

AirSim360: A Panoramic Simulation Platform within Drone View

Anonymous CVPR submission

Paper ID 5082

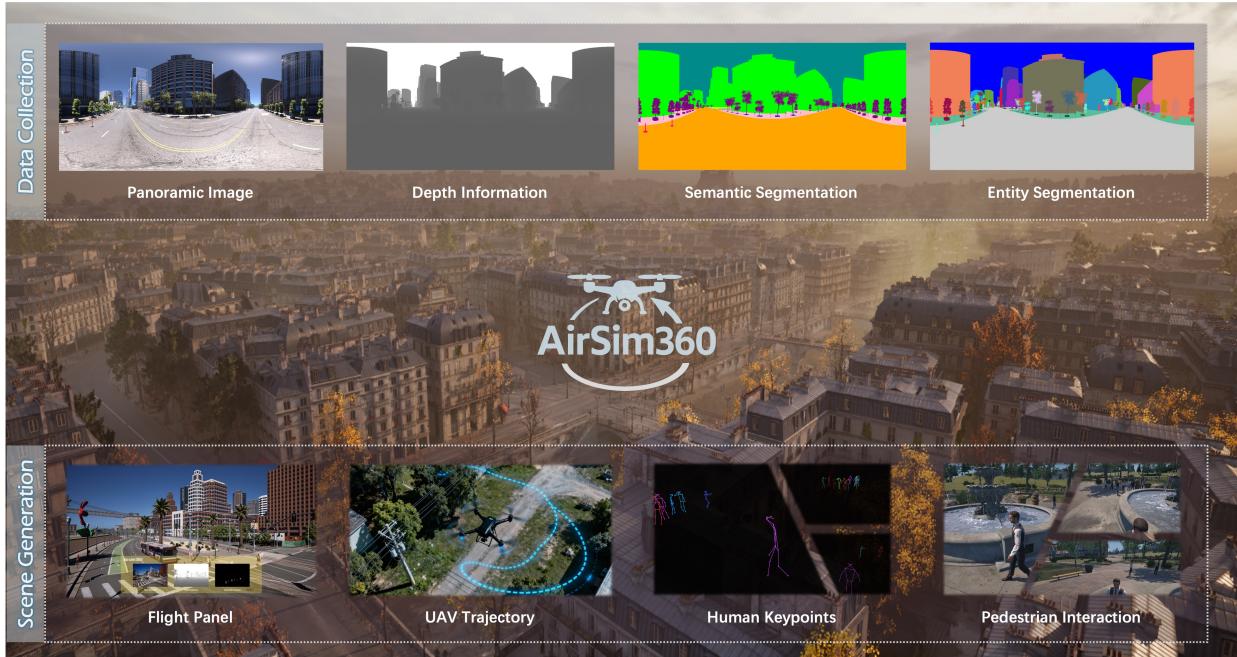


Figure 1. Overview of Airsim360. This work introduces a panoramic UAV simulation platform based on a cutting-edge rendering engine, enabling closed-loop simulation for omnidirectional aerial systems and offering an integrated toolkit for intelligent data acquisition across diverse flight scenarios.

Abstract

The field of 360-degree omnidirectional understanding has been receiving increasing attention for advancing spatial intelligence. However, the lack of large-scale and diverse data remains a major limitation. In this work, we propose AirSim360, a simulation platform for omnidirectional data from aerial viewpoints, enabling wide-ranging scene sampling with drones. Specifically, AirSim360 focuses on three key aspects: a render-aligned data and labeling paradigm for pixel-level geometric, semantic, and instance-level understanding; an interactive pedestrian-aware system for modeling human behavior; and an automated trajectory generation paradigm to support navigation tasks. Furthermore, we collect more than 60K panoramic samples and conduct

extensive experiments across various tasks to demonstrate the effectiveness of our simulator. Unlike existing simulators, our work is the first to systematically model the 4D real world under an omnidirectional setting. The entire platform, including the toolkit, plugins, and collected datasets, will be made publicly available.

1. Introduction

Embodied intelligence with omnidirectional perception has gained increasing attention due to the 360° full view in spatial intelligence. It can benefit various robotic applications, such as omnidirectional obstacle avoidance during navigation tasks.

Different from large-scale perspective image datasets, omnidirectional data remain scarce due to the

015
016
017
018
019
020

021
022
023
024
025
026
027
028

Table 1. Comparison of Simulators. This table compares recent simulation platforms across nine key dimensions, including rendering quality, interface support, and data generation flexibility, highlighting the advantages of our generative-AI-oriented simulator. Our **specific characters** are highlighted. “Ent” indicate the entity segmentation that cover the whole image.

Platform	AirSim [33]	CARLA [12]	Cosys-AirSim [18]	OmniGibson [21]	UnrealZoo [43]	OpenFly [14]	AirSim360 (ours)
<i>Realism level</i>	Medium	High	High	Medium	High	High	High
<i>Configurable Dynamics</i>	No	No	No	No	No	No	Yes
<i>API type</i>	Python	Python	Python	Python	Python	Python	Python / Blueprint
<i>Scenario type</i>	Outdoor/Indoor	Outdoor	Outdoor/Indoor	Outdoor/Indoor	Outdoor/Indoor	Outdoor	Outdoor/Indoor
<i>UAV simulator</i>	Yes	No	Yes	No	Yes	Yes	Yes
<i>Annotation capability</i>	Dep/Seg	Dep/Seg/Ins	Dep/Seg/Ins	Dep/Seg/Ins	Dep/Seg/Ins	Dep/Seg/Ins	Dep/Seg/Ent
<i>Custom label support</i>	No	No	Yes	Yes	No	No	Yes
<i>Panoramic image</i>	Yes	No	No	No	No	No	Yes
<i>Rendering engine</i>	UE 4.27	UE 5	UE 5.2	IsaacSim	UE 5	UE 5	UE 4.27 - UE 5.6

029 limited use of 360° cameras in daily life, not to mention
030 the exhaustive human labeling required for many tasks.
031 As a result, most panoramic methods are restricted by
032 small datasets and have scarcely explored data scaling.

033 Inspired by recent advances in simulation plat-
034 forms [34], a straightforward solution is to rotate agent
035 across multiple angles in simulator to capture an omni-
036 directional view, including both images and correspond-
037 ing ground truth. However, this approach introduces two
038 major issues. First, it is computationally inefficient, re-
039quiring repeated rendering and significantly increasing
040 data collection time. Second, the definitions of ground-
041 truth signals are not aligned with those in the perspective
042 domain. For example, omnidirectional depth represents
043 slant range along the viewing ray rather than the orthog-
044 onal z-axis distance used in perspective projection.

