state-of-the-art methods. For the dynamic sequences, we selected methods specifically designed to handle the interference of dynamic objects, Dynamic-VINS [33] and ViQu-SLAM [16]. For the static sequences, we included methods that emphasize accuracy in static environments, Point-SLAM [25] and GS-SLAM [23].

For evaluation, we used the Root Mean Square Error (RMSE) from the Absolute Trajectory Error (ATE). Prior

275 to error computation, all estimated trajectories were aligned  
 276 to the ground truth using the SE(3) Umeyama alignment  
 277 method. The Residual-DtACI network employs three error  
 278 estimators with learning rates  $\gamma$  of 0.1, 0.2, and 0.4. The  
 279 weight update parameter  $\eta$  is set to 2.72, and the smoothing  
 280 coefficient  $\sigma$  is fixed at 0.001.

281 **TUM RGB-D Static Dataset** Table 1 reports the tracking  
 282 results on the TUM RGB-D static dataset. Our method  
 283 achieved consistently better robot pose tracking perfor-  
 284 mance than the baseline ORB-SLAM3. In most sequences,  
 285 the Residual-DtACI refinement yielded the best RMSE per-  
 286 formance. For example, on the *fr2\_xyz* sequence, our  
 287 method achieved an RMSE of 0.0034 meters, correspond-  
 288 ing to a 10.53% improvement over ORB-SLAM3. Simi-  
 289 larly, on the *fr3\_office* sequence, we obtained an RMSE of  
 290 0.0066 meters, resulting in a 36.54% improvement over the  
 291 baseline. Beyond improving upon ORB-SLAM3, the re-  
 292 fined trajectories also demonstrated superiority over several  
 293 state-of-the-art methods, highlighting the effectiveness of  
 294 our Residual-DtACI network in enhancing localization ac-  
 295 curacy under static environments.

296 **TUM RGB-D Dynamic Dataset** Table 2 presents the  
 297 tracking results on the TUM RGB-D dynamic sequences.  
 298 These sequences are particularly challenging due to the  
 299 presence of moving objects, which often cause drift in the  
 300 tracking process. Our method achieved notable improve-  
 301 ments in RMSE over ORB-SLAM3 across all dynamic se-  
 302 quences. For instance, on the *fr3\_walking\_half* sequence,  
 303 our method achieved an 85.38% improvement over the  
 304 baseline, while on the *fr3\_walking\_rpy* sequence, the im-  
 305 provement reached 88.53%. Furthermore, compared with  
 306 state-of-the-art methods that rely on detection-and-removal  
 307 strategies for dynamic objects, our Residual-DtACI net-  
 308 work achieved a remarkable RMSE of 0.0307 meters on the  
 309 *fr3\_walking\_rpy* sequence.

310 The above results demonstrate that our Residual-DtACI  
 311 network is effective in both static and dynamic environ-  
 312 ments. In static sequences, pose tracking errors mainly arise  
 313 from sensor noise and long-term drift, whereas in dynamic  
 314 sequences, the interference of moving objects often leads to  
 315 severe tracking drift. Our refinement approach effectively  
 316 reduces these errors in both cases, achieving remarkable im-  
 317 provements in robot pose trajectory accuracy over the base-  
 318 line and current methods.

## 319 4.2. Dynamic Target Localization

We evaluated our proposed point cloud-based method for  
 dynamic target localization on the YCB-Video [32] and BE-  
 HAVE [4] datasets. As the baseline, we adopted a 2D detec-  
 tion-based method from prior work [20]. The YCB-Video  
 dataset is a widely used benchmark for object pose estima-  
 tion, providing accurate 6D poses of objects. The BEHAVE  
 dataset focuses on tracking human-object interactions. Al-

Table 1. RESULTS OF RMSE OF ATE [m] ON TUM RGB-D STATIC DATASETS. **BEST IS BOLD**; SECOND-BEST IS UNDERLINED.

