

# HE-Nav: A High-Performance and Efficient Navigation System for Aerial-Ground Robots

Anonymous Review. Paper-ID [123]

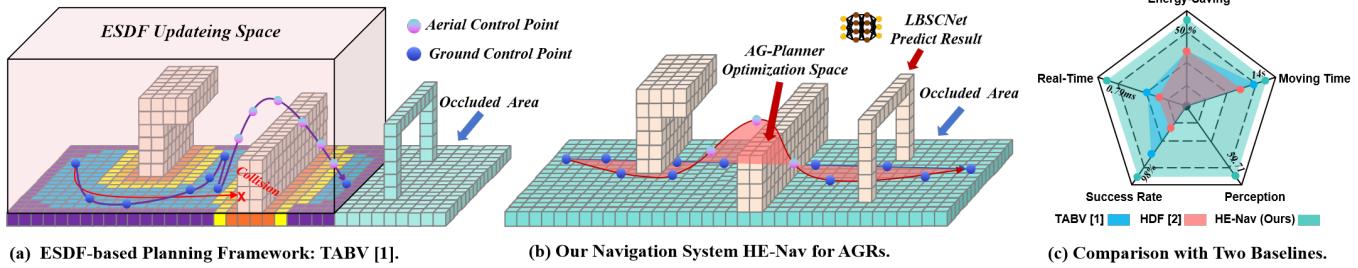


Fig. 1: (a) ESDF-based methods face the dual predicaments of reduced path planning performance and increased energy consumption. (b) Our HE-Nav system can generate energy-saving, collision-free aerial-ground paths in real-time with the help of the LBSCNet model and AG-Planner. (c) HE-Nav Comparison with two baselines.

**Abstract**—Aerial-ground robots (AGRs) have unique dual-mode capabilities (i.e., flying and driving), making them well-suited for search and rescue tasks. While existing Euclidean Signed Distance Field (ESDF)-based navigation systems have made progress in structured indoor scenarios by constructing an ESDF map to search collision-free hybrid paths, these systems often show suboptimal performance and energy consumption in complex, occluded settings (e.g., forests) due to their inability to perceive unknown areas that are occluded resulting in incomplete local maps coupled with inherent limitations of the path planner.

In this paper, we present HE-Nav, the first high-performance and energy-efficient navigation system tailored for AGRs. This novel system incorporates a lightweight Bird's Eye View (BEV)-guided semantic scene completion network called LBSCNet within its perception module. Assisted by an exquisitely designed SCB-Fusion module and attention mechanisms that facilitate real-time prediction of obstacle distributions within occluded regions ahead of time. Furthermore, HE-Nav introduces an advanced AGR path planner which integrates both gradient-based path optimiser and Kinodynamic A\* algorithms to produce safe and energy-saving trajectories.

Extensive simulations and real-world experiments demonstrate that HE-Nav significantly outperforms two recent AGR navigation systems, achieving up to 50% reduction in overall energy consumption while maintaining superior performance metrics including a remarkable 98% success rate alongside shorter average planning times and moving times. The code and hardware configuration will be released.

## I. INTRODUCTION

In recent years, aerial-ground robots (AGRs) [1, 2, 3, 4] have emerged as a promising solution for search [5, 6], exploration [7, 8], and rescue tasks [9, 10]. This is attributed to their exceptional mobility and long endurance, which enable them to seamlessly switch between aerial and ground modes, allowing for hybrid locomotion (i.e., flying and driving) in the above challenging tasks. As a key enabler for autonomous

navigation in AGRs, the perception module provides a local map to the path planner, which in turn searches for aerial-ground hybrid trajectories that are both high performance (i.e., planning success rate and moving times) and efficient (i.e., real-time and lower energy consumption).

Existing ESDF-based navigation systems [1, 4] utilize sensors (e.g., cameras) to perceive surrounding environments and establish Euclidean Signed Distance Field (ESDF) maps (in Fig. 1a), subsequently the path planner to search for collision-free trajectories that favour ground paths and only switch to the aerial mode when necessary (e.g., encountering impassable obstacles), thereby promoting energy efficiency. Unfortunately, While these methods have proven successful in structured indoor scenarios, they face two limitations when navigating in complex and occluded environments (e.g., disaster zones or wilderness areas).

Firstly, limited perception results in incomplete local maps (i.e., containing occlusion-induced unknown areas) due to the narrow field of view in sensor-based mapping. This not only raises collision risks (i.e., *red path* in Fig. 1a.) but also prolongs movement time with redundant paths (i.e., *purple path* in Fig. 1a). To improve the performance of AGRs during navigation, semantic scene completion (SSC) networks [11] show promise for predicting obstacle distribution and semantics in occluded areas. However, existing networks struggle to strike a balance between completion accuracy and real-time inference. Some use 3D convolution [12, 13] to improve accuracy but unsuitable for resource-limited AGR devices to ensure real-time inference. Others propose lightweight network structures [14] but with significantly reduced accuracy.

Secondly, the existing path planners are inefficient. Specifically, building the ESDF maps generates redundant calculation

times that do not meet the real-time requirements of path planning since it takes up about 70% of the time [15], and obstacles only take up 30% of the entire space (in Fig. 1a). Moreover, while the energy costs of flying are typically accounted for in planning, the path-searching algorithms often overlook the energy implications of ground movements, such as steering adjustments. This oversight results in overall energy consumption being inefficient.

In this study, our key insight to tackle these limitations is the co-design of a novel lightweight perception module and efficient path planner. The former achieves accurate completion and real-time reasoning to construct a complete local map, which ensures the latter minimizes movement time. Moreover, By improving the gradient calculation method and path search algorithm, we also reduce planning time and energy consumption.

Based on our insights into the perception and planning modules, we present **HE-Nav**, the first *high-performance* and *efficient* navigation system tailored for AGRs, as illustrated in Fig. 2. The system consists of three asynchronous modules: perception, planning, and control. Initially, the perception module employs the lightweight LBSCNet network for real-time 3D scene completion and semantic predictions. Subsequently, these results are seamlessly updated to the local map for path planning. In the planning phase, Our AGR motion planner (AG planner) first generates an initial path (i.e., *blue path* in Fig.1b.), then uses a gradient-based spline optimizer and a post-refinement procedure to optimize the air-ground trajectory, ultimately obtaining an energy-saving, safe, smooth and dynamically feasible trajectory (i.e., *brown path* in Fig.1b.). This trajectory will be sent to the controller for precise tracking.

However, the design of the perception module and path planner is confronted with two major challenges. The first challenge involves designing a lightweight SSC network using a sparse convolution that can ensure real-time reasoning on AGR devices. However, such a network may lack the ability to capture contextual information or detailed features, particularly in occlusion areas, which leads to a significant drop in completion accuracy. Moreover, dense feature fusion in 3D space continues to impede the real-time performance of reasoning. High-latency reasoning hinders timely local map updates, which in turn adversely affects path planning performance. The second challenge arises from the fact that while Zhou et al. [15] have developed an ESDF-free path planner specifically for quadcopters, it does not cater to the unique requirements of AGRs, particularly concerning energy efficiency (i.e., flight-centric trajectory generation results in higher energy consumption) and dynamic constraints.

To tackle these challenges, our LBSCNet decouple the learning of semantics and geometry into two distinct branches. This separation is crucial, as semantic context and geometric structure are complementary in SSC tasks [16]. By integrating attention mechanisms, not only is the network's ability to learn semantics and geometry accelerated, but it also captures rich and dense contextual information as well as features of oc-

cluded areas. For feature fusion, inspired by [17], we migrated feature fusion to BEV space and proposed the innovative component SCB-Fusion to fuse BEV features, semantic and geometric features in BEV space. This fusion method not only reduces the complexity of calculations but also improves the final completion accuracy. Subsequently, for the path palwner, we design a novel path search algorithm that adds extra energy cost to motion primitives that require sharp turns on the ground or have destinations in the air to generate energy-efficient collision-free guidance paths. In addition, considering the ground movement of AGRs faced non-holonomic constraints, we enforced a cost on ground control points to limit the curvature of the terrestrial trajectory.