045 In this work, we focus on building AirSim360, an
046 omnidirectional simulation platform in the drone view
047 to model the 4D real world that consists of high-quality
048 static environments and movable pedestrians. The rea-
049 son we choose UAV as our agent because it can sample
050 much more data by exploring a wider range of spaces
051 than a ground-based one. Based on the Unreal Engine
052 (UE) 5 series for scene rendering, AirSim360 integrates
053 custom dynamics and communication modules into UE,
054 enabling UAVs to execute actions driven by an external
055 physical modeling engine and serving as the core engine
056 for our entire data-collection toolkit. Table 1 summa-
057 rizes the core capabilities of Airsim360, which offers a
058 comprehensive API suite ranging from data acquisition
059 (top) to flight control interfaces (bottom). Unlike exist-
060 ing platforms, our simulator supports full runtime inter-
061 action, enabling the generation of video-level panoptic
062 segmentation annotations.

063 Specifically, our AirSim360 has three main charac-
064 ters including render-aligned data and label generation,
065 interactive pedestrian-aware system, and automated tra-
066 jectory generation paradigm. For data collection, we
067 adopt the equirectangular projection (ERP) as an om-

nidirectional representation, which can compress mul-
068 tiple perspective views in a single continuous image.
069 However, this representation introduces several chal-
070 lenges for both image content and human- or pixel-
071 level annotations that should be carefully addressed. On
072 the one hand, we propose a GPU-side texture copying
073 mechanism based on Render Hardware Interface (RHI)
074 to enable fast and seamless stitching of six cube-face
075 views captured simultaneously, while maintaining light-
076 ing consistency across all faces. On the other hand,
077 the depth information is re-calculated to the slant range
078 along the viewing ray and segmentation annotations
079 have semantic and entity-level labels across frames, en-
080 suring complete panoramic coverage and consistent en-
081 tity identities over time. Furthermore, we develop point-
082 to-point path planning schemes for automatic data col-
083 lection and introduce an interactive pedestrian-aware
084 system that simulates movable pedestrians with various
085 actions and annotated keypoints, thus enhancing human-
086 centric perception and analysis.

087 Finally, extensive experiments across five panoramic
088 tasks, including depth estimation, semantic/entity seg-
089 mentation, human keypoint detection, and vision-
090 language navigation, demonstrate that our simulated
091 data transfers effectively to real-world scenarios.

092 Our contributions are fivefold.

- We propose a simulation platform, Airsim360, with native support for 360-degree aerial scenarios. Built on the UE 5 Series with an aerial dynamics module for high-quality scene rendering, the platform provides offline data generation tools capable of capturing panoramic images and corresponding ground truth, including semantic segmentation, entity segmentation, depth estimation, and 3D human keypoints.
- For efficient data collection, we develop point-to-point path planning schemes that enable both automatic sampling and user-designed flight paths. Moreover, we introduce an interactive pedestrian-aware system

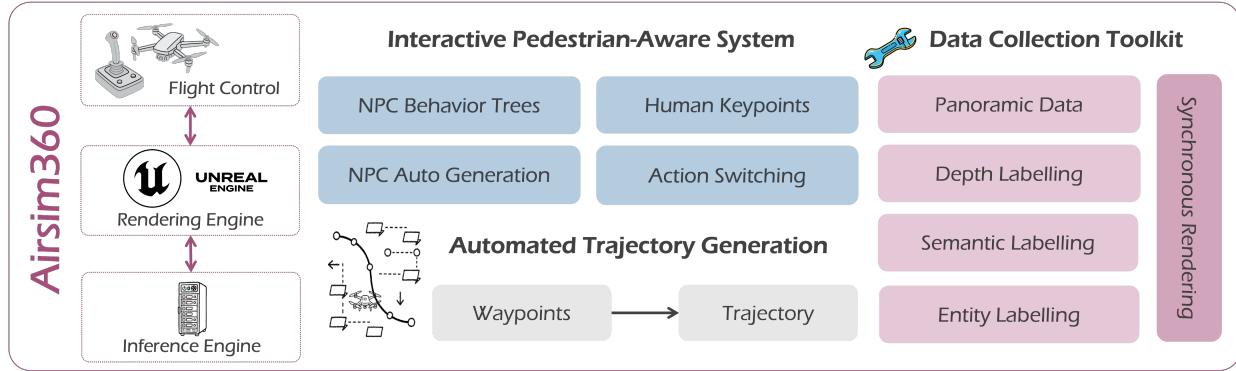


Figure 2. The left panel depicts the core interaction architecture of Airsim360, comprising the flight control module, rendering engine, and inference engine. The middle and right panels showcase three purpose-built data generation modules: the Interactive Pedestrian-Aware System, Automated Trajectory Generator, and Data Collection Toolkit.

that synthesizes pedestrian behaviors with detailed annotations and models realistic interactions.

- We test our simulated data for various perception and navigation tasks. Extensive experiments demonstrate our data significantly improve performance and robustness when evaluated on real-world validation sets.
- To our knowledge, AirSim360 is the first platform to support omnidirectional navigation for UAVs. Compared with conventional monocular systems, panoramic UAVs provide superior perceptual coverage, enabling more efficient target search with minimal additional motion cost.
- AirSim360 has excellent backward compatibility with ranging from latest UE 5.6 down to 4.27, our toolkit and benchmark tasks can run across these versions.

UAVScenes [36] and UnrealZoo. Among them, AirSim is an earlier simulator that utilizes UE for aerial scene rendering, but its latest officially supported version remains UE 4.27. UnrealCV and UnrealZoo are plugins built on top of UE that provide socket-based interfaces to capture RGB images, depth estimation, and semantic segmentation. However, they only support conventional perspective views and provide limited geometric and semantic information at relatively low frame rates, while lacking entity-level discrimination and making them insufficient for complex panoramic UAV tasks. In Table 1, our UAV360 platform additionally enables omnidirectional perception with full 360-degree coverage, and offers entity-level segmentation and 3D keypoint annotations to support advanced panoramic UAV applications at high frame rate.