Sequence	ORB-SLAM3	Point-SLAM	GS-SLAM	Ours	Improvement
fr1_desk	<u>0.0186</u>	0.0434	0.0150	0.0193	–
fr1_desk2	0.0296	0.0454	–	0.0246	16.89%
fr1_room	0.0695	0.3092	–	0.0296	57.41%
fr2_xyz	<u>0.0038</u>	0.0131	0.0144	0.0034	10.53%
fr3_office	<u>0.0104</u>	0.0348	0.0149	0.0066	36.54%

Table 2. RESULTS OF RMSE OF ATE [m] ON TUM RGB-D DYNAMIC DATASETS. **BEST IS BOLD**; SECOND-BEST IS UNDERLINED.

Sequence	ORB-SLAM3	Dynamic-VINS	ViQu-SLAM	Ours	Improvement
fr3_walking_xyz	0.4048	0.0486	0.1258	<u>0.0885</u>	78.14%
fr3_walking_static	0.0176	0.0077	<u>0.0124</u>	0.0135	23.30%
fr3_walking_rpy	0.2676	0.0629	<u>0.0720</u>	0.0307	88.53%
fr3_walking_half	0.3920	0.0608	0.0261	0.0573	85.38%

though it does not directly provide the human center position, it includes the 3D coordinates of the left hip joint  $H_{hip,l}$  and right hip joint  $H_{hip,r}$ . We derived the ground-truth human center position  $H_c$  as the average of the two hip joints:

$$H_c = \frac{H_{hip,l} + H_{hip,r}}{2} \quad (11)$$

320 The evaluation metrics were defined with reference to  
 321 6D pose estimation [32], including the position error  $T_e$   
 322 and the Average Distance (ADD) metric. The position error  $T_e$   
 323 is defined as the Euclidean distance between the estimated  
 324 and ground-truth center points. The ADD metric uses po-  
 325 sition error  $T_e$  and considers a prediction correct if it is  
 326 smaller than a predefined threshold, set to 10% of the 3D  
 327 model diameter.

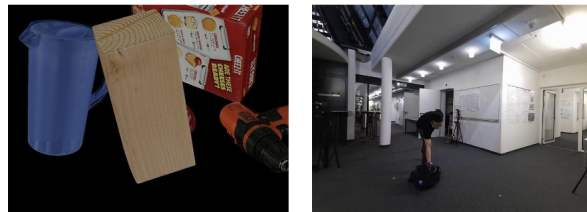


Figure 5. Examples from benchmark datasets of target localization. (Left) YCB-Video dataset, which contains objects with significant occlusions. (Right) BEHAVE dataset, which includes humans in diverse poses.

Table 3. RESULTS OF ADE [m] AND FDE [m] ON ETH/UYC DATASETS. **BEST IS BOLD**; SECOND-BEST IS UNDERLINED.

Method	ETH		HOTEL		UNIV		ZARA1		ZARA2		AVG	
	ADE↓	FDE↓	ADE↓	FDE↓	ADE↓	FDE↓	ADE↓	FDE↓	ADE↓	FDE↓	ADE↓	FDE↓
PPT [18]	0.36	0.49	<u>0.11</u>	<u>0.14</u>	0.22	0.40	0.17	0.30	0.18	0.33	0.21	0.33
NPSN [2]	0.36	0.59	0.16	0.25	0.23	0.39	0.18	0.32	0.14	0.25	0.21	0.36
MID [12]	0.39	0.66	0.13	0.22	0.22	0.45	0.17	0.30	<u>0.13</u>	0.27	0.21	0.38
ET [3]	0.36	0.53	0.12	0.19	0.24	0.43	0.19	0.33	0.14	0.24	0.21	0.34
LED [22]	0.39	0.58	<u>0.11</u>	0.17	0.26	0.43	0.18	<u>0.26</u>	<u>0.13</u>	<u>0.22</u>	0.21	0.33
MART [15]	0.35	0.47	0.14	0.22	0.25	0.45	0.17	0.29	<u>0.13</u>	<u>0.22</u>	0.21	0.33
TUTR [26]	0.40	0.61	<u>0.11</u>	0.18	0.23	0.42	0.18	0.34	<u>0.13</u>	0.25	0.21	0.36
TrajFine[29]	0.35	0.60	<u>0.11</u>	0.18	0.22	0.48	<u>0.15</u>	0.30	<b>0.12</b>	0.25	0.19	0.36
V <sup>2</sup> -Net [31]	<b>0.23</b>	<b>0.37</b>	<u>0.11</u>	0.16	0.21	0.35	0.19	0.30	0.14	0.24	0.18	<u>0.28</u>
V <sup>2</sup> -Net+TR [38]	<u>0.26</u>	<b>0.37</b>	<u>0.11</u>	0.15	<u>0.19</u>	<u>0.34</u>	0.17	0.30	<u>0.13</u>	0.23	<u>0.17</u>	<u>0.28</u>
<b>Ours</b>	0.27	<u>0.40</u>	<b>0.08</b>	<b>0.13</b>	<b>0.15</b>	<b>0.23</b>	<b>0.12</b>	<b>0.17</b>	<u>0.13</u>	<b>0.21</b>	<b>0.15</b>	<b>0.23</b>
Improvement	25.00%	18.37%	27.27%	7.14%	31.82%	42.50%	29.41%	43.33%	27.78%	36.36%	28.57%	30.30%