We first evaluated the LBSCNet network on the SemanticKITTI benchmark, comparing its completion accuracy and inference speed to a state-of-the-art SSC network. Next, we tested HE-Nav in simulated and real indoor and outdoor environments, comparing its performance and energy consumption against two baseline approaches (TABV [1] and HDF [2]). Our evaluation shows that:

- **HE-Nav is high-performance.** HE-Nav achieves a 98% success rate in complex and occluded simulation environments, with a reduced average moving time ( $\approx 12s$ ).
- **HE-Nav is planning-efficiency.** AG-planner demonstrates an 8x improvement in planning time compared to the ESDF-based method.
- **HE-Nav is energy-efficiency.** Utilizing a novel energy-efficient Kinodynamics A\* algorithm, AG-planner results in a 50
- **LBSCNet is accurate and real-time.** LBSCNet enables real-time inference (20.08 FPS) and low-latency updates, achieving state-of-the-art performance (IoU = 59.71) on the SemanticKITTI benchmark.

Our main contributions are LBSCNet, a novel lightweight SSC network and a novel hierarchical path planner (i.e., AG-Planner) tailored for AGRs. The former, with its lightweight network structure and the incorporation of innovative components such as SCB-Fusion and criss-cross attention, can predict obstacle distribution in occluded areas in real-time. The latter uses a gradient-based method to achieve ESDF-free, which significantly reduces the planning time. It is assisted by the energy-saving Kinodynamic A\* algorithm to ensure energy saving of the overall path. In addition, considering that AGRs face inherent non-integrity constraints in ground motion, a cost is also added to the ground control points to limit curvature. Ultimately our HE-Nav produces energy-efficient, safe, smooth, and dynamically feasible hybrid trajectories.

## II. RELATED WORK

### A. Motion Planning for AGRs

Numerous researchers have explored various aerial-ground robot configurations, such as incorporating passive wheels [1, 18, 3, 8, 4], cylindrical cages [19], or multi-limb [20] onto drones, while others [9, 6, 5, 10, 21] have integrated rotors with wheeled robots to achieve dual-mode (i.e., flying

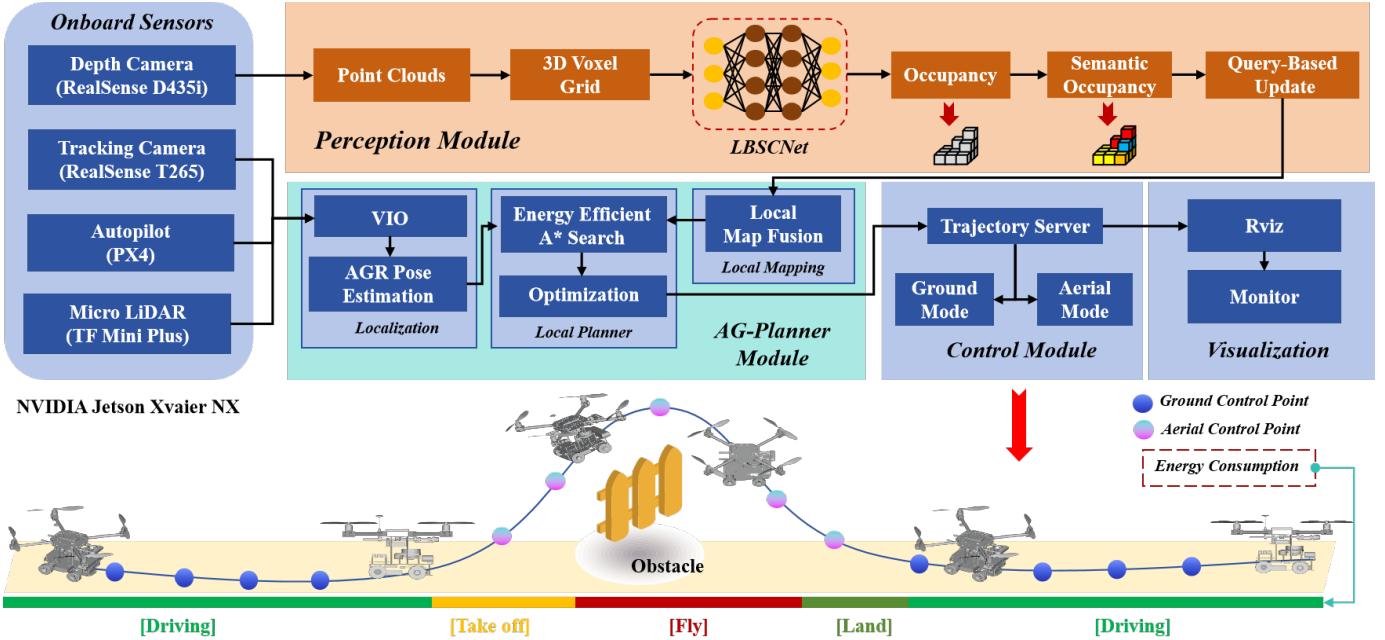


Fig. 2: HE-Nav system architecture. The perception, planning and control modules are deployed on the onboard computer (i.e., NVIDIA Jetson Xavier NX) and run asynchronously.

and driving) locomotion. These designs facilitate enhanced stability and control in both locomotion modes. Consequently, we have adopted a second mechanical structure to further customize our Aerial-Ground Robotic (AGR) system, which has four wheels and four rotors.

Although existing research primarily focuses on innovative mechanical structure designs, the area of AGR autonomous navigation remains underexplored. To the best of our knowledge, *Fan et al.* [2] address terrestrial-aerial motion planning. Their approach initially employs the A\* algorithm to search for a geometric path as guidance, favouring terrestrial paths by adding extra energy costs to aerial nodes. However, this path-searching method is limited due to its lack of dynamic models. Additionally, the local planner's trajectories lack post-refinement, resulting in potential issues with smoothness and dynamic feasibility. *Zhang et al.* [1] proposed a path planner and controller capable of efficient and adaptive path searching, but it relies on an ESDF map. The intensive computation and limited perception of occluded areas lead to a low success rate in path planning and increased energy consumption.

In our navigation system, we propose a new motion planner, namely AG-Palnner, which focuses on estimating the gradient of collision trajectory segments without building ESDF maps, resulting in a significant reduction in computational effort. It also contains an energy-efficient Kinodynamic A\* algorithm that searches for energy-efficient guidance paths, and finally models collision, smoothing and dynamic feasibility through a gradient-based spline optimizer to obtain the optimal trajectory, assisted by a post-refinement process to further improve the trajectory of robustness. In addition, for the dynamic constraints of AGR itself, we also added the curvature limit cost of the ground trajectory in the optimization formula to

deal with non-holonomic constraints.

### B. Occlusion-aware for AGRs

In recent years, the field of semantic scene completion has witnessed significant advancements, particularly in addressing the challenges posed by the obstruction of obstacles in complex and unknown environments. These advancements have led to the development of various point-cloud-based and camera-based methods for predicting obstacle distribution in occluded areas.

In the realm of camera-based methods, *Cao et al.* [12] introduced MonoScene, a groundbreaking approach that infers scene structure and semantics from a single monocular RGB image. Their key contribution lies in a novel 2D-3D feature projection bridging, inspired by optics, that enforces spatial semantic consistency. Another notable work by *Li et al.* [13] is VoxFormer, a Transformer-based semantic scene completion framework capable of generating complete 3D volume semantics using only 2D images.

On the other hand, point-cloud-based methods have also made significant strides. *Cheng et al.* [22] proposed S3CNet, a method designed to predict semantically complete scenes from a single unified LiDAR point cloud. *Roldao et al.* [14] introduced LMSCNet, a multiscale 3D semantic scene completion approach that uses a 2D UNet backbone with comprehensive multiscale skip connections to enhance feature flow, along with a 3D segmentation head. *Xia et al.* [11] devised a method that employs knowledge distillation from a multi-frame model to improve the performance of single-frame semantic scene completion. Lastly, *Zuo et al.* [23] proposed PointOcc, which introduces a cylindrical three-perspective view for effective

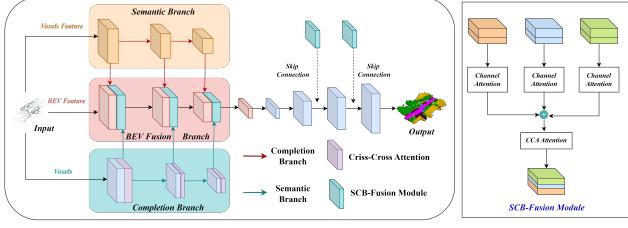


Fig. 3: The overview of the proposed LBSCNet.

and comprehensive representation of point clouds, along with a PointOcc model for efficient processing.