2. Related Work

UAV Dataset UAV datasets typically contain raw images and task-related ground truth, such as trajectory waypoints, depth maps, semantic labels, and point clouds. In general, such data are obtained through acquisition of real-world flights with manual annotation [13, 20, 37] or all-in-one simulation platforms [24, 26, 32]. However, existing datasets are designed primarily for perspective views with limited fields of view. In contrast, we focus on an omnidirectional setting with 360° coverage in aerial scenarios, which enables panoramic UAV tasks such as omnidirectional obstacle avoidance.

Panoramic Task Panoramic visual tasks in understanding [38, 40] and generation [7, 8, 17] have rapidly advanced in support of spatial intelligence [23], where equirectangular projection (ERP) remains the most prevalent representation due to its straightforward mapping from spherical coordinates to a rectangular image. However, existing methods typically rely on real-world datasets with limited scale and diversity [3, 10, 11, 28, 39, 41], leaving their generalization ability across domains underexplored. In this work, we systematically benchmark depth estimation, pixel-level scene parsing, human keypoint detection, and UAV navigation under both in-domain and out-of-domain settings, demonstrating the merit of the proposed UAV360 platform.

3. AirSim360 Platform

As shown in Figure 2, AirSim360 is a large-scale omnidirectional simulation platform for drone scenarios. Unlike other platforms such as AirSim and UnrealZoo, we

106
107
108
109
110
111
112
113
114
115
116
117
118
119
120

140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155

121
122
123
124
125
126
127
128
129
130
131
132

156
157
158
159
160
161
162
163
164
165
166
167
168
169

133
134
135
136
137
138
139

Embodied Simulator Embodied simulators [2, 6, 30] are crucial in robotics, such as autonomous driving and robotic grasping. It is because real-world data collection is often challenging, especially for rare long-tail scenarios. As UAVs play an increasingly important role, a range of UAV-oriented simulation platforms have emerged, including AirSim, UnrealCV [31], OpenFly,

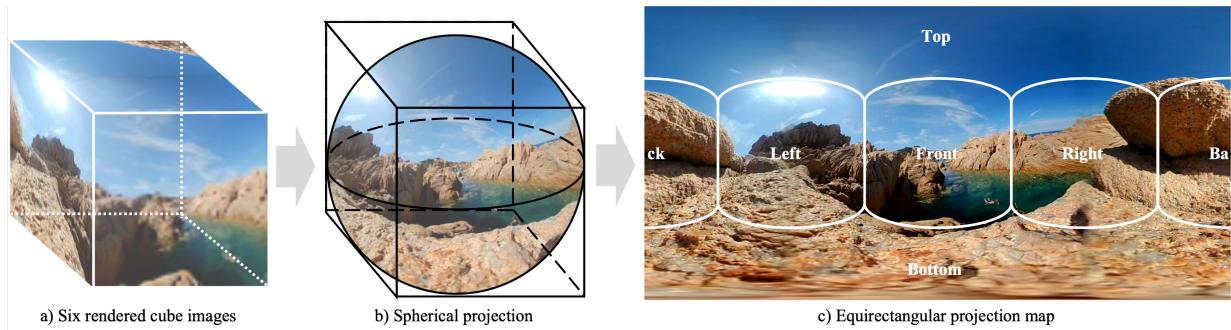


Figure 3. The equirectangular projection serves as a standard representation for panoramic imagery. In our implementation, we configured six identically calibrated cameras facing different directions. Since these are idealized pinhole cameras without lens distortion in the simulation environment, their outputs can be seamlessly stitched into a corresponding panoramic image through camera stitching algorithms.

emphasize consistency and compatibility among three key components in low-altitude ground environments, including static surroundings, the UAV, and human actors, which aligns closely with the main characteristics of the human-centric real world.

In the following subsections, we first introduce our simulation architecture and its core design, including the simulation environment, flight control, and communication system for efficient panoramic data sampling. Next, we describe the off-line data collection system, which generates panoramic images and pixel-level ground truth such as semantic segmentation, entity segmentation, and human keypoints using an automated trajectory generation tool. Finally, we present the interactive pedestrian-aware system, which enables realistic human behaviors and interactions within the simulation environment.

3.1. Overview

Built upon high-quality rendering capabilities of Unreal Engine, AirSim360 is an online closed-loop simulator compatible with multiple versions ranging from UE 4.27 to UE 5.6. By default, we adopt UE5 as our setting, considering its significant improvements in dynamic illumination, geometric detail, and overall scalability.

To improve user accessibility, we further optimize the communication module between virtual sensors and the UAV flight controller while providing an open external interface such as Vision-Language-Action (VLA) [19] for flexible integration. In the current system, the external model can return high-level control commands or a direct target position. These values are then parsed by our custom flight control module into the corresponding thrust and torque for each of the four rotors.

Compared with other platforms such as AirSim and UnrealZoo, we independently compile a flight control module and integrate it deeply into UE. This design allows us to simulate various types of drone through a simple user prompt that can enhance the generalization and flexibility of our system. More details of the implemen-

tation are provided in the appendix, as they fall beyond our primary focus on omnidirectional simulation.

3.2. Data Collection Toolkit

Apart from online simulation, AirSim360 provides a powerful toolkit for offline data collection to overcome scale limitations of existing omnidirectional datasets, particularly in pixel-level annotations that are closely related to semantic and geometric understanding.

3.2.1. Render-Aligned Data and Label Generation

Inspired by the stitching process of omnidirectional cameras, our data, whether the original images or their corresponding ground truth, are generated by stitching six cube-face views together. Each view corresponds to one of the six directions, namely front, back, left, right, up, and down, with a 90-degree field of view. Specifically, we obtain six non-overlapped cube images $I_c^{6 \times H_c \times W_c}$ as input and then produce an equirectangular projection map (ERP) $I_e^{H_e \times W_e}$. Here, H_c , W_c , H_e and W_e are the height and width of the cube and erp images, respectively. For clear illustration of such a process, we use the spherical projection which serves as a bridge between the input and output representations in Figure 3. The detailed equation is in our appendix.

In the following, we introduce the unique challenges posed by our stitching process for each data type, along with the corresponding solutions.