Table 4. RESULTS OF  $T_e$  [m] AND ADD [%] ON YCB-Video DATASETS. **BEST RESULTS ARE BOLD**; SECOND-BEST IS UNDERLINED.

Object	Baseline		Ours		Improvement	
	$T_e$ ↓	ADD ↑	$T_e$ ↓	ADD ↑	$T_e$ ↓	ADD ↑
003_cracker_box	0.0672	4.83	<b>0.0350</b>	<b>63.13</b>	47.92%	+58.30%
004_sugar_box	0.0459	0.48	<b>0.0204</b>	<b>80.73</b>	55.56%	+80.25%
008_pudding_box	0.0454	2.26	<b>0.0182</b>	<b>71.32</b>	59.91%	+60.06%
009_gelatin_box	0.0324	3.56	<b>0.0165</b>	<b>81.05</b>	49.07%	+77.49%
036_wood_block	0.0673	0.67	<b>0.0414</b>	<b>42.59</b>	38.48%	+41.92%
061_foam_brick	0.0372	<b>99.50</b>	<b>0.0281</b>	96.72	24.46%	–

**YCB-Video Dataset** Quantitative results of object localization on the YCB-Video dataset are provided in Table 4. Our method achieved lower position errors  $T_e$  and higher ADD scores than the baseline on most objects. For example, on *004\_sugar\_box*, our method improved ADD by 80.25%, while on *009\_gelatin\_box*, the position error was reduced by nearly half and the ADD increased by 77.49%. On the other hand, for *061\_foam\_brick*, the ADD metric did not show improvement, and the baseline method already achieved a position error of 0.0372 meters. We note that due to its simple geometry and limited occlusion, the 2D baseline was sufficient for handling such simple cases. However, in more complex scenarios with occlusions, our point cloud-based method demonstrates notable advantages, as shown in Fig. 5.

**BEHAVE Dataset** The person localization performance on the BEHAVE dataset is reported in Table 5. In this dataset, which focuses on human-object interaction scenarios, our method consistently achieved improvements across all sequences. For example, on the *01\_basketball* sequence,

Table 5. RESULTS OF  $T_e$  [m] AND ADD [%] ON BEHAVE DATASETS. **BEST RESULTS ARE BOLD**; SECOND-BEST IS UNDERLINED.

Sequence	Baseline		Ours		Improvement	
	$T_e$ ↓	ADD ↑	$T_e$ ↓	ADD ↑	$T_e$ ↓	ADD ↑
01_backpack_back	0.1993	65.85	<b>0.1383</b>	<b>78.57</b>	30.61%	+12.72%
01_basketball	0.2250	68.89	<b>0.1324</b>	<b>80.85</b>	41.16%	+11.96%
01_boxtiny_hand	0.2023	60.00	<b>0.1648</b>	<b>70.83</b>	18.54%	+10.83%
01_toolbox	0.2009	73.33	<b>0.1471</b>	<b>74.47</b>	26.78%	+1.14%
01_yogamat_hand	0.2470	68.75	<b>0.1497</b>	<b>72.92</b>	39.39%	+1.17%

the position error  $T_e$  was 0.1324 meters, representing a 41.16% reduction compared to the baseline, while the ADD improved by 11.96%. On the *01\_yogamat\_hand* sequence, the position error  $T_e$  decreased by 39.39%. The BEHAVE dataset involves humans performing diverse postures, such as bending, leaning, and crouching, as shown in Fig. 5. These results demonstrate the adaptability of our method to various human poses in complex interaction scenarios.