Despite the remarkable advancements in camera-based and point-cloud-based methods for semantic scene completion, these approaches often demand significant computational resources, rendering them unsuitable for real-time execution on resource-constrained robotic platforms. To address this limitation, we propose a lightweight semantic scene completion network guided by Bird's Eye View (BEV) features, which serves as the perception module for the HE-Nav navigation system. This network efficiently predicts the distribution of obstacles in occluded areas, ensuring seamless navigation in complex environments while maintaining low computational overhead, making it an ideal solution for resource-limited robotic devices.

### C. Energy-Efficient for AGRs

Energy efficiency is vital for aerial-ground robots since it directly impacts their endurance and overall performance. By utilizing energy efficiently, these robots can operate for extended periods, reducing downtime and optimizing flight and ground navigation.

Although the path planning frameworks proposed by *Fan et al.* [2] and *Zhang et al.* [1] take into account the energy consumption of AGRs, their approach is not comprehensive enough. They primarily add additional penalties to aerial flight, encouraging the robot to favour ground paths. While this strategy somewhat reduces energy consumption, it overlooks other factors that contribute to energy inefficiency. For instance, when moving on the ground, the robot's turning angle and travelling speed can lead to additional energy consumption. Frequent steering demands extra energy from the motor, resulting in suboptimal energy usage.

## III. PERCEPTION MODULE OF HE-NAV

In this section, we introduce a lightweight BEV-guided three-branch SSC network (LBSCNet), depicted in Fig. 3. LBSCNet consists of a semantic branch, a completion branch, and a BEV fusion branch, serving as an alternative to conventional memory-intensive SSC networks that concurrently predict geometry and semantics. By employing pre-trained models offline on ARG devices, LBSCNet is capable of real-time obstacle distribution prediction in occluded areas. Subsequently, these predictions are integrated into a local map, which is utilized for efficient path planning.

### A. LBSCNet Network Structure

LBSCNet decoupling the learning process of semantics and completion (or geometry), allows the network to concentrate on specific features (i.e., semantics and geometry), resulting in more efficient learning. The specific structures are as follows:

**Semantic Branch.** Point clouds  $P \in \mathbb{R}^{n \times 3}$  undergo processing via a voxelization layer to extract voxel features. Initially, the point cloud is partitioned based on the voxel resolution  $s$ . Points are mapped into voxel space, and their features are aggregated using an aggregation function (e.g., the max function) to obtain a single voxel feature. A multi-layer perceptron (MLP) is then utilized to reduce the dimensionality of this feature vector, yielding the final voxel features  $V_f$  with a dense spatial resolution of  $L \times W \times H$ .

Upon completing voxelization and entering the semantic branch, voxel features  $V_f$  are fed into three sparse encoder blocks to obtain sparse semantic features  $\{Sem_f^1, Sem_f^2, Sem_f^3\}$ . Each sparse encoder block comprises a residual block [24] with sparse convolutions and an SGFE module developed in [25]. The integration of the SGFE module not only enriches voxel features through multi-scale sparse projection and attention mechanisms, capturing both local and global features but also alleviates the computational burden by reducing feature resolution. The semantic branch is optimized using lovasz loss [26] and cross-entropy loss [27]. The semantic loss  $L_s$  is the sum of the loss at each stage, expressed as follows:

$$L_{sem} = \sum_{i=1}^3 (L_{lovasz,i} + L_{ce,i}) \quad (1)$$

**Completion Branch.** The completion branch takes occupancy voxels  $V \in \mathbb{R}^{1 \times L \times W \times H}$  generated by the depth camera point cloud, indicating whether the voxels are occupied. This branch outputs multi-scale dense completion features  $\{Com_f^1, Com_f^2, Com_f^3\}$ , providing more intricate geometric information. As depicted in Fig. 3, the completion branch comprises three residual blocks and three GPU memory-efficient criss-cross attention modules. The residual blocks incorporate dense 3D convolutions with a kernel size of  $3 \times 3 \times 3$ , capturing local geometric nuances. Conversely, the criss-cross attention module is designed to capture long-range dependencies by gathering contextual information in both horizontal and vertical directions, thereby enriching the completion features with a global context. The training loss  $L_c$  for this branch is calculated as follows:

$$L_{com} = \sum_{i=1}^3 (L_{lovasz,i} + L_{bce,i}) \quad (2)$$

**BEV Feature Fusion Branch.** Prior research has utilized dense 3D convolutions to achieve semantic scene completion by fusing dense 3D features. However, this approach is memory-intensive and requires significant GPU resources, rendering it impractical for deployment on resource-constrained robotic devices. Leveraging recent advancements in BEV perception, we propose a lightweight BEV feature fusion branch

for SSC tasks. By projecting learned semantic and geometric features into the BEV space for fusion, the computational overhead is substantially reduced, enhancing scene completion performance and ensuring real-time inference capabilities.

To project the three-dimensional semantic features  $\{Sem_f^1, Sem_f^2, Sem_f^3\}$  into the two-dimensional BEV space, we generate the BEV index based on the voxel index. Features sharing the same BEV index are aggregated using an aggregation function (i.e., the max function) to obtain sparse BEV features. Employing the feature densification function provided by spconv [28], dense BEV features  $\{B_f^{sem,0}, B_f^{sem,1}, B_f^{sem,2}, B_f^{sem,3}\}$  are generated based on the BEV index and sparse BEV features.

For geometric features  $\{Com_f^1, Com_f^2, Com_f^3\}$ , we stack dense 3D features along the  $z$ -axis and apply 2D convolution to reduce the feature dimension, generating dense BEV features  $\{B_f^{com,0}, B_f^{com,1}, B_f^{com,2}, B_f^{com,3}\}$ . With semantic and geometric BEV features sharing the same dimensions, our BEV feature fusion network adopts a U-Net architecture with 2D convolutions. The encoder comprises an input layer and four residual blocks. To fully utilize geometric and semantic features at different scales, we designed an SCB-Fusion module to fuse current semantic features, geometric features, and BEV features from the previous layer. The fused features can be expressed as:

$$\begin{aligned} F_{SCB} = \Phi & \{ \lambda [N(F_{bev})] \times F_{bev} \\ & + \lambda [N(F_{com})] \times F_{com} \\ & + \lambda [N(F_{sem})] \times F_{sem} \end{aligned} \quad (3)$$

where  $\lambda$  denotes the sigmoid function.  $\Phi$  is the  $1 \times 1$  convolution.

**Total Loss Function.** We train the whole network end-to-end. The multi-task loss  $L_{total}$  is expressed as :

$$L_{total} = 3 \times L_{bev} + L_{sem} + L_{com} \quad (4)$$

#### IV. AERIAL-GROUND MOTION PLANNING

In this section, we introduce the novel AG-Planner. It is built on EGO-Planner [15] and consists of 1) an energy-efficient Kinodynamic A\* path searching front-end, 2) a gradient-based trajectory optimization back-end and 3) a post-refinement procedure. Our AG-Planner evaluates and projects gradient information directly from obstacles instead of a pre-built ESDF like [1].