Panoramic Image: Capturing multimodal data simultaneously from six directions significantly increases GPU rendering load and storage pressure, resulting in substantial frame rate drops and difficulty in maintaining high throughput at high resolutions. For multi-view image rendering, we have redesigned the composition mechanism for six camera-associated images from ground up. After establishing corresponding camera rendering views, we utilize internal functions of Render Hardware Interface (RHI) to perform GPU-side texture resource copying, enabling one-time stitching of the six

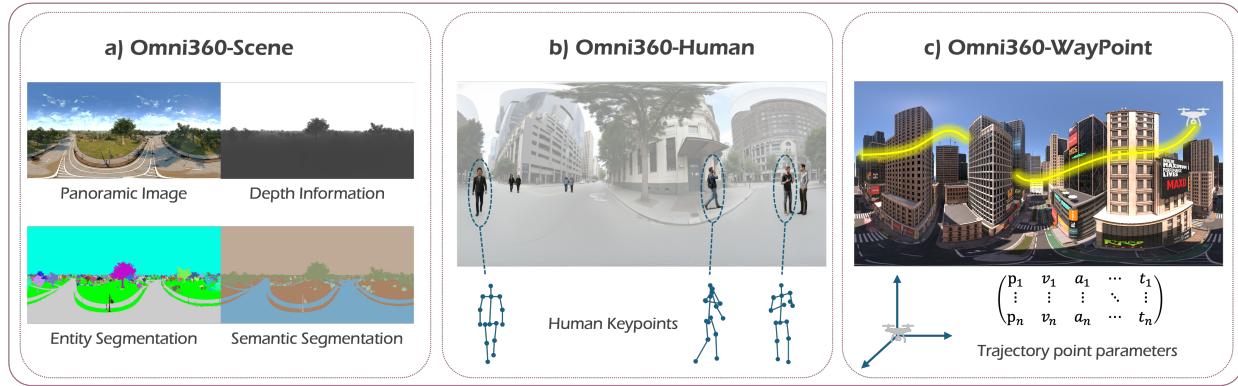


Figure 4. Visualization of Omni360-X. a) Visualizes Render-Aligned Data with corresponding labels (Sec. 3.2.1). b) Demonstrates pedestrian-aware label generation under panoramic views (Sec. 3.2.2). c) Shows trajectory synthesis from sparse waypoints using the Minimum Snap (Sec. 3.2.3).

249 images. This approach thus avoids the drawbacks of sec-
250 ondary stitching within material nodes via blueprints.
251 **Depth Information:** In a perspective image, depth is
252 defined as the distance between a 3D point and the cam-
253 era center along the optical axis. This definition differs
254 in the omnidirectional view, where no tangent plane ex-
255 exists. Therefore, in our setting, depth is defined as the
256 distance from the center of the camera to the point along
257 the viewing ray.

258 Since Z-Depth can be directly acquired from the pre-
259 computed depth buffer (Z-Buffer) in the Unreal Engine,
260 we proposed a material-based pipeline to extract depth
261 data and render the results into a Render Target. This ap-
262 proach allows writing the depth information directly into
263 the alpha channel of the corresponding image. During
264 the actual rendering process, with known camera world
265 coordinates and corresponding intrinsic and extrinsic pa-
266 rameters, the precise distance from the camera to each
267 spatial point can be computationally determined.

268 **Semantic Segmentation:** We assign a semantic label to
269 each pixel to enhance scene understanding. In our im-
270 plementation within Unreal Engine, we represent a set
271 of semantic categories by assigning specific RGB val-
272 ues. To achieve this, we utilize the Stencil Buffer from
273 the graphics rendering pipeline to assign colors to the
274 static mesh actuator. The Stencil Buffer stores an inte-
275 ger value between 0 and 255 for each pixel in the scene.
276 Through custom-designed post-process material, these
277 values are transformed into designated color outputs, fa-
278 cilitating the acquisition of semantic labels.

279 **Entity Segmentation:** Unlike other simulators that fo-
280 cus only on partial instances, we consider all entities
281 present in the scene, ensuring complete coverage across
282 the entire image. However, this design presents a chal-
283 lenge related to quantity limitations. The Stencil Buffer
284 mechanism restricts the total number of assignable cat-
285 egories to 256, which is insufficient for the large num-
286 ber of objects typically involved in entity segmentation

287 tasks, particularly in complex scenarios. To overcome
288 this limitation, we develop a dedicated entity segmen-
289 tation method capable of labeling all static mesh actors,
290 skeleton mesh actors, and landscape elements within the
291 scene. This approach allows us to obtain complete and
292 fine-grained segmentation results for any given frame.

293 **Synchronous Rendering among Various Sensors:**
294 Various sensors should capture data synchronously
295 without latency. Considering that UAV-mounted cam-
296 eras are in constant motion, we have deactivated the
297 *Capture Every Frame* option for all rendering cam-
298 eras, significantly reducing GPU resource consumption.
299 Since we have custom-designed the cameras and various
300 image sensors, we introduced an Event Dispatcher to
301 synchronize all sensors with a unified trigger signal, en-
302 abling simultaneous acquisition of multiple data types.

3.2.2. Interactive Pedestrian-Aware System

303 Scenarios involving humans are critically important in
304 low-altitude real-world environments. Therefore, sim-
305 ulating realistic human behavior is a key aspect of our
306 platform. We begin it from three perspectives. First, we
307 allow users to create a customizable number of pedes-
308 trians within a defined active area, automatically as-
309 signing various behaviors to them. Second, we enable
310 autonomous interactions by combining NPC Behavior
311 Trees with State Machines. For example, a pedestrian
312 state machine can trigger transitions based on a multi-
313 actor message dispatch and receive mechanism, allow-
314 ing an agent to switch from a walking state to a chat-
315 ting state when meeting another agent, or to randomly
316 activate an OnPhoneCall state while walking (see Fig-
317 ure 5). Third, human keypoints are generated in real
318 time as pedestrians interact, ensuring temporally consis-
319 tent annotations for downstream perception tasks.

320 Generating pedestrian keypoint data presents two
321 main challenges. First, a framework capable of support-
322 ing autonomous pedestrian motion is required to accom-

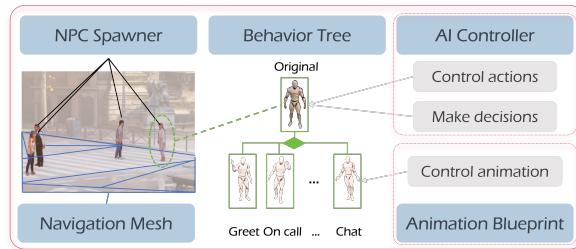


Figure 5. Conceptual diagram of the IPAS logic, showing the interplay between the Behavior Tree, State Machine, and the message-passing system for autonomous agent interaction.

modate diversity in body movements. Second, a method for binding identical skeletal keypoints across different characters is essential to avoid positional inaccuracies introduced by manual keypoint annotation.

