

Hestia: Hierarchical Next-Best-View Exploration for Systematic Intelligent Autonomous Data Collection

CHENG-YOU LU, University of Technology Sydney

ZHUOLI ZHUANG, University of Technology Sydney

TRUNG LE, University of Technology Sydney

DA XIAO, University of Technology Sydney

YU-CHENG FRED CHANG, University of Technology Sydney

THOMAS DO, University of Technology Sydney

SRINATH SRIDHAR, Brown University

CHIN-TENG LIN, University of Technology Sydney

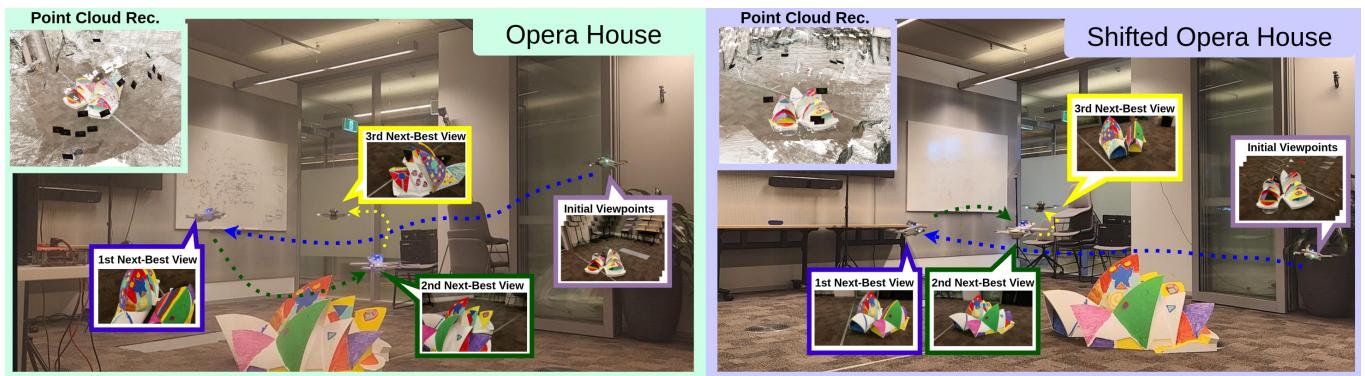


Fig. 1. Hestia is an active autonomous data collection system designed to facilitate flexible and efficient data collection. Hestia combines a next-best-view planner for viewpoint prediction with a real-world data collection system utilizing a drone equipped with an RGB camera. By integrating a pointmap estimation model [11, 68] and initializing with three viewpoints, Hestia ensures synchronization [64] between real-world and virtual environments. Hestia demonstrates robust performance, accurately predicting viewpoints in real-world scenarios even when objects are shifted. We hope that this prototype data collection system will advance the field of large-scale scene data collection using drones. *Please refer to our demonstration video in the supplementary materials for further details.*

Advances in 3D reconstruction and novel view synthesis methods have enabled efficient and photorealistic rendering. However, the data collection process for the methods is typically manual, making it time-consuming and labor-intensive. To address the challenges, this study introduces Hierarchical Next-Best-View Exploration for Systematic Intelligent Autonomous Data Collection (Hestia)¹. Hestia leverages reinforcement learning (RL) to develop a generalizable policy capable of predicting the next-best view in a 5-Dof

¹In Greek mythology, Hestia is the goddess of the hearth, symbolizing home, foundation, and structure, representing stability and guidance in complex systems

Authors' Contact Information: Cheng-You Lu, University of Technology Sydney; Zhuoli Zhuang, University of Technology Sydney; Trung Le, University of Technology Sydney; Da Xiao, University of Technology Sydney; Yu-Cheng Fred Chang, University of Technology Sydney; Thomas Do, University of Technology Sydney; Srinath Sridhar, Brown University; Chin-Teng Lin, University of Technology Sydney.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference acronym 'XX, June 03–05, 2018, Woodstock, NY'

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06

<https://doi.org/XXXXXXXX.XXXXXXXX>

space. Unlike prior voxel-based RL approaches, Hestia adopts the concept that a voxel is worth more than a ray, treating each voxel as a cube rather than a point for observation and reward calculation to achieve a higher coverage ratio with less data. For observation, Hestia is trained on the largest dataset for this task, ensuring robust generalization. For action, Hestia employs a hierarchical structure to separately predict the look-at point and camera position, simplifying the high-dimensional next-best-view task. For training, Hestia incorporates a greedy formulation to mitigate spurious correlations, improving training efficiency. The experimental results tested in the NVIDIA IsaacLab environment show that Hestia can achieve robust performance across datasets with objects in different positions. Hestia is integrated into the proposed real-world autonomous data collection system, demonstrating its feasibility. Unlike traditional systems, where sensors are either passive or manually controlled, Hestia's real-world deployment highlights its potential to advance active, autonomous data collection systems for capturing complex scenes. *The code and data processing script for Hestia are included in the supplementary materials and will be released after the paper is published.*

CCS Concepts: • Computing methodologies → Vision for robotics; Markov decision processes; Point-based models; • Computer systems organization → Robotic autonomy.