The above results confirm the effectiveness of our proposed point cloud-based method for dynamic target localization. In addition to human-centered scenarios, we also evaluated on an object dataset to extend the diversity of dynamic targets. These results highlight the potential of our point cloud-based approach to generalize across diverse dynamic targets and to remain robust in realistic scenarios involving occlusions and complex human object interactions.

### 4.3. Trajectory Refinement for Person Prediction Trajectory

For the trajectory prediction module, we used PPT [18] as the baseline predictor and applied our Residual-DtACI net-

328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347

348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367

work for trajectory refinement. The network structure and parameters were kept the same as in the robot pose trajectory refinement experiment. We evaluated the quantitative benefits of our method on the ETH/UCY dataset, in comparison with a wide range of current state-of-the-art approaches. The ETH [24]/UCY [14] dataset is a widely used benchmark for human trajectory prediction, consisting of two datasets with five different scenes, ETH and HOTEL from the ETH dataset, and UNIV, ZARA1, and ZARA2 from the UCY dataset.

For evaluation, we followed the previous trajectory prediction metrics: Average Displacement Error (ADE) and Final Displacement Error (FDE). ADE computes the mean Euclidean distance between the predicted and ground-truth positions across all time steps in the trajectory, thereby reflecting the overall prediction accuracy. In contrast, FDE measures the Euclidean distance between the predicted final position and the ground-truth final destination, capturing the accuracy of the final destination prediction. We followed the standard setting by taking the past 8 timestamps (3.2s) as the observed trajectory and predicting the subsequent 12 timestamps (4.8s) as the future trajectory.

**ETH/UCY Dataset** As shown in Table 3, we compared our method with ten state-of-the-art trajectory prediction methods. Our method achieved the best prediction performance, with an average ADE of 0.15 meters and an average FDE of 0.23 meters across all scenes. Per scene results remain competitive, and are best in several cases. For example, in the *HOTEL* scene, our method achieved an ADE of 0.08 meters and an FDE of 0.13 meters; in the *UNIV* scene, an ADE of 0.15 meters and an FDE of 0.23 meters; and in the *ZARA* scenes, an ADE of 0.12 meters and an FDE of 0.17 meters. As a trajectory refinement network, our approach demonstrated significant improvements over the baseline PPT model, reducing ADE from 0.21 to 0.15 meters (a 28.57% improvement) and FDE from 0.33 to 0.23 meters (a 30.30% improvement).

The experiment of trajectory prediction further demonstrates the effectiveness of our Residual-DtACI network for trajectory refinement. When applied to trajectories affected by different sources of error, our network was able to significantly reduce these errors. These results validate that the Residual-DtACI network provides a robust refinement mechanism for trajectory prediction tasks, achieving consistent improvements over the baseline as well as competitive or superior performance compared to state-of-the-art methods.

#### 4.4. Ablation Study

To further validate the overall effectiveness of our framework, we conducted an ablation study on THUD++ [17] dataset. In this study, we removed the error reduction methods from each module of our framework to evaluate their

Table 6. RESULTS OF ABLATION STUDY ON THUD++ DATASET. **BEST RESULTS ARE BOLD**; SECOND-BEST IS UNDERLINED.