To address these issues, we design Interactive Pedestrian-aware System (IPAS), a system that enables autonomous and interactive movements among pedestrians. Additionally, we implement a Blueprint-based approach that invokes pre-existing universal skeletal points in the Skeletal Mesh via blueprint functions. For keypoints not included in the standard skeletal framework, we use the *Add Socke* method to incorporate them into the Skeleton Tree, thereby enabling users to access and output all designed body keypoint coordinates.

3.2.3. Automated Trajectory Generation Paradigm

The ground truth trajectories in previous platforms have largely relied on human control of the drone. To make the sampling process more efficient, we adopt Minimum Snap trajectory planning [25] into our simulation pipeline. A collector only needs to specify a few key waypoints within the scene, after which the system automatically generates a smooth and realistic trajectory that adheres to the dynamic constraints of the UAV. By adjusting parameters such as maximum velocity and acceleration, the resulting polynomial coefficients of the trajectories can be applied directly to real quadrotors.

4. Datasets Collection

Based on our proposed AirSim360 simulator, we collect several large-scale datasets for omnidirectional scene understanding. Then, we introduce the datasets gathered from a variety of Unreal Engine 5 (UE5) scenarios. Specifically, Omni360-X is a large-scale panoramic dataset containing 60K nonduplicated frames filtered using SSCD [29]. To serve diverse research objectives, Omni360-X is organized into three subsets, including Scene, Pedestrian, and Trajectory, where each emphasizes a specific aspect of 360-degree understanding. The visualization is shown in Figure 4.

More detailed information about the Omni360-X dataset is provided in the supplementary material due

Table 2. Composition of the Omni360-Scene Dataset. The dataset includes four Unreal Engine scenes, each depicting a unique urban environment. For scenes of varying scales, we collected panoramic data with corresponding labels, while the final column reports the total number of labels per scene.

Scenarios Name	Area (m ²)	Nums	Label types	Sem. Cat.
<i>City Park</i>	800,000	25,600	Dep/Seg/Ins	25
<i>Downtown West</i>	60,000	6,800	Dep/Seg/Ins	29
<i>SF City</i>	250,000	22,000	Dep/Seg/Ins	20
<i>New York City</i>	44,800	6,600	Dep/Seg/Ins	25

to page limitations.

364

Omni360-Scene We focus primarily on scene parsing within Omni360-Scene. Inspired by ADE20K [44], Omni360-Scene contains more than 60K images annotated with depth and panoptic segmentation. Panoptic segmentation has two components: semantic and entity segmentation. For semantic labeling, we first extract semantic nouns from UE 5, which can be directly used for understanding open-world or open-vocabulary scenes. We then design a hierarchical semantic tree to ensure semantic consistency across different scenes. For entity segmentation, we further decompose stuff categories, such as trees and buildings, into individual entities. This is made possible by the simulator, which enables detailed entity separation that would be extremely challenging to achieve manually.

365
366
367
368
369
370
371
372
373
374
375
376
377
378
379

In Table 2, we present the data statistics for the four scenarios within Omni360-Scene. The dataset contains a total of 61,000 images captured from four scenarios, which demonstrate significant diversity in both physical scale (from 44,800 m² to 800,000 m²) and semantic complexity (20 to 29 categories). All scenarios provide comprehensive ground truth for Depth, Semantic Segmentation, and Instance Segmentation

380
381
382
383
384
385
386
387

Omni360-Human To better understand pedestrian behavior, we adopt 3D monocular human localization as our primary task, as it measures both human positions and postures in the 3D world, which are closely related to behavioral understanding. Accordingly, Omni360-Human contains approximately 100K samples covering more than 10 pedestrian behaviors across diverse camera distances and viewpoints in about 6 scenes.

388
389
390
391
392
393
394
395

As shown in Figure 4 (b), each sample contains both camera information, including the camera's absolute position and rotation angles in the world coordinate system, and pedestrian information, such as the locations and rotation angles of human keypoints.

396
397
398
399
400

Omni360-WayPoint Building upon our proposed automatic trajectory synthesis scheme, we will release

401
402

Table 3. Summary of the Omni360-WayPoint dataset. Flight paths of different lengths and sampling rates were generated across four outdoor scenes with physically consistent UAV dynamics.

Thumbnail	Scenario	Scene type	Count	Spacing	Length range
	<i>City Park</i>	outdoor	20 000	0.5	[50, 150]
	<i>Downtown West</i>	outdoor	5 000	0.2	[20, 50]
	<i>New York City</i>	outdoor	5 000	0.2	[20, 50]
	<i>SF City</i>	outdoor	20 000	0.5	[50, 150]

Table 4. Performance metrics for frame capture settings.

Capture Every Frame	FPS	GPU Time
Enable	20	54 ms
Disable	29	35 ms

403 Omni360-WayPoint, an open dataset containing over
404 100,000 UAV waypoints, as described in Table 3. The
405 trajectories adhere to realistic flight dynamics and in-
406 clude route variants parameterized by different maxi-
407 mum level-flight speeds, enabling analyses across di-
408 verse kinematic regimes. Such datasets can serve as
409 various targets. For example, they can provide physics-
410 consistent supervision for perception and state estima-
411 tion, and enable trajectory prediction and system iden-
412 tification for model-based control and reinforcement
413 learning. Moreover, they align instruction–video–action
414 for VLA training, strengthen 3D reconstruction and
415 mapping with accurate poses.

416 5. Experiments

417 In this section, we conduct several ablation studies on
418 the AirSim360 platform to demonstrate the effectiveness
419 of our design choices. We then present and benchmark
420 the sampled dataset in multiple tasks, indicating the ben-
421 efits of our proposed simulator for real-world scenarios.

422 5.1. Ablation Study on AirSim360

423 We conduct several comparison experiments on a work-
424 station equipped with an NVIDIA RTX 4060 Ti (8 GB),
425 an Intel i7-14700, and 32 GB of RAM. The evaluation
426 covers two operational modes: perspective camera ren-
427 dering and panoramic image stitching and transmission.
428 Across all settings, our system consistently achieves
429 higher data throughput and improved frame rates than
430 conventional approaches.

431 Each camera in the rendering engine is processed as
432 an independent rendering pass. Because the cameras
433 move in real time, we optimize their scene-scanning pat-

Table 5. Comparison of frame rates across different platforms with 6 cameras.