Additional Key Words and Phrases: Next-best-view planner, Reinforcement learning, Data collection system, Unmanned aerial vehicle

ACM Reference Format:

Cheng-You Lu, Zhuoli Zhuang, Trung Le, Da Xiao, Yu-Cheng Fred Chang, Thomas Do, Srinath Sridhar, and Chin-Teng Lin. 2018. Hestia: Hierarchical Next-Best-View Exploration for Systematic Intelligent Autonomous Data Collection. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX')*. ACM, New York, NY, USA, 12 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

Throughout the years, multiview-based 3D scene reconstruction [1, 11, 21, 25, 45, 47, 48, 59, 60, 66, 68–71, 73] and novel view synthesis [3, 4, 8, 12, 13, 24, 40, 44, 56–58, 75, 79, 80, 82] have been critical research fields in computer vision and graphics. These methods leverage multiview information (e.g., correspondences) from images to reconstruct high-fidelity scenes or synthesize photorealistic novel views. However, collecting multiview data is often time-consuming and labor-intensive, as these methods usually depend on manually captured data. The situation becomes more challenging for generalizable approaches [3, 8, 65, 75, 79, 80, 82] that require diverse and extensive data [20, 34, 65, 76, 81]. From a data collection perspective, one way to reduce the effort is to develop machines (e.g., multiview camera rigs) [2, 29, 32, 38, 53, 72, 74] to minimize human intervention. Although these multiview data collection machines make the data collection process more scalable, their sensors are usually passive (e.g., non-movable), lacking the ability to actively explore the scene. This lack of active exploration may result in missing occluded parts during data capture. To enable active capture and alleviate the effort required for manual capture, we propose Hestia, a Hierarchical Next-Best-View Exploration for Systematic Intelligent Autonomous Data Collection. The goal of Hestia is to actively collect data from object-centric scenes with a higher coverage ratio and less amount of data, demonstrating that the proposed prototype has the potential to operate in real-world settings and to advance the development of active, intelligent data collection systems. We address this problem by introducing two main components, which are a next-best-view planner that actively determines a five-degree-of-freedom (5 DoF) viewpoint (e.g., camera position, yaw, and pitch) based on voxel face observations and a systematic intelligent, autonomous data collection system that uses a drone equipped with an RGB camera as the agent for data collection.

In recent decades, next-best-view planners have been developed for active capture, demonstrating promising potential for reducing the human workload. This is especially important when a professional human pilot is unavailable to control the robot (e.g., a drone) in time-critical tasks, such as search and rescue (SaR), where human resources are limited or urgently needed elsewhere, and in routine tasks (e.g., environmental documentation) that require frequent data collection. Traditional next-best-view planners [17, 18, 35, 42, 77] rely on heuristic rules to predict the viewpoint. Although these approaches show promising and awesome performance for specific scenarios, the same handcrafted rules or hyperparameters may not be the best setting for other scenes as discussed by [7]. Thanks to advancements in deep learning and computational capabilities, several learning-based (e.g., online learning or generalizable methods) next-best-view planners [7, 16, 22, 23, 27, 31, 49, 54, 55, 62, 78] have been developed. Among these approaches, a generalizable

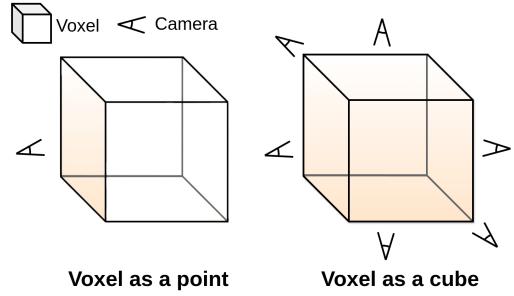


Fig. 2. A voxel is worth more than a ray. Unlike the RL-based generalizable method [7], Hestia treats each voxel as a cube by considering its six faces, rather than a point. This reduces the information loss inherent in point approximations, ensuring a more accurate representation of the voxel.

RL-based method [7] that utilizes an occupancy grid as the observation demonstrates effective performance in terms of coverage ratio, viewpoint flexibility, and generalization capability.