##### A. Energy-Efficient Kinodynamic Hybrid Path Searching

Our AG-Planner first creates a naive “initial trajectory” (in Fig.4) that overlooks obstacles by randomly adding coordinate points, considering the positions of both the starting and target points. Following that, for the “collision trajectory segment” (i.e., the trajectory inside the obstacle), the back end of our planner employs an energy-efficient kinodynamic A\* path search (in Alg.1) to establish a safe “guidance trajectory segment”  $\tau$ , which uses motion primitives instead of straight lines as graph edges in the searching loop. In this study, we add extra flying and ground-steering energy consumption to the

---

##### Algorithm 1: Energy-Efficient Kinodynamic A\* Search

---

```

Input: Start State  $x_s$  and Target State  $x_g$ 
Output: Energy-Efficient Valid Path between  $x_s$  and  $x_g$ 
Data:  $O = \emptyset$  and  $C = \emptyset$ ;  $f(x_s) = g(x_s) + h(x_s)$ ;  $O.push(x_s)$ 

1 while  $O.empty()$  do
2    $x \leftarrow O.popMin()$ 
3   if  $x == x_g$  then
4     return path
5   end
6   else
7      $C.push(x)$ 
8     foreach  $n \in neig(x)$  do
9        $g_n \leftarrow (um.squaredNorm() + w_{time}) * \tau + g(x)$ 
10      // next node flying
11      if  $z \geq ground\_judge$  then
12         $g_n -= x.fly\_penalty\_g$ 
13        // add fly penalty cost
14         $g_n += fly\_cost * z + f\_cost\_base$ 
15         $fly\_penalty\_g = fly\_cost * z + f\_cost\_base$ 
16         $steer\_penalty\_g = 0$ 
17         $next\_motion\_state = true$ 
18      end
19      // next node driving
20      else
21         $g_n -= x.steer\_penalty\_g$ 
22        // add steer penalty cost
23         $steer\_cost = steer\_cost * pow(\omega_z, 2)$ 
24         $g_n += steer\_cost + ground\_cost\_base$ 
25         $steer\_penalty\_g = steer\_cost + g\_cost\_base$ 
26         $fly\_penalty\_g = 0$ 
27         $next\_motion\_state = false$ 
28      end
29       $f_n = g_n + \lambda * estimateHeuristic(n, x_g)$ 
30      if  $n \notin O \cup C$  then
31         $n.updateCost(g_n, fly\_penalty, steer\_penalty, f_n)$ 
32         $O.push(n)$ 
33      end
34    end
35  end
36 return null // Cannot find a valid path

```

---

motion primitives (in Fig. 4). Consequently, the path searching not only tends to plan ground trajectories and avoid large turns but also switches to aerial mode and flies over them only when AGRs encounter huge obstacles, thereby promoting energy-saving.

##### B. Gradient-Based B-spline Trajectory Optimization

**B-spline Trajectory Formulation:** In trajectory optimization, the trajectory is parameterized by a uniform B-spline curve  $\Theta$ , which is uniquely determined by its degree  $p_b$ ,  $N_c$  control points  $\{Q_1, Q_2, Q_3, \dots, Q_{N_c}\}$ , and a knot vector  $\{t_1, t_2, t_3, \dots, t_{M-1}, t_M\}$ , where  $Q_i \in \mathbb{R}^3$ ,  $t_m \in \mathbb{R}$ ,  $M = N + p_b$ . Following the matrix representation of the [29] the value of a B-spline can be evaluated as:

$$\Theta(u) = [1, u, \dots, u^p] \cdot M_{p_b+1} \cdot [Q_{i-p_b}, Q_{i-p_b+1}, \dots, Q_i]^T \quad (5)$$

where  $M_{p+1}$  is a constant matrix depends only on  $p_b$ . And  $u = (t - t_i)/(t_{i+1} - t_i)$ , for  $t \in [t_i, t_{i+1}]$ . In particular, in ground mode, we assume that AGR is driving on flat

ground so that the vertical motion can be omitted and we only need to consider the control points in the two-dimensional horizontal plane, denoted as  $Q_{\text{ground}} = \{Q_{t0}, Q_{t1}, \dots, Q_{tM}\}$ , where  $Q_{ti} = (x_{ti}, y_{ti})$ ,  $i \in [0, M]$ . In aerial mode, the control points are denoted as  $Q_{\text{aerial}}$ . According to the properties of B-spline: the  $k^{\text{th}}$  derivative of a B-spline is still a B-spline with order  $p_{b,k} = p_b - k$ , since  $\Delta t$  is identical along  $\Theta$ , the control points of the velocity  $V_i$ , acceleration  $A_i$  and jerk  $J_i$  curves are obtained by:

$$V_i = \frac{Q_{i+1} - Q_i}{\Delta t}, A_i = \frac{V_{i+1} - V_i}{\Delta t}, J_i = \frac{A_{i+1} - A_i}{\Delta t} \quad (6)$$

**Collision Avoidance Force Estimation:** Inspired by [15], For each control point on the collision trajectory segment, vector  $v$  (i.e., a safe direction pointing from inside to outside of that obstacle) is generated from  $\iota$  to  $\tau$  and  $p$  is defined at the obstacle surface. With generated  $\{p, v\}$  pairs, the planner maximizes  $D_{ij}$  and returns an optimized trajectory. The obstacle distance  $D_{ij}$  if  $i^{\text{th}}$  control point  $Q_i$  to  $j^{\text{th}}$  obstacle is defined as:

$$D_{ij} = (Q_i - p_{ij}) \times v_{ij} \quad (7)$$

Because the guide path  $\iota$  is energy-saving, the generated path is also energy efficient.

**B-spline Trajectory Optimization:** The basic requirements of the re-planned B-spline are three-folds: smoothness, safety, and dynamical feasibility. Based on the special properties of AGR bimodal, we firstly adopt the following cost terms designed by Zhou et al. [15]:

$$\min J_1 = \lambda_s J_s + \lambda_c J_c + \lambda_f (J_v + J_a + J_j) \quad (8)$$

where  $J_s$  is the smoothness penalty,  $J_c$  is for collision, and  $J_v, J_a, J_j$  are dynamical feasibility costs that limit velocity, acceleration and jerk.  $\lambda_s, \lambda_c, \lambda_f$  are weights for each cost terms. Detailed explanations can be found in [15]. Subsequently, based on our observations, AGR faces non-holonomic constraints when driving on the ground, which means that the ground velocity vector of AGR must be aligned with its yaw angle. Additionally, AGR needs to deal with curvature limitations that arise due to minimizing tracking errors during sharp turns. Therefore, a cost for curvature needs to be added, and  $J_n$  can be formulated as:

$$J_n = \sum_{i=1}^{M-1} F_n(Q_{ti}) \quad (9)$$

where  $F_n(Q_{ti})$  is a differentiable cost function with  $C_{\max}$  specifying the curvature threshold:

$$F_n(Q_{ti}) = \begin{cases} (C_i - C_{\max})^2, & C_i > C_{\max}, \\ 0, & C_i \leq C_{\max} \end{cases} \quad (10)$$

where  $C_i = \frac{\Delta \beta_i}{\Delta Q_{ti}}$  is the curvature at  $Q_{ti}$ , and the  $\Delta \beta_i = \left| \tan^{-1} \frac{\Delta y_{ti+1}}{\Delta x_{ti+1}} - \tan^{-1} \frac{\Delta y_{ti}}{\Delta x_{ti}} \right|$ . In general, the overall objec-

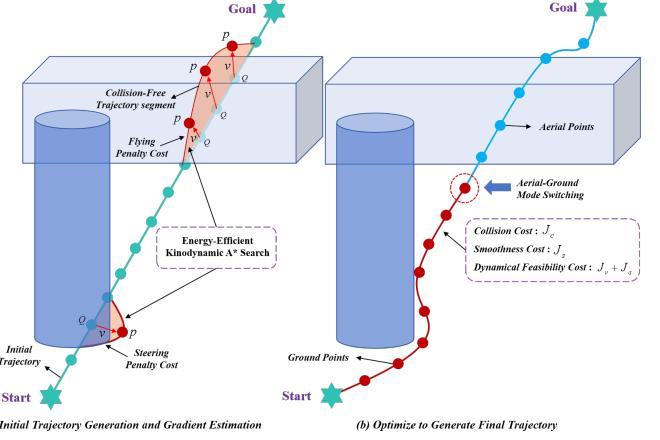


Fig. 4: Illustration of AG-Planner and topological trajectory generation.

tive function is formulated as follows:

$$\begin{aligned} \min J_{\text{all}} = & \lambda_s J_s + \lambda_c J_c + \lambda_f (J_v + J_a + J_j) + \lambda_n J_n \\ \text{s.t. } & \left\{ \begin{array}{l} J_s = \sum_{i=1}^{N_c-1} \|A_i\|_2^2 + \sum_{i=1}^{N_c-2} \|J_i\|_2^2 \\ J_c = \sum_{i=1}^{N_c} j_c(Q_i) \\ J_v = \sum_{i=1}^{N_c} \omega_v F(V_i) \\ J_a = \sum_{i=1}^{N_c-1} \omega_a F(A_i) \\ J_j = \sum_{i=1}^{N_c-2} \omega_j F(J_i) \\ J_n = \sum_{i=1}^{M-1} F_n(Q_{ti}) \end{array} \right. \end{aligned} \quad (11)$$

The optimization problem is solved by a non-linear optimization solver NLOpt. After path planning is completed, a setpoint on the generated trajectory is selected according to the current timestamp and then sent to the controller. An aerial setpoint includes the yaw angle and 3D position, velocity, and acceleration. A terrestrial one includes the yaw angle and 2D position and velocity. In addition, when the Z-axis coordinate of the next control point is greater than the ground threshold, that is, when mode switching is required, an additional trigger signal will be sent to the controller (i.e., PX4 Autopilot). The controller will automatically switch to *Offboard Mode* to enter the flight state.