Platform	Nums of Camera	FPS
Not Optimized	6	14
Airsim360	6	18

Table 6. The MPDE results trained on two datasets, evaluated across three public benchmarks and one control set, *Omni360*. The *Omni360* refers to the dataset Omni360-Human.

Training Set	Test Set	Dist. Err	Ang. Err	Ang. Err (Pub)	Dist. Err (Pub)
<i>nuScenes</i>	<i>nuScenes</i>	1.078	31.90		
	<i>KITTI</i>	0.822	31.50	21.21	0.484
	<i>FreeMan</i>	0.260	17.00		
	<i>Omni360</i>	2.439	33.30		
<i>nuScenes</i> + <i>Omni360</i>	<i>nuScenes</i>	1.068	30.70		
	<i>KITTI</i>	0.809	31.20	17.02	0.458
	<i>FreeMan</i>	0.228	11.60		
	<i>Omni360</i>	1.779	15.20		

tern to improve efficiency. As shown in Table 4, this optimized scanning strategy yields a 45% increase in per-camera frame rate.

Given the panoramic design of our platform, we restructure the composition and data-transmission pipeline for the six camera feeds in the rendering engine. Leveraging the Render Hardware Interface (RHI) described in Section 3.2.1, we reimplement the C++ pipeline, increasing the frame rate from 14 to 18 FPS.

5.2. Monocular Pedestrian Distance Estimation

To assess the impact of synthetic data, we apply IPAS (view in Sec. 3.2.2) to Monocular Pedestrian Distance Estimation (MPDE) in Table 6. The generated Omni360-Human dataset provides 3D keypoint annotations in the world coordinate system. Using diverse simulated scenes for data augmentation, our approach achieves improved results on several MPDE benchmarks. The model is trained with Omni360-Human and nuScenes [5], and evaluated on Omni360-Human, nuScenes, KITTI [16], and FreeMan [35], reporting both Euclidean and angular distance errors following MonoLoco++ [4].

As shown in Table 6, training with the Omni360-Human dataset consistently improves performance across all three public test sets. The average angular error decrease from 21.21° to 17.02° , and the mean distance error from 0.484 m to 0.458 m. To better reflect real-world camera configurations, synthetic data are generated with a 20° pitch angle. The largest test set, FreeMan, exhibit the most significant improvement, with the angular error reduced from 11.6° to 11.6° . These

Table 7. Quantitative comparison for depth estimation. We fine-tune on the UniK3D model and compare the performance of training on Deep360 versus our Omni360 dataset. ↓ and ↑ denote lower or higher is better. Best results are in bold.

Experiment Type	Training Data	Evaluation Data	AbsRel ↓	RMSE ↓	$\delta_1 \uparrow$
Out-of-Domain	Deep360	SphereCraft	8.2570	0.0566	0.3490
	Omni360 (Ours)	SphereCraft	5.4372	0.0435	0.3990
Cross-Domain	Deep360	Omni360 (Ours)	0.3600	0.0714	0.4896
	Omni360 (Ours)	Deep360	0.1762	0.0229	0.6672

Table 8. Results of panoramic segmentation. The mIoU and mAP metrics are used to evaluate performance on the semantic and entity segmentation tasks, respectively.

Task	WildPASS	Omni360-Scene	Performance
Semantic	✓	○	58.0
	✓	✓	67.4
Entity	✓	○	24.6
	✓	✓	38.9

results highlight the effectiveness of synthetic data in enhancing monocular pedestrian distance estimation.

5.3. Panoramic Depth Estimation

To evaluate the generalization and domain adaptability of our Omni360 dataset, we conduct out-of-domain and cross-domain experiments for panoramic depth estimation. All datasets consist of outdoor synthetic data. The UniK3D model [27] is used as the baseline and fine-tuned after its official implementation. We compare the results with the Deep360 dataset [22], using Absolute Relative Error(AbsRel), Root Mean Squared Error(RMSE), and a percentage metrics δ_1 , where $i = 1.25$, as evaluation metrics.

Out-of-Domain: Models trained in Deep360 and Omni360 are tested on SphereCraft [15]. In Table 7, the model trained on Omni360 achieves better results, showing improved generalization to unseen environments.

Cross-Domain: We perform cross-domain tests between Deep360 and Omni360. The model trained on Omni360 performs best, indicating stronger robustness and more transferable representations.

5.4. Panoramic Segmentation

In Table 8, we evaluate semantic and entity segmentation in the WildPASS validation set [41] using OOOPS [42] and Mask2Former [9]. With the incorporation of our Omni360-Scene segmentation data, both tasks achieve substantial performance gains, highlighting the effectiveness of large-scale data scaling for pixel-level semantic understanding.

5.5. Benchmark on Omni360-X

5.5.1. Panoramic Vision-Language Navigation

Vision-Language Navigation (VLN) [1] for unmanned aerial vehicles (UAVs) represents a challenging yet crucial direction toward embodied intelligence in aerial

Table 9. Results of panoramic VLN. Three standard metrics, Success Rate (SR), Success weighted by Path Length (SPL), and Navigation Error (NE), are used to evaluate model performance.

Model	SR	SPL	NE
<i>qwen2.5-vl-72b-instruct</i>	0.4	0.3843	18099.73
<i>qwen3-vl-plus</i>	0.0	0.0	11436.97
<i>qwen3-vl-flash</i>	0.2	0.1945	9506.26
<i>doubaos-1-6-251015</i>	0.5	0.4813	10573.89

agents, , requiring joint understanding of spatial scenes and language instructions. Most existing UAV-VLN methods mostly rely on forward-facing cameras, resulting in a narrow field of view and large blind areas that limit target search efficiency and situational awareness.

We introduce a panoramic UAV-VLN task built on our simulation platform, where the UAV perceives the environment through panoramic vision instead of a single forward camera.

A set of preliminary experiments is conducted using several vision–language models to assess their ability to interpret panoramic observations for navigation.

In this section, to demonstrate the You Only Move Once (YOMO) capability of panoramic UAVs in decision-making tasks, we place most targets as salient objects and designed short flight trajectories. Under these conditions, the panoramic UAV can directly move toward the goal without additional yaw rotations or exploration. The similar SPL and SR values reported in Table 9 confirm that panoramic perception enables near-optimal, single-step decision making when the UAV is close to the target, validating the YOMO principle.