Inspired by the approach [7], we develop an RL-based generalizable next-best-view planner utilizing a voxel grid as the observation. While both the prior work [7] and Hestia adopt RL to learn a next-best-view policy and use voxel grids as observations, our method introduces a concept, "A voxel is worth more than a ray", by incorporating the visibility of the six faces of each voxel into the observation and reward function (see Fig. 2). By explicitly representing each voxel as a cube rather than a single point, Hestia can account for the visibility of individual voxel faces, thus achieving more comprehensive data capture.

To further help the learning process, Hestia focuses on three key aspects of generalizable RL: the observation (e.g., the diversity and scale of the dataset), the action (e.g., the dimension of the action space), and the learning process (e.g., strategies for training). For observation, Hestia utilizes the largest dataset containing 30,000 shapes, processed from Objaverse [9, 10], for the next-best-view task. This ensures that Hestia encounters a more diverse range of surface geometries, rather than just cubic shapes (e.g., buildings) as shown in Fig. S1 and Fig. S2. Such diversity is potentially critical for real-world scenarios, which may involve a wide variety of shapes. For action, instead of directly predicting the 5 DoF next-best viewpoint, a task that is inherently complex and challenging, Hestia adopts a hierarchical structure to simplify the process. The approach begins by predicting the look-at point, which serves as the target of focus, followed by determining the viewpoint position based on this point. Finally, Hestia formulates the next-best-view task as a greedy optimization problem. Given an occupancy grid, the goal is to find the viewpoint that maximizes the coverage ratio. To achieve this, Hestia relies solely on the previous image, camera pose, and the given occupancy grid to predict the next-best-view, rather than depending on a sequence of images [7, 54] and camera poses [7]. Additionally, the reward discount factor γ is set to a small value to prioritize current improvements without relying on a disproportionately large goal reward. Following a greedy algorithm, Hestia can

mitigate the spurious correlations² (see Fig. S3) arising from false associations between the current action and a large future reward (e.g., the final goal reward). By incorporating these components, Hestia can successfully learn a generalizable next-best-view policy and reach a better coverage ratio with fewer images than previous generalizable methods [7, 23, 54].

The proposed next-best-view planner is integrated into a systematic, intelligent, autonomous data collection system, with the planner running on a ground station, to demonstrate its effectiveness in real-world scenarios. The system employs a drone as an agent, equipped with a camera mounted on a gimbal, to capture data. Using drones instead of robot arms [27] practically enables data collection in regions that are difficult for humans to access and supports operations involving large-scale outdoor objects (e.g., buildings). The contributions of this work are:

- A generalizable RL-based next-best-view planner achieves a better coverage ratio with fewer images by explicitly considering voxels as cubes rather than points.
- Help the learning process through hierarchical handling of the high-dimensional action space, training on the largest dataset to ensure diversity, and adopting a greedy algorithm to minimize spurious correlations.
- Integrate the next-best-view planner into an intelligent, autonomous data collection prototype to demonstrate its potential for active data capture under real-world scenarios.

2 Related Work

Data collection systems. Advances in 3D reconstruction [11, 45, 46, 48, 50, 59, 63, 66, 68, 69, 71] and novel view synthesis [3–5, 8, 12, 13, 24, 30, 36, 40, 43, 44, 56–58, 67, 75, 79, 80, 82, 83] techniques have enabled the rendering of high-fidelity and photo-realistic scenes. However, most approaches depended on manually captured data [14, 51, 52, 74], which is both time-consuming and labor-intensive. This limitation makes the data collection process less practical for generalizable approaches [3, 5, 8, 11, 30, 50, 65, 67, 68, 75, 79, 80, 82]. To overcome this, some methods [9, 10, 33, 34, 37, 65, 81] collected data from the internet, leveraging contributions from people worldwide. While these approaches enable the rapid collection of large datasets and benefit generalizable methods, they face potential drawbacks such as inconsistent data quality and the risk of saturation due to the inability to generate their own data. Instead, many people developed their data collection machines [2, 6, 15, 28, 29, 32, 38, 72, 74] to capture high-quality data, which is a more sustainable approach in the long term. For instance, multiview camera rigs have been used to capture face-forward scenes [2, 28, 29, 32, 74], 360° outward-facing scenes [72], and 360° inward-facing scenes [38]. Other studies [6] have combined multiple camera devices (e.g., spectacles camera and 360° camera) to capture scenes from different perspectives. Additionally, ego-centric approaches [15] developed wearable camera glasses to collect data from a first-person perspective while the user was in motion. Although these methods can collect large amounts of data using multiple cameras or passively gather data through

²Spurious correlations [19, 26] refer to certain groups contributing to model errors. In this study, spurious correlations refer to large positive future rewards assigned to suboptimal current next-best-view decisions, leading to ineffective policy learning.

human-operated devices, they lack the capability to actively explore the scene. In contrast, our proposed data collection system is intelligent, enabling autonomous exploration of the scene using a drone equipped with a mounted camera. While demonstrated in an indoor environment, drones as agents, unlike robotic arms, have shown the potential to capture data in hard-to-reach areas for humans and large-scale scenes (e.g., buildings).