## V. EVALUATION

In this section, we first evaluate the perception module (i.e., LBSCNet) on the SemanticKITTI benchmark for its accuracy in SSC tasks, as well as its real-time inference and update capabilities. We then integrate the perception module and the planning module by deploying a pre-trained model offline, forming a complete HE-Nav system. Subsequently, we conduct experiments in both simulated and real-world environments to assess the performance of the aerial-ground robot (AGR) when using HE-Nav for autonomous navigation, focusing on **performance** metrics (i.e. planning success rate, total movement time) and **efficiency** (planning time and energy consumption)..

### A. Evaluation setup

**Perception Module:** For the training and testing of LBSCNet, we carried out the process on a server equipped with 4 NVIDIA RTX 3090 GPUs, 128GB of memory, and an Intel(R) Xeon(R) Silver 4216 CPU @ 2.10GHz. The large-scale outdoor scenarios dataset SemanticKITTI [30] is the dataset used. We trained the LBSCNet model for 80 epochs on a single NVIDIA 3090 with a batch size of 12, employing the Adam optimizer [31] with an initial learning rate of 0.001, and augmenting the input point cloud by randomly flipping along the x-y axis during the training process. Ultimately, we deployed the pre-trained model offline with the best completion accuracy to predict occlusion areas.

**Navigation Simulation Experiment:** Simulation experiments were conducted on a laptop with Ubuntu 20.04, i9-13900HX CPU, and NVIDIA RTX 4060 GPU. We simulated aerial-ground robot navigation in complex scenarios, consisting of a  $20m \times 20m$  room and a  $3m \times 30m$  corridor with numerous random obstacles, creating occluded spaces and unknown areas. The AGR's task was to navigate from a starting point to a designated destination without collision. We also record the total movement and path planning time to compare with the baseline.

**Indoor and Outdoor Real-world Experiment:** We deployed HE-Nav on a custom AGR platform (in Fig. 5) for real-world indoor and outdoor environment experiments. This platform utilizes the Prometheus software [32] and is equipped with a RealSense D435i depth camera and a T265 camera. Additionally, it features a Jetson Xavier NX onboard computer to run the HE-Nav. More detailed hardware specifications are provided in the supplementary materials.

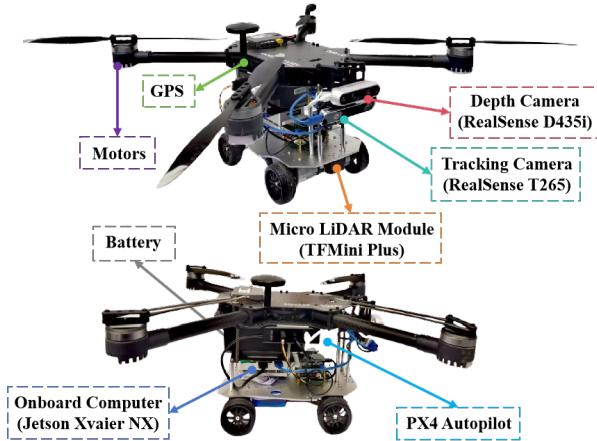


Fig. 5: The detailed composition of the robot platform.

**Metrics:** For the perception module, we use intersection over union (IoU) to evaluate scene completion quality and the mean IoU (mIoU) of 19 semantic classes to assess the performance of semantic segmentation. Moreover, we also focus on LBSCNet's inference speed to ensure it meets the real-time requirements for autonomous navigation. Regarding navigation, we pay attention to performance metrics such as

planning success rate (%) and total moving time (t), planning time (ms) as well as energy consumption (W).

**Baseline methods.** For the perception module, we compare LBSCNet against the state-of-the-art SSC methods with public resources: (1) a camera-based SSC method MonoScene [12] and VoxFormer [13], (2) point-cloud-based SSC methods including LMSCNet [14], and SSCNet [33] and SCPNet [11]. To evaluate the performance and efficiency of HE-Nav, we compared HE-Nav with two state-of-the-art AGRs navigation systems: TABV [1], Fan et al. [2].

Parameter	Value
Battery Capacity	10000 mAh
Battery Weight	1008 g
Rated Power	231 Wh
Operating Voltage	23.05 V
Driving Energy Consumption	$\approx 252.00$ W
Hovering Energy Consumption	$\approx 533.08$ W
Flying Energy Consumption	$\approx 990.00$ W

TABLE I: Battery and Energy Consumption Parameters

Method	IoU	mIoU	Prec.	Recall	FPS
SSCNet [33]	53.20	14.55	59.13	<b>84.15</b>	12.00
LMSCNet [14]	55.32	17.01	77.11	66.19	13.50
LMSCNet-SS [14]	56.72	17.62	<b>81.55</b>	65.07	13.50
S3CNet [22]	45.60	29.50	48.79	77.13	1.20
Monoscene [12]	38.55	12.22	51.96	59.91	< 1
VoxFromer-T [13]	57.69	18.42	69.95	76.70	< 1
VoxFromer-S [13]	57.54	16.48	70.85	75.39	< 1
SCPNet [11]	56.10	<b>36.70</b>	72.43	78.61	< 1
<b>LBSCNet (Ours)</b>	<b>59.71</b>	23.58	77.60	71.29	<b>20.08</b>

TABLE II: Quantitative comparison against the state-of-the-art SSC methods on the official SemanticKITTI benchmark.

### B. LBSCNet Comparison against the state-of-the-art.

**Quantitative Results:** We evaluated our proposed LBSCNet against state-of-the-art SSC methods on the SemanticKITTI test datasets by submitting results to the official test server. Table I demonstrates that LBSCNet not only achieves the highest completion metric IoU (59.71%) but also ranks third in the semantic segmentation metric mIoU (23.58%). Although SCPNet's semantic segmentation accuracy surpasses ours, its dense network design renders it incapable of real-time operation (i.e., FPS < 1). In contrast, LBSCNet outperforms SCPNet by 6.45% in IoU and runs approximately 20 times faster.

The exceptional accuracy and real-time inference performance of our LBSCNet can be attributed to the innovative semantic and completion decoupling network structure, which

leverages contextual semantic information to enhance scene understanding and completion. Additionally, the incorporation of the novel SCB-fusion module and CCA module allows the network to remain lightweight while significantly improving completion accuracy by capturing contextual features and learning long-distance dependencies. Furthermore, our LBSCNet exhibits low latency and real-time operation (20.08 FPS) due to the utilization of sparse 3D convolutions and lightweight BEV feature fusion within the network. Consequently, LBSCNet is well-suited for real-time perception in ARG navigation systems.

**Qualitative Results:** We provide visualizations on the SemanticKITTI validation set, as depicted in Fig. 5, and include results from LMSCNet [14], Monoscene[12], VoxFormer[13], and SCPNet [11] for a comprehensive comparison. As illustrated in Fig. 6, our LBSCNet demonstrates superior SSC predictions, particularly for “wall” classes and larger objects like cars, aligning with the results in Table I. Importantly, the occlusion areas we target, such as vegetation and trees behind walls, are accurately completed, proving vital for subsequent path planning applications.

More qualitative and quantitative results are provided in the supplementary material, i.e., in Section VII-A.