Model Ablation	ADE↓	Degrad.	FDE↓	Degrad.
Pose w/o Residual-DtACI	0.2595	16.26%	0.2163	-0.1645
Loc. w/o point cloud	0.2904	30.11%	0.1762	-0.1244
Pred. w/o Residual-DtACI	0.3228	44.62%	0.4042	-0.3524
<b>Ours</b>	<b>0.2232</b>	-	<b>0.0518</b>	-

individual contributions. (a. Pose w/o Residual-DtACI) The robot pose trajectory was directly obtained from ORB-SLAM3 without refinement by Residual-DtACI; (b. Loc. w/o Point Cloud) Dynamic target localization was performed using the baseline method; (c. Pred. w/o Residual-DtACI) Future trajectories of human were predicted using PPT without refinement by Residual-DtACI. As shown in Table 6, the full framework achieves the best performance, while removing any component leads to a noticeable degradation.

#### 4.5. Evaluation on a Quadruped Robot

We implemented our framework on a Unitree Go2 robot with an onboard NVIDIA Jetson Orin computer and an Intel RealSense D435i depth camera. For the perception module, both the baseline components and our proposed optimization methods were directly deployed. For the path planning module, we adopted a cost map-based approach with the Fast Marching Method (FMM) as the foundation [8, 37], and further enhanced the navigation capability by incorporating local search in the command generator. The final linear and angular velocity commands were then transmitted to the Unitree Go2 through its high-level SDK. We conducted experiments in a laboratory environment, where the robot was assigned multiple fixed targets while a person approached it at different speeds.

### 5. Conclusion and Future work

This paper proposed a framework that systematically addresses accumulated perception errors in dynamic obstacle avoidance scenarios. By integrating our proposed Residual-DtACI network for trajectory refinement and the proposed point cloud-based method for target localization, our approach reduces errors across different perception modules. It enhances the performance of dynamic autonomous navigation. Future work will extend the framework to outdoor environments while considering the impact of illumination changes. Evaluations on benchmark datasets and real-world tests on the Unitree Go2 robot verify the effectiveness of the proposed framework.