Scene-specific next-best-view planners. Next-best-view planners have demonstrated promising results in active 3D reconstruction by predicting the optimal viewpoint for data capture based on the current state (e.g., the reconstruction progress). Traditional approaches [17, 18, 35, 42, 77] rely on hand-crafted rules to determine the next-best viewpoint. For instance, the method [77] selected the next-best viewpoint by maximizing a rating function that prioritizes smooth regions and partially overlapping patches. Similarly, the approaches [18, 42] proposed a method to identify object boundaries and collect data along the boundary between seen and unseen object surfaces, which was later extended to multi-agent settings [17], to enable more efficient data capture. Another approach [35], inspired by human behavior, scanned segmented objects sequentially, using a predefined database of similar objects to guide the selection of the next-best viewpoint. Recent advances in deep learning and increased computational power have given rise to learning-based next-best-view methods [7, 23, 27, 31, 49, 54, 55, 62, 78]. In particular, some studies [31, 62] proposed training multiple neural radiance field (NeRF) [40] models and leveraging ensemble learning to measure the disagreement among these models as an indicator of uncertainty, which was then used to select the viewpoint that maximizes this uncertainty. Other works [49, 55] incorporated Bayesian-based NeRF [39, 61] by treating radiance estimates as Gaussian distributions and using an active learning scheme to choose the viewpoint that achieves the highest information gain through variance reduction. Meanwhile, another line of approaches [27, 78] defined the next-best viewpoint as the viewpoint that maximizes the entropy of the density field along the camera rays. Although these methods have demonstrated outstanding performance in collecting data, all require sampling candidate viewpoints to determine the next-best viewpoint. In addition, because these approaches build upon NeRF, they rely on online learning and are therefore scene-specific.

Generalizable next-best-view planners. Unlike the aforementioned online-learning approaches, the generalizable methods [7, 23, 54] avoided the training process for new scenes, thereby enabling faster next-best-view selection. The prediction time is important for real-world tasks where a robot (e.g., a drone) may run out of battery within a few minutes. Among the generalizable methods, prior work [23] developed a Bayesian-based generalizable NeRF, building upon generalizable NeRF [75] and Bayesian approaches [39, 61]. Hence, the approach [23] can select the next-best viewpoint that maximizes the variance of the view for an unknown scene without requiring further training. Another line of generalizable next-best-view approaches [7, 54] utilized reinforcement learning (RL) to learn a next-best-view planner to bypass the needs of sampling candidates viewpoints. Prior work [54] proposed learning a 3-DoF next-best-view policy using a series of grayscale images as observations and

the coverage ratio between the ground truth point cloud and the reconstructed point cloud as the reward. Subsequently, prior work [7] improved upon this method by incorporating occupancy grids into the observations, which provide explicit geometric information. This enhancement enabled the development of a 5-DoF next-best-view planner, achieving an outstanding coverage ratio for unknown scenes. Although the approach [7] achieved impressive results in viewpoint flexibility and coverage ratio, its simplification of treating each voxel as a point overlooked the geometry of the surface.

Different from voxel-based generalizable RL methods, Hestia integrates the concept, "*A voxel is worth more than a ray*", into the voxel grid to enable more comprehensive data collection (see Fig. 2). Unlike prior works [7, 54], which were trained on the cubic-like House3K dataset, Hestia is trained on a larger, more diverse dataset of approximately 30,000 shapes from Objaverse [9, 10] (see Fig. S1 and Fig. S2), improving robustness in real-world scenarios. Hestia predicts the camera's look-at point first and determines the viewpoint position later, simplifying the 5-DoF action space and surpassing the 3-DoF constraints of prior works [54]. It employs a greedy optimization scheme with a small discount factor to maximize the coverage ratio and mitigate spurious correlations. Unlike methods [7, 54] that depend on sequences of images and large goal rewards, Hestia operates effectively with minimal dependencies. Integrated into a real-world data collection system, Hestia demonstrates its feasibility in practical scenarios.