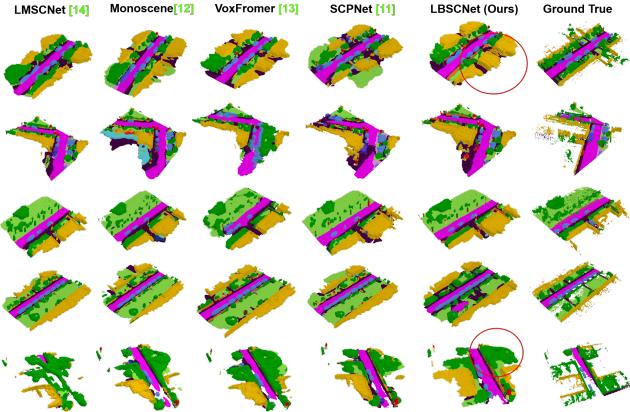


Fig. 6: Qualitative results of our method and others. LBSCNet better captures the scene layout in large-scale scenarios.

**Ablation Study:** Ablation studies on the SemanticKITTI validation set (Table 4) highlight the significance of two key components in our network: self-attention mechanisms and SCB-Fusion Module. The CCA mechanism substantially impacts completion and semantic prediction by effectively aggregating context across rows and columns. *Without CCA* causes a 3.86% and 7.48% drop for completion and semantic completion, respectively. Meanwhile, SCB-Fusion captures local scene features, such as occluded areas, with low computational overhead. *Without SCB-Fusion* leads to a 2.47% decline in IoU.

### C. Simulated Air-Ground Robot Navigation

We conducted a comparative analysis of our HE-Nav navigation system against TABV [1] and HDF [2] in a square room and corridor scenario. 100 trials with varying obstacle

Method	IoU $\uparrow$	mIoU $\uparrow$
LBSCNet (ours)	54.92	17.69
w/o SCB-Fusion Module	54.15	17.26
w/o Criss-Cross Attention	52.80	16.37

TABLE III: Ablation study of our model design choices on the SemanticKITTI validation set.

placements, we recorded the moving time, length, energy consumption, planning time and success rate (i.e., no collisions). As depicted in Fig. 7, our HE-Nav system boasts

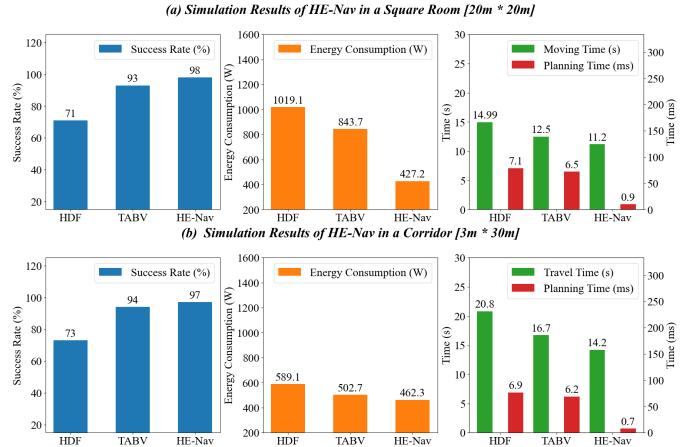


Fig. 7: Quantitative results of HE-Nav in two simulation scenarios.

remarkable planning success rates of 98% and 97%, with total movement times of 11.2s and 14.2s in square rooms and corridors, respectively. This exceptional performance stems from our cutting-edge LBSCNet’s capacity to predict obstacle distribution in occluded areas, enabling planners to effectively circumvent these zones and substantially minimize collision risks. Furthermore, our path planner integrates seamlessly with the KinoDynamic A\* algorithm, resulting in the lowest overall energy consumption (i.e., 427.2W and 462.3W). In square scenarios, HE-Nav achieves a 4x reduction in energy consumption compared to state-of-the-art TABV planning methods. Regarding planning time, our HE-Nav system experiences a 6x acceleration compared to TABV, due to the elimination of redundant ESDF calculations. As depicted in Fig. 8, the

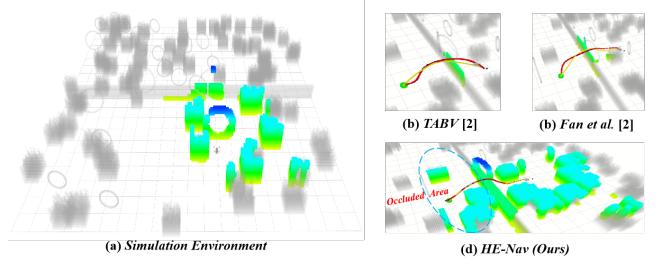
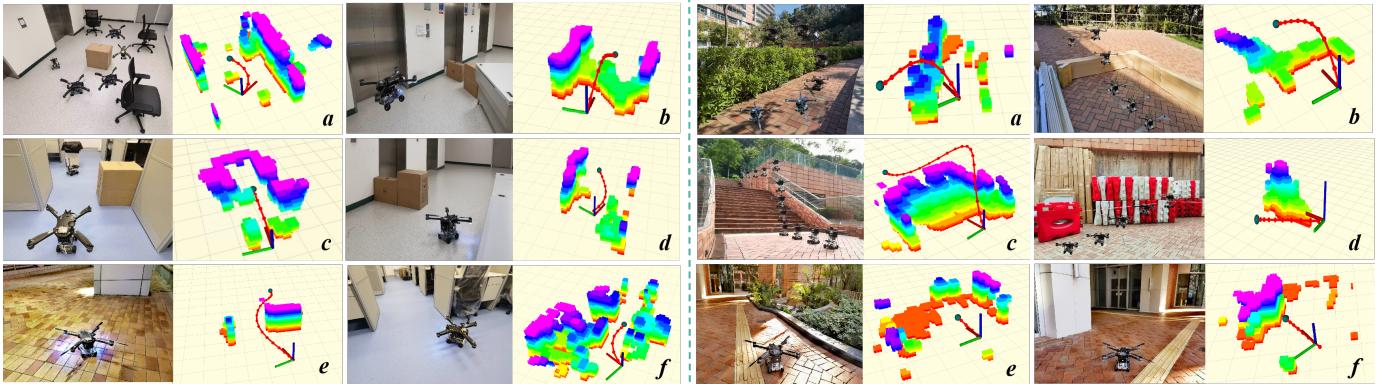


Fig. 8: Qualitative results of path planning and occlusion prediction in simulation environment.

path generated by the HDF fails to effectively consider both



(a) Indoor Real-World Experiments.

(b) Outdoor real-world Experiments.

Fig. 9: Four methods were used to plan paths in a simulated square room. AGRNav demonstrates the ability to predict the distribution of obstacles in occluded areas.

smoothness and dynamic feasibility. Additionally, the TABV path primarily focuses on the energy consumption associated with flight, which results in premature flight actions and consequently leads to increased energy consumption, rendering the overall energy consumption suboptimal. This lack of perception causes TABV to encounter difficulties in pathfinding, further exacerbating energy consumption. In contrast, our HE-Nav system adeptly addresses this shortcoming through its ability to perceive and predict occlusions, thereby optimizing both path planning and energy consumption. For a more comprehensive understanding of the qualitative results, please refer to the supplementary material provided.

#### D. Real-world Air-Ground Robot Navigation

The real-world experiment was conducted using a customized AGR, with average energy consumption during driving, flying, and hovering at 197.52W, 532.07W, and 987.61W, respectively. We deployed HE-Nav on the NVIDIA Jetson Xavier NX onboard computer, where LBSCNet utilized Deploy after TensorRT optimization. Initially, the depth camera captures sparse point cloud data, serving as LBSCNet's input. The dense prediction results are then updated to the local map in real-time, allowing the AG-Planner to plan the path.

We assessed the performance and energy efficiency of HE-Nav across six indoor and six outdoor scenes. Figure X displays the real environment and path planning results. As per the quantitative results in Figure XX, HE-Nav exhibits the lowest overall energy consumption indoors, primarily due to the inclusion of an additional turning penalty term for the ground segment. This term minimizes large-angle turning paths, effectively constraining energy usage in ground mode. In outdoor scenes, energy consumption is reduced by XXX compared to HDF. The elimination of ESDF also significantly simplifies path planning time. On the Jetson Xavier NX, planning time is reduced to 1.2ms, which is five times faster than ESDF-based TABV.

We provide additional qualitative and quantitative results in the supplementary material.