459

**References**

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

- [1] Inhwon Bae, Jin-Hwi Park, and Hae-Gon Jeon. Learning pedestrian group representations for multi-modal trajectory prediction. In *European Conference on Computer Vision*, pages 270–289. Springer, 2022. 2
- [2] Inhwon Bae, Jin-Hwi Park, and Hae-Gon Jeon. Non-probability sampling network for stochastic human trajectory prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6477–6487, 2022. 7
- [3] Inhwon Bae, Jean Oh, and Hae-Gon Jeon. Eigentrajectory: Low-rank descriptors for multi-modal trajectory forecasting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10017–10029, 2023. 7
- [4] Bharat Lal Bhatnagar, Xianghui Xie, Ilya A Petrov, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Behave: Dataset and method for tracking human object interactions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15935–15946, 2022. 6
- [5] Carlos Campos, Richard Elvira, Juan J Gómez Rodríguez, José MM Montiel, and Juan D Tardós. Orb-slam3: An accurate open-source library for visual, visual-inertial, and multimap slam. *IEEE transactions on robotics*, 37(6):1874–1890, 2021. 3
- [6] Sehwan Choi, Jungho Kim, Junyong Yun, and Jun Won Choi. R-pred: Two-stage motion prediction via tube-query attention-based trajectory refinement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8525–8535, 2023. 2
- [7] Jakob Engel, Thomas Schöps, and Daniel Cremers. Lsd-slam: Large-scale direct monocular slam. In *European conference on computer vision*, pages 834–849. Springer, 2014. 1, 2
- [8] Zipeng Fu, Ashish Kumar, Ananye Agarwal, Haozhi Qi, Jitendra Malik, and Deepak Pathak. Coupling vision and proprioception for navigation of legged robots. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17273–17283, 2022. 8
- [9] Xiang Gao, Rui Wang, Nikolaus Demmel, and Daniel Cremers. Ldso: Direct sparse odometry with loop closure. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2198–2204. IEEE, 2018. 1, 2
- [10] Isaac Gibbs and Emmanuel J Candès. Conformal inference for online prediction with arbitrary distribution shifts. *Journal of Machine Learning Research*, 25(162):1–36, 2024. 2, 4
- [11] Carlos Gómez-Huélamo, Marcos V Conde, Rafael Barea, and Luis M Bergasa. Improving multi-agent motion prediction with heuristic goals and motion refinement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5323–5332, 2023. 2
- [12] Tianpei Gu, Guangyi Chen, Junlong Li, Chunze Lin, Yongming Rao, Jie Zhou, and Jiwen Lu. Stochastic trajectory prediction via motion indeterminacy diffusion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 17113–17122, 2022. 7
- [13] Glenn Jocher, Ayush Chaurasia, and Jing Qiu. Ultralytics yolo, 2023. [Online; accessed 11-Sep-2025]. 2, 3
- [14] Laura Leal-Taixé, Michele Fenzi, Alina Kuznetsova, Bodo Rosenhahn, and Silvio Savarese. Learning an image-based motion context for multiple people tracking. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3542–3549, 2014. 8
- [15] Seongju Lee, Junseok Lee, Yeonguk Yu, Taeri Kim, and Kyoo-bin Lee. Mart: Multiscale relational transformer networks for multi-agent trajectory prediction. In *European Conference on Computer Vision*, pages 89–107. Springer, 2024. 2, 7
- [16] Chengyang Li, Yulai Zhang, Zhiqiang Yu, Xinming Liu, and Qing Shi. A robust visual slam system for small-scale quadruped robots in dynamic environments. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 321–326. IEEE, 2024. 1, 2, 5
- [17] Zeshun Li, Fuhao Li, Wanting Zhang, Zijie Zheng, Xueping Liu, Yongjin Liu, and Long Zeng. Thud++: Large-scale dynamic indoor scene dataset and benchmark for mobile robots. *arXiv preprint arXiv:2412.08096*, 2024. 5, 8
- [18] Xiaotong Lin, Tianming Liang, Jianhuang Lai, and Jianfang Hu. Progressive pretext task learning for human trajectory prediction. In *European Conference on Computer Vision*, pages 197–214. Springer, 2024. 1, 2, 3, 7
- [19] Mengmeng Liu, Hao Cheng, Lin Chen, Hellward Broszio, Jiangtao Li, Runjiang Zhao, Monika Sester, and Michael Ying Yang. Laformer: Trajectory prediction for autonomous driving with lane-aware scene constraints. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2039–2049, 2024. 2
- [20] Yunfei Luan, Muhang He, Yudong Tian, Chengjie Lin, Yunhan Fang, Zihao Zhao, Jianxin Yang, and Yao Guo. Intelligent disinfection robot with high-touch surface detection and dynamic pedestrian avoidance. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3595–3601. IEEE, 2024. 2, 5, 6
- [21] Karttikeya Mangalam, Harshayu Girase, Shreyas Agarwal, Kuan-Hui Lee, Ehsan Adeli, Jitendra Malik, and Adrien Gaidon. It is not the journey but the destination: Endpoint conditioned trajectory prediction. In *European conference on computer vision*, pages 759–776. Springer, 2020. 1, 2
- [22] Weibo Mao, Chenxin Xu, Qi Zhu, Siheng Chen, and Yanfeng Wang. Leapfrog diffusion model for stochastic trajectory prediction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5517–5526, 2023. 7
- [23] Hidenobu Matsuki, Riku Murai, Paul HJ Kelly, and Andrew J Davison. Gaussian splatting slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18039–18048, 2024. 1, 2, 5
- [24] Stefano Pellegrini, Andreas Ess, Konrad Schindler, and Luc Van Gool. You’ll never walk alone: Modeling social behavior for multi-target tracking. In *2009 IEEE 12th international conference on computer vision*, pages 261–268. IEEE, 2009. 8