3 Methods

3.1 RL Problem Definition

Our next-best-view task is to identify a 5-DoF viewpoint that maximizes the coverage ratio of an incomplete occupancy grid of the scene. The task's goal is similar to greedy methods, which always seek the locally optimal solution. We formulate the problem as a Markov Decision Process (MDP), denoted by the tuple $\{S, A, P, R, \gamma\}$. At each time step t , the agent (e.g., a drone) with RGB-D camera observes a state s_t from the set of all possible states S and chooses an action a_t from the action space A . The environment then transitions to the subsequent state s_{t+1} according to the probabilities described by P , and provides a reward r_t . The magnitude of this reward is determined by the reward function

$$R(\cdot | s, a) : S \times A \rightarrow r. \quad (1)$$

In reinforcement learning (RL), the main goal is to discover an optimal policy π that maximizes the expected sum of discounted rewards, given by

$$E_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}, \quad (2)$$

where $\gamma \in (0, 1)$ is the discount factor. We set the discount factor, γ , to 0.01 to align with the greedy-like objective and to avoid spurious correlations from large positive future rewards (see Fig. S3).

State space. The state space of Hestia is defined as

$$S = \left\{ s_t \mid s_t = \{I_t, M_t, G_t, L_t\}, t \in \mathbb{N} \right\} \quad (3)$$

where $I_t \in \mathbb{R}^{h \times w}$ is the grayscale image with height h and width w , and $L_t \in \mathbb{R}^3$ is the camera look-at point. The vector $M_t \in \mathbb{R}^6$ consists of $X_t \in \mathbb{R}^3$, which is the camera position, pitch, and yaw, as well as $H_t \in \mathbb{R}^1$, representing the maximum flyable height for the capture. Meanwhile, $G_t \in \mathbb{R}^{g \times g \times g \times 10}$ includes the aggregated grid information at resolution g , consisting of $O_t \in \mathbb{R}^{g \times g \times g \times 1}$ for the cumulative occupancy grid, $C_t \in \mathbb{R}^{g \times g \times g \times 3}$ for the positional encoding, and $F_t \in \{0, 1\}^{g \times g \times g \times 6}$ for the cumulative face visibility. The cumulative face visibility is updated iteratively as

$$F_t = f_t \vee F_{t-1} \quad (4)$$

where F_t represents the cumulative face visibility for all voxels up to time t , and $f_t \in \{0, 1\}^{g \times g \times g \times 6}$ denotes the face visibility for the current timestep t . To compute f_t , the depth image D_t is unprojected into a voxelized point cloud $V = \{v_i \mid i \in \mathbb{N}\}$, where v_i is the i -th voxel. Each voxel v_i is associated with a viewing direction vector $d_{v_i} \in \mathbb{R}^3$, defined as the vector pointing from the voxel center to the collision-free camera position a'_t . The vector d_{v_i} is computed as

$$d_{v_i} = \frac{a'_t - p_{v_i}}{\|a'_t - p_{v_i}\|} \quad (5)$$

where p_{v_i} is the center of voxel v_i . For each voxel v_i and its six outward-facing face normals $n_{i,j} \in \mathbb{R}^3$, the face visibility is determined and aggregated as

$$f_t(v_i, j) = \mathbb{1}(d_{v_i} \cdot n_{i,j} > 0), \quad \forall v_i \in V, j \in \{1, \dots, 6\} \quad (6)$$

where $\mathbb{1}(\cdot)$ is the indicator function. By iterating over all voxels and their respective faces, f_t is constructed for the current timestep, and the cumulative visibility F_t is updated accordingly. Although this method cannot handle all face visibilities, the approximation enables efficient computation of face visibility. Unlike prior works [7], which consider only O_t and C_t and thereby treat voxels as points, we treat voxels as cubes to mitigate the information loss caused by approximating voxels as points (see Fig. 2). For details regarding O_t and C_t , please refer to the work [7].

Action space. The action space

$$A = \left\{ a_t \mid a_t \in [-1, 1]^3, t \in \mathbb{N} \right\} \quad (7)$$

represents the set of possible 3-DoF viewpoints (e.g., camera positions) at each time step t , where each coordinate is initially bounded within $[-1, 1]$. These coordinates are subsequently normalized to the environment's scale to ensure appropriate positioning within the scene. Additionally, the camera's pitch and yaw are derived from the look-at point and the collision-free action a'_t converted from a_t (see Sec. 3.2).

Reward. The reward function is defined as

$$r_t = R(s_t, a_t) = r_{\text{coverage}}(s_t, a_t) + r_{\text{constraint}}(s_t, a_t) \quad (8)$$

where $r_{\text{coverage}}(s_t, a_t)$ encourages the observation of new voxel faces and is expressed as