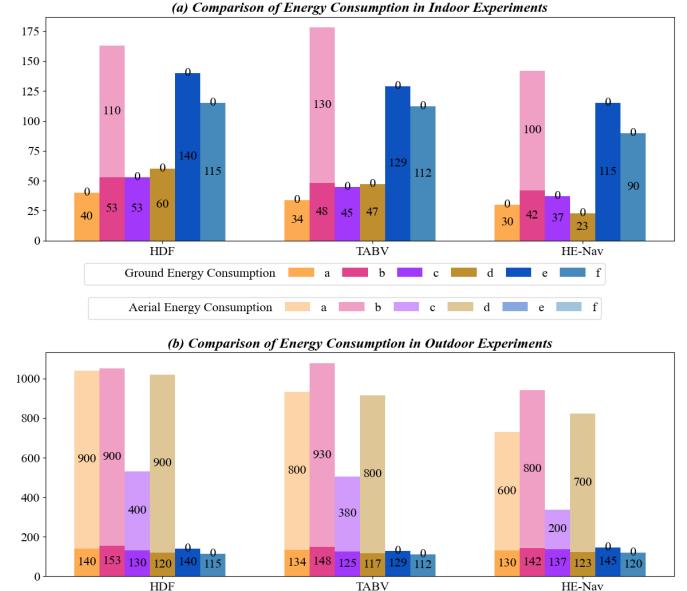


Fig. 10: Qualitative results of path planning and occlusion prediction in simulation environment.

## VI. CONCLUSION

we have presented HE-Nav, a high-performance and efficient navigation system specifically designed for aerial-ground robots (AGR). By integrating innovative features such as the lightweight BEV-guided semantic scene completion network (LBSCNet) and the aerial-ground motion planner (AG-planner), our system is capable of predicting obstacle distributions in occluded areas and generating low-collision risk, energy-efficient aerial-ground hybrid trajectories in real-time. Through extensive simulations and real experiments, HE-Nav has been shown to significantly outperform recent planning frameworks in performance (i.e., planning success rate and total movement time) and efficiency (i.e., planning time and energy consumption).

## VII. ACKNOWLEDGMENTS

We thank all anonymous reviewers for their helpful comments.

## REFERENCES

- [1] Ruibin Zhang, Yuze Wu, Lixian Zhang, Chao Xu, and Fei Gao. Autonomous and adaptive navigation for terrestrial-aerial bimodal vehicles. *IEEE Robotics and Automation Letters*, 7(2):3008–3015, 2022.
- [2] David D Fan, Rohan Thakker, Tara Bartlett, Meriem Ben Miled, Leon Kim, Evangelos Theodorou, and Ali-akbar Agha-mohammadi. Autonomous hybrid ground/aerial mobility in unknown environments. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3070–3077. IEEE, 2019.
- [3] Neng Pan, Jinqi Jiang, Ruibin Zhang, Chao Xu, and Fei Gao. Skywalker: A compact and agile air-ground omnidirectional vehicle. *IEEE Robotics and Automation Letters*, 8(5):2534–2541, 2023.
- [4] Ruibin Zhang, Junxiao Lin, Yuze Wu, Yuman Gao, Chi Wang, Chao Xu, Yanjun Cao, and Fei Gao. Model-based planning and control for terrestrial-aerial bimodal vehicles with passive wheels. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1070–1077, 2023. doi: 10.1109/IROS55552.2023.10342188.
- [5] Xinyu Zhang, Yuanhao Huang, Kangyao Huang, Xiaoyu Wang, Dafeng Jin, Huaping Liu, and Jun Li. A multi-modal deformable land-air robot for complex environments. *arXiv preprint arXiv:2210.16875*, 2022.
- [6] Xinyu Zhang, Yuanhao Huang, Kangyao Huang, Ziqi Zhao, Jingwei Li, Huaping Liu, and Jun Li. Coupled modeling and fusion control for a multi-modal deformable land-air robot. *arXiv preprint arXiv:2211.04185*, 2022.
- [7] Eric Sihite, Arash Kalantari, Reza Nemovi, Alireza Ramezani, and Morteza Gharib. Multi-modal mobility morphobot (m4) with appendage repurposing for locomotion plasticity enhancement. *Nature communications*, 14(1):3323, 2023.
- [8] Youming Qin, Yihang Li, Xu Wei, and Fu Zhang. Hybrid aerial-ground locomotion with a single passive wheel. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1371–1376. IEEE, 2020.
- [9] Qifan Tan, Xinyu Zhang, Huaping Liu, Shuyuan Jiao, Mo Zhou, and Jun Li. Multimodal dynamics analysis and control for amphibious fly-drive vehicle. *IEEE/ASME Transactions on Mechatronics*, 26(2):621–632, 2021.
- [10] Xiaoyu Wang, Kangyao Huang, Xinyu Zhang, Honglin Sun, Wenzhuo Liu, Huaping Liu, Jun Li, and Pingping Lu. Path planning for air-ground robot considering modal switching point optimization. In *2023 International Conference on Unmanned Aircraft Systems (ICUAS)*, pages 87–94. IEEE, 2023.
- [11] Zhaoyang Xia, Youquan Liu, Xin Li, Xinge Zhu, Yuexin Ma, Yikang Li, Yuenan Hou, and Yu Qiao. Scpnet: Semantic scene completion on point cloud. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17642–17651, 2023.
- [12] Anh-Quan Cao and Raoul de Charette. Monoscene: Monocular 3d semantic scene completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3991–4001, 2022.
- [13] Yiming Li, Zhiding Yu, Christopher Choy, Chaowei Xiao, Jose M Alvarez, Sanja Fidler, Chen Feng, and Anima Anandkumar. Voxformer: Sparse voxel transformer for camera-based 3d semantic scene completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9087–9098, 2023.
- [14] Luis Roldao, Raoul de Charette, and Anne Verroust-Blondet. Lmscnet: Lightweight multiscale 3d semantic completion. In *2020 International Conference on 3D Vision (3DV)*, pages 111–119. IEEE, 2020.
- [15] Xin Zhou, Zhepei Wang, Hongkai Ye, Chao Xu, and Fei Gao. Ego-planner: An esdf-free gradient-based local planner for quadrotors. *IEEE Robotics and Automation Letters*, 6(2):478–485, 2020.
- [16] Xu Yan, Jiantao Gao, Jie Li, Ruimao Zhang, Zhen Li, Rui Huang, and Shuguang Cui. Sparse single sweep lidar point cloud segmentation via learning contextual shape priors from scene completion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 3101–3109, 2021.
- [17] Zhijian Liu, Haotian Tang, Alexander Amini, Xinyu Yang, Huizi Mao, Daniela L Rus, and Song Han. Bevfusion: Multi-task multi-sensor fusion with unified bird’s-eye view representation. In *2023 IEEE international conference on robotics and automation (ICRA)*, pages 2774–2781. IEEE, 2023.
- [18] Tong Wu, Yimin Zhu, Lixian Zhang, Jianan Yang, and Yihang Ding. Unified terrestrial/aerial motion planning for hytaqs via nmpc. *IEEE Robotics and Automation Letters*, 8(2):1085–1092, 2023.
- [19] Arash Kalantari and Matthew Spenko. Design and experimental validation of hytaq, a hybrid terrestrial and aerial quadrotor. In *2013 IEEE International Conference on Robotics and Automation*, pages 4445–4450. IEEE, 2013.
- [20] Mikhail Martynov, Zhanibek Darush, Aleksey Fedoseev, and Dzmitry Tsetserukou. Morphogear: An uav with multi-limb morphogenetic gear for rough-terrain locomotion. In *2023 IEEE/ASME International Conference on Advanced Intelligent Mechatronics (AIM)*, pages 11–16. IEEE, 2023.
- [21] Muqing Cao, Xinhang Xu, Shanghai Yuan, Kun Cao, Kangcheng Liu, and Lihua Xie. Doublebee: A hybrid aerial-ground robot with two active wheels. *arXiv preprint arXiv:2303.05075*, 2023.
- [22] Ran Cheng, Christopher Agia, Yuan Ren, Xinhai Li, and Liu Bingbing. S3cnet: A sparse semantic scene

- completion network for lidar point clouds. In *Conference on Robot Learning*, pages 2148–2161. PMLR, 2021.
- [23] Sicheng Zuo, Wenzhao Zheng, Yuanhui Huang, Jie Zhou, and Jiwen Lu. Pointocc: Cylindrical tri-perspective view for point-based 3d semantic occupancy prediction. *arXiv preprint arXiv:2308.16896*, 2023.
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [25] Shuangjie Xu, Rui Wan, Maosheng Ye, Xiaoyi Zou, and Tongyi Cao. Sparse cross-scale attention network for efficient lidar panoptic segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 2920–2928, 2022.
- [26] Maxim Berman, Amal Rannen Triki, and Matthew B Blaschko. The lovász-softmax loss: A tractable surrogate for the optimization of the intersection-over-union measure in neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4413–4421, 2018.
- [27] Zhilu Zhang and Mert Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. *Advances in neural information processing systems*, 31, 2018.
- [28] Spconv Contributors. Spconv: Spatially sparse convolution library. <https://github.com/traveller59/spconv>, 2022.
- [29] C de Boor. Subroutine package for calculating with b-splines, 1971.
- [30] 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.
- [31] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [32] Amovlab. Prometheus UAV open source project. <https://github.com/amov-lab/Prometheus>.
- [33] Shuran Song, Fisher Yu, Andy Zeng, Angel X Chang, Manolis Savva, and Thomas Funkhouser. Semantic scene completion from a single depth image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1746–1754, 2017.