- 573 [25] Erik Sandström, Yue Li, Luc Van Gool, and Martin R Oswald. Point-slam: Dense neural point cloud-based slam. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 18433–18444, 2023. 1, 2, 5 630
- 574 631
- 575 [26] Liushuai Shi, Le Wang, Sanping Zhou, and Gang Hua. Trajectory unified transformer for pedestrian trajectory prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9675–9684, 2023. 7 632
- 576 633
- 577 [27] Shaoshuai Shi, Li Jiang, Dengxin Dai, and Bernt Schiele. Motion transformer with global intention localization and local movement refinement. *Advances in Neural Information Processing Systems*, 35:6531–6543, 2022. 2 634
- 578 635
- 579 [28] Jürgen Sturm, Nikolas Engelhard, Felix Endres, Wolfram Burgard, and Daniel Cremers. A benchmark for the evaluation of rgb-d slam systems. In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pages 573–580. IEEE, 2012. 5 636
- 580 637
- 581 [29] Kuan-Lin Wang, Li-Wu Tsao, Jih-Ciang Wu, Hong-Han Shuai, and Wen-Huang Cheng. Trajfine: Predicted trajectory refinement for pedestrian trajectory forecasting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4483–4492, 2024. 3, 7 638
- 582 639
- 583 [30] Rui Wang, Martin Schworer, and Daniel Cremers. Stereo dso: Large-scale direct sparse visual odometry with stereo cameras. In *Proceedings of the IEEE international conference on computer vision*, pages 3903–3911, 2017. 1, 2 640
- 584 641
- 585 [31] Conghao Wong, Beihao Xia, Ziming Hong, Qinmu Peng, Wei Yuan, Qiong Cao, Yibo Yang, and Xinge You. View vertically: A hierarchical network for trajectory prediction via fourier spectrums. In *European Conference on Computer Vision*, pages 682–700. Springer, 2022. 7 642
- 586 643
- 587 [32] Yu Xiang, Tanner Schmidt, Venkatraman Narayanan, and Dieter Fox. Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes. *arXiv preprint arXiv:1711.00199*, 2017. 6 644
- 588 645
- 589 [33] Linhui Xiao, Jinge Wang, Xiaosong Qiu, Zheng Rong, and Xudong Zou. Dynamic-slam: Semantic monocular visual localization and mapping based on deep learning in dynamic environment. *Robotics and Autonomous Systems*, 117:1–16, 2019. 1, 2, 5 646
- 590 647
- 591 [34] Chenxin Xu, Maosen Li, Zhenyang Ni, Ya Zhang, and Siheng Chen. Groupnet: Multiscale hypergraph neural networks for trajectory prediction with relational reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6498–6507, 2022. 2 648
- 592 649
- 593 [35] Chenxin Xu, Weibo Mao, Wenjun Zhang, and Siheng Chen. Remember intentions: Retrospective-memory-based trajectory prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6488–6497, 2022. 1, 2 649
- 594 650
- 595 [36] Qi Yin, Yang Ye, Meida Chen, and Yangming Shi. A Teacher–Student Learning Approach to Improve Quadruped Robots’ Autonomous Locomotion and Obstacle Avoidance on Construction Sites. *Available at SSRN 5907605*, 2026. 8 630
- 596 631
- 597 [38] Hanbit Yoon, Usman Ali, Joonhee Choi, and Eunbyung Park. Rethinking convolutional neural networks for trajectory refinement. *Pattern Recognition*, 157:110883, 2025. 2, 7 632
- 598 633
- 599 [39] Chao Yu, Zuxin Liu, Xin-Jun Liu, Fugui Xie, Yi Yang, Qi Wei, and Qiao Fei. Ds-slam: A semantic visual slam towards dynamic environments. In *2018 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 1168–1174. IEEE, 2018. 1, 2 634
- 600 635
- 601 [40] Yang Zhou, Hao Shao, Letian Wang, Steven L Waslander, Hongsheng Li, and Yu Liu. Smartrefine: A scenario-adaptive refinement framework for efficient motion prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15281–15290, 2024. 2 636
- 602 637
- 603 [41] Zikang Zhou, Jianping Wang, Yung-Hui Li, and Yu-Kai Huang. Query-centric trajectory prediction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 17863–17873, 2023. 2 638
- 604 639
- 605 640
- 606 641
- 607 642
- 608 643
- 609 644
- 610 645
- 611 646
- 612 647
- 613 648
- 614 649
- 615 650
- 616 651
- 617 652
- 618 653
- 619 654
- 620 655
- 621 656
- 622 657
- 623 658
- 624 659
- 625 660
- 626 661
- 627 662
- 628 663
- 629 664