6. Conclusion

To address the issue of lacking large-scale omnidirectional data, we propose AirSim360, a panoramic simulation platform built on Unreal Engine. Inspired by the key elements in 4D real world, AirSim360 focuses on three aspects. First, we introduce a render-aligned data pipeline for generating 360-degree images with pixel-level depth and segmentation annotations in ERP representation. Second, we develop an interactive pedestrian-aware system that enables the study of human behavior through 3D monocular human localization. Third, we present an automated trajectory-generation paradigm that produces realistic aerial trajectories for navigation tasks. Built on top of AirSim360, we collect more than 60K non-duplicate frames. Ablation studies on the UE 5-based design, along with extensive experiments on the collected data, demonstrate the effectiveness of our platform and the benefits of the resulting dataset.

References

- [1] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton Van Den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3674–3683, 2018. 8
- [2] Dragomir Anguelov, Carole Dulong, Daniel Filip, Christian Frueh, Stéphane Lafon, Richard Lyon, Abhijit Ogale, Luc Vincent, and Josh Weaver. Google street view: Capturing the world at street level. *Computer*, 43(6):32–38, 2010. 3
- [3] Iro Armeni, Sasha Sax, Amir R Zamir, and Silvio Savarese. Joint 2d-3d-semantic data for indoor scene understanding. *arXiv preprint arXiv:1702.01105*, 2017. 3
- [4] Lorenzo Bertoni, Sven Kreiss, and Alexandre Alahi. Perceiving humans: from monocular 3d localization to social distancing. *IEEE Transactions on Intelligent Transportation Systems*, 2021. 7
- [5] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020. 7
- [6] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. *arXiv preprint arXiv:1709.06158*, 2017. 3
- [7] Shih-Hsiu Chang, Ching-Ya Chiu, Chia-Sheng Chang, Kuo-Wei Chen, Chih-Yuan Yao, Ruen-Rone Lee, and Hung-Kuo Chu. Generating 360 outdoor panorama dataset with reliable sun position estimation. In *SIGGRAPH Asia*. 2018. 3
- [8] Zhaoxi Chen, Guangcong Wang, and Ziwei Liu. Text2light: Zero-shot text-driven hdr panorama generation. *ACM Transactions on Graphics (TOG)*, 41(6):1–16, 2022. 3
- [9] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *CVPR*, 2022. 8
- [10] Shih-Han Chou, Cheng Sun, Wen-Yen Chang, Wan-Ting Hsu, Min Sun, and Jianlong Fu. 360-indoor: Towards learning real-world objects in 360° indoor equirectangular images. In *WACV*, 2020. 3
- [11] Xin Deng, Hao Wang, Mai Xu, Yichen Guo, Yuhang Song, and Li Yang. Lau-net: Latitude adaptive upscaling network for omnidirectional image super-resolution. In *CVPR*, 2021. 3
- [12] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla: An open urban driving simulator. In *Conference on robot learning*, 2017. 2
- [13] Yue Fan, Winson Chen, Tongzhou Jiang, Chun Zhou, Yi Zhang, and Xin Wang. Aerial vision-and-dialog navigation. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 3043–3061, 2023. 3
- [14] Yunpeng Gao, Chenhai Li, Zhongrui You, Junli Liu, Zhen Li, Pengan Chen, Qizhi Chen, Zhonghan Tang, Liansheng Wang, Penghui Yang, et al. Openfly: A comprehensive platform for aerial vision-language navigation. In *arXiv*, 2025. 2
- [15] Christiano Gava, Yunmin Cho, Federico Raue, Sebastian Palacio, Alain Pagan, and Andreas Dengel. Spherecraft: A dataset for spherical keypoint detection, matching and camera pose estimation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 4408–4417, 2024. 8
- [16] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3354–3361. IEEE, 2012. 7
- [17] Yannick Hold-Geoffroy, Akshaya Athawale, and Jean-François Lalonde. Deep sky modeling for single image outdoor lighting estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6927–6935, 2019. 3
- [18] Wouter Jansen, Erik Verreycken, Anthony Schenck, Jean-Edouard Blanquart, Connor Verhulst, Nico Huebel, and Jan Steckel. Cosys-airsim: a real-time simulation framework expanded for complex industrial applications. *arXiv preprint arXiv:2303.13381*, 2023. 2
- [19] Sicong Jiang, Zilin Huang, Kangan Qian, Ziang Luo, Tianze Zhu, Yang Zhong, Yihong Tang, Menglin Kong, Yunlong Wang, Siwen Jiao, et al. A survey on vision-language-action models for autonomous driving. *arXiv preprint arXiv:2506.24044*, 2025. 4
- [20] Jungdae Lee, Taiki Miyanishi, Shuhei Kurita, Koya Sakamoto, Daichi Azuma, Yutaka Matsuo, and Nakama Inoue. Citynav: Language-goal aerial navigation dataset with geographic information. *arXiv preprint arXiv:2406.14240*, 2024. 3
- [21] Chengshu Li, Ruohan Zhang, Josiah Wong, Cem Gokmen, Sanjana Srivastava, Roberto Martín-Martín, Chen Wang, Gabrael Levine, Wensi Ai, Benjamin Martinez, Hang Yin, Michael Lingelbach, Minjune Hwang, Ayano Hiranaka, Sujay Garlanka, Arman Aydin, Sharon Lee, Jiankai Sun, Mona Anvari, Manasi Sharma, Dhruva Bansal, Samuel Hunter, Kyu-Young Kim, Alan Lou, Caleb R Matthews, Ivan Villa-Renteria, Jerry Huayang Tang, Claire Tang, Fei Xia, Yunzhu Li, Silvio Savarese, Hyowon Gweon, C. Karen Liu, Jiajun Wu, and Li Fei-Fei. Behavior-1k: A human-centered, embodied ai benchmark with 1,000 everyday activities and realistic simulation. *arXiv preprint arXiv:2403.09227*, 2024. 2
- [22] Ming Li, Xueqian Jin, Xuejiao Hu, Jingzhao Dai, Sidan Du, and Yang Li. Mode: Multi-view omnidirectional depth estimation with 360° cameras. In *European Conference on Computer Vision (ECCV)*, 2022. 8
- [23] Xin Lin, Xian Ge, Dizhe Zhang, Zhaoliang Wan, Xianshun Wang, Xiangtai Li, Wenjie Jiang, Bo Du, Dacheng Tao, Ming-Hsuan Yang, and Lu Qi. One flight over