# SUPPLEMENTARY MATERIAL FOR OUR “HE-NAV”

In the supplementary material, we discuss additional implementation details and provide more qualitative and quantitative results about *LBSCNet*, *Simulation Experiments* and *Real-World Experiments*. We encourage the reader to browse these results and videos.

## A. LBSCNet

In Table IV, we present an extensive array of quantitative results, encompassing completion accuracy and semantic segmentation accuracy. Moreover, the visualization of outcomes in the SemantiKITTI dataset validation set is depicted in Fig. 11. It is evident that LBSCNet excels in comparison to other methods with respect to completion and semantic representation of roads, vehicles, buildings, and vegetation, which is in alignment with the findings displayed in Table 3. Despite our semantic segmentation results ranking third among all approaches, we possess superior completion accuracy and real-time performance. This is of paramount importance for Autonomous Ground Robots (AGRs) to accurately and promptly predict the distribution of obstacles in occluded areas during navigation.

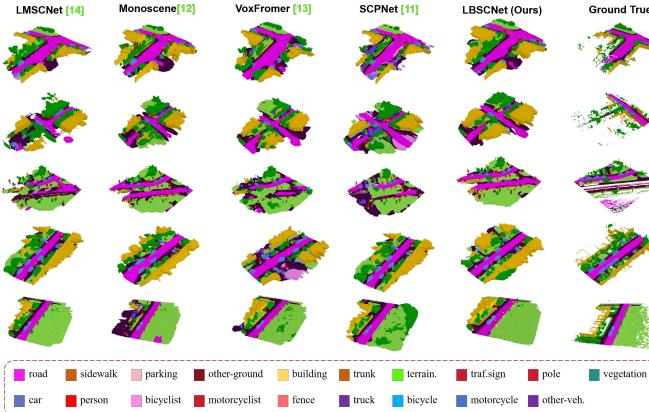


Fig. 11: Qualitative results of our method and others. LBSCNet better captures the scene layout in large-scale scenarios.

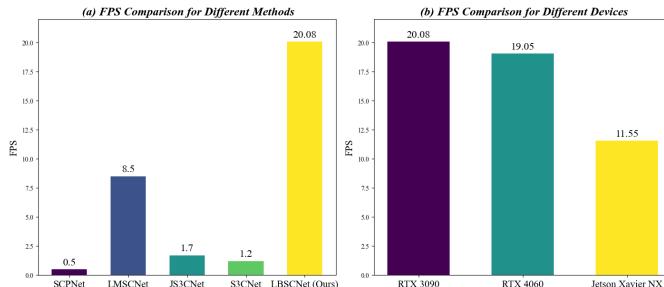


Fig. 12: Qualitative results of our method and others. LBSCNet better captures the scene layout in large-scale scenarios.

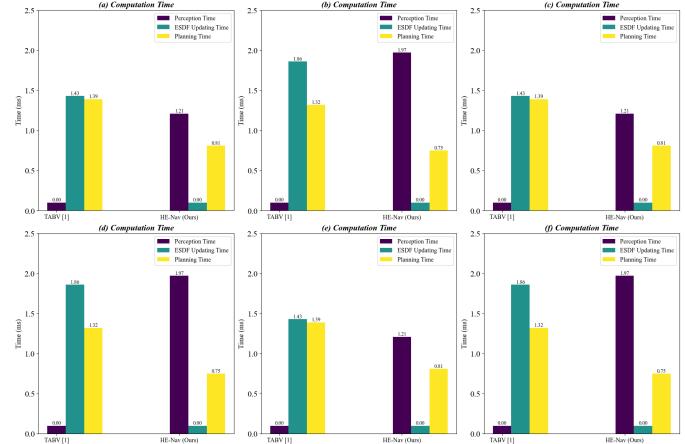


Fig. 13: Qualitative results of path planning and occlusion prediction in simulation environment.

In addition, as illustrated in Fig. 12a, the inference speed comparison of LBSCNet highlights its performance advantages. Owing to their dense 3D convolution design, existing point cloud-based SSC methods are unable to achieve real-time inference. Concurrently, Fig. 12b demonstrates the inference speed of LBSCNet on various devices. It achieves 20.08 FPS on an RTX 3090 GPU and 19.05 FPS on an RTX 4060 GPU (i.e., in a simulated experiment). Furthermore, when optimized by TensorRT on a Jetson Xavier NX, LBSCNet attains a real-time performance of 11.55 FPS (i.e., in a real-world experiment).

<b>Method</b>	<b>LBSCNet (Ours)</b>	<b>SCPNet [11]</b>	<b>VoxFormer [13]</b>	<b>MonoScene [12]</b>	<b>LMSCNet [14]</b>
<b>IoU (%)</b>	<b>59.71</b>	56.10	57.69	38.55	54.89
<b>Precision (%)</b>	78.60	68.13	69.95	51.96	<b>82.21</b>
<b>Recall (%)</b>	71.29	74.92	<b>76.70</b>	59.91	62.29
<b>mIoU</b>	23.58	<b>36.70</b>	18.42	12.22	14.13
<b>car</b> (3.92%)	35.80	<b>46.40</b>	37.46	24.64	35.41
<b>bicycle</b> (0.03%)	8.00	<b>33.20</b>	2.87	0.23	0.00
<b>motorcycle</b> (0.03%)	4.10	<b>34.90</b>	1.24	0.20	0.00
<b>truck</b> (0.16%)	4.90	13.80	<b>10.38</b>	13.84	3.49
<b>other-veh.</b> (0.20%)	8.10	<b>29.10</b>	10.61	2.13	0.00
<b>person</b> (0.07%)	3.40	<b>28.20</b>	3.50	1.37	0.00
<b>bicyclist</b> (0.07%)	2.70	<b>24.70</b>	3.92	1.00	0.00
<b>motorcyclist</b> (0.05%)	1.80	1.80	0.00	0.00	0.00
<b>road</b> (15.30%)	<b>71.30</b>	68.50	66.15	57.11	67.56
<b>parking</b> (1.12%)	39.40	<b>51.30</b>	23.96	18.60	13.22
<b>sidewalk</b> (11.13%)	42.90	<b>49.80</b>	34.53	27.58	34.20
<b>other-grnd</b> (0.56%)	16.70	<b>30.70</b>	0.76	2.00	0.00
<b>building</b> (14.10%)	<b>43.40</b>	38.80	29.45	15.97	27.83
<b>fence</b> (3.90%)	31.50	<b>44.70</b>	11.15	7.37	4.42
<b>vegetation</b> (39.3%)	45.10	<b>46.40</b>	38.07	19.68	33.32
<b>trunk</b> (0.51%)	26.20	<b>40.10</b>	12.75	2.57	3.01
<b>terrain</b> (9.17%)	40.90	<b>48.70</b>	39.61	31.59	41.51
<b>pole</b> (0.29%)	15.00	<b>40.40</b>	15.56	3.79	4.43
<b>traf.-sign</b> (0.08%)	6.80	<b>25.10</b>	8.09	2.54	0.00

TABLE IV: Quantitative comparison against the state-of-the-art SSC methods.