- 655 the gap: A survey from perspective to panoramic vision,
656 2025. 3
- 657 [24] Shubo Liu, Hongsheng Zhang, Yuankai Qi, Peng Wang,
658 Yanning Zhang, and Qi Wu. Aerialvln: Vision-and-
659 language navigation for uavs. In *Proceedings of the*
660 *IEEE/CVF International Conference on Computer Vi-*
661 *sion*, pages 15384–15394, 2023. 3
- 662 [25] Daniel Mellinger and Vijay Kumar. Minimum snap
663 trajectory generation and control for quadrotors. In *2011*
664 *IEEE international conference on robotics and automa-*
665 *tion*, pages 2520–2525. IEEE, 2011. 6
- 666 [26] Dipendra Misra, Andrew Bennett, Valts Blukis, Eyyvind
667 Niklasson, Max Shatkhnin, and Yoav Artzi. Mapping in-
668 structions to actions in 3d environments with visual goal
669 prediction. *arXiv preprint arXiv:1809.00786*, 2018. 3
- 670 [27] Luigi Piccinelli, Christos Sakaridis, Mattia Segu, Yung-
671 Hsu Yang, Siyuan Li, Wim Abbeloos, and Luc Van
672 Gool. Unik3d: Universal camera monocular 3d estimation,
673 2025. 8
- 674 [28] Giovanni Pintore, Eva Almansa, Marco Agus, and Enrico
675 Gobbetti. Deep3dlayout: 3d reconstruction of an indoor
676 layout from a spherical panoramic image. *TOG*, 2021. 3
- 677 [29] Ed Pizzi, Sreya Dutta Roy, Sugosh Nagavara Ravindra,
678 Priya Goyal, and Matthijs Douze. A self-supervised de-
679 scriptor for image copy detection. In *CVPR*, 2022. 6
- 680 [30] Xavier Puig, Eric Undersander, Andrew Szot,
681 Mikael Dallaire Cote, Tsung-Yen Yang, Ruslan Partsey,
682 Ruta Desai, Alexander William Clegg, Michal Hlavac,
683 So Yeon Min, et al. Habitat 3.0: A co-habitat for humans,
684 avatars and robots. *arXiv preprint arXiv:2310.13724*,
685 2023. 3
- 686 [31] Weichao Qiu, Fangwei Zhong, Yi Zhang, Siyuan Qiao,
687 Zihao Xiao, Tae Soo Kim, and Yizhou Wang. Unrealcv:
688 Virtual worlds for computer vision. In *ACM MM*, 2017.
689 3
- 690 [32] Shital Shah, Debadatta Dey, Chris Lovett, and Ashish
691 Kapoor. Airsim: High-fidelity visual and physical sim-
692 ulation for autonomous vehicles. In *Field and Service*
693 *Robotics*, 2017. 3
- 694 [33] Shital Shah, Debadatta Dey, Chris Lovett, and Ashish
695 Kapoor. Airsim: High-fidelity visual and physical sim-
696 ulation for autonomous vehicles. In *Field and service*
697 *robotics: Results of the 11th international conference*,
698 2017. 2
- 699 [34] Oleksandra Sobchyshak, Santiago Berrezueta-Guzman,
700 and Stefan Wagner. Pushing the boundaries of immer-
701 sion and storytelling: A technical review of unreal en-
702 gine. *Displays*, page 103268, 2025. 2
- 703 [35] Jiong Wang, Fengyu Yang, Bingliang Li, Wenbo Gou,
704 Danqi Yan, Ailing Zeng, Yijun Gao, Junle Wang, Yan-
705 qing Jing, and Ruimao Zhang. Freeman: Towards bench-
706 marking 3d human pose estimation under real-world con-
707 ditions. In *Proceedings of the IEEE/CVF Conference on*
708 *Computer Vision and Pattern Recognition*, pages 21978–
709 21988, 2024. 7
- 710 [36] Sijie Wang, Siqi Li, Yawei Zhang, Shangshu Yu, Sheng-
711 hai Yuan, Rui She, Quanjiang Guo, JinXuan Zheng,
712 Ong Kang Howe, Leonrich Chandra, et al. Uavscenes:
713 A multi-modal dataset for uavs. In *ICCV*, 2025. 3
- [37] Xiangyu Wang, Donglin Yang, Ziqin Wang, Hohin
714 Kwan, Jinyu Chen, Wenjun Wu, Hongsheng Li, Yue
715 Liao, and Si Liu. Towards realistic uav vision-language
716 navigation: Platform, benchmark, and methodology.
717 *arXiv preprint arXiv:2410.07087*, 2024. 3
- [38] Hang Xu, Qiang Zhao, Yike Ma, Xiaodong Li, Peng
719 Yuan, Bailan Feng, Chenggang Yan, and Feng Dai. Pan-
720 dora: A panoramic detection dataset for object with ori-
721 entation. In *European conference on computer vision*,
722 pages 237–252. Springer, 2022. 3
- [39] Menghan Xu et al. Predicting head movement in
724 panoramic video: A deep reinforcement learning ap-
725 proach. In *CVPR*, 2018. 3
- [40] Kailun Yang, Xinxin Hu, Luis M Bergasa, Eduardo
727 Romera, and Kaiwei Wang. Pass: Panoramic annular se-
728 mantic segmentation. *IEEE Transactions on Intelligent*
729 *Transportation Systems*, 21(10):4171–4185, 2019. 3
- [41] Kailun Yang, Jiaming Zhang, Simon Reiß, Xinxin Hu,
731 and Rainer Stiefelhagen. Capturing omni-range context
732 for omnidirectional segmentation. In *CVPR*, 2021. 3, 8
- [42] Junwei Zheng, Ruiping Liu, Yufan Chen, Kunyu Peng,
734 Chengzhi Wu, Kailun Yang, Jiaming Zhang, and Rainer
735 Stiefelhagen. Open panoramic segmentation. In *ECCV*,
736 2024. 8
- [43] Fangwei Zhong, Kui Wu, Churan Wang, Hao Chen, Hai
738 Ci, Zhoujun Li, and Yizhou Wang. Unrealzoo: Enriching
739 photo-realistic virtual worlds for embodied ai. In *ICCV*,
740 2025. 2
- [44] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja
742 Fidler, Adela Barriuso, and Antonio Torralba. Seman-
743 tic understanding of scenes through the ade20k dataset.
744 *International Journal of Computer Vision*, 127(3):302–
745 321, 2019. 6
- [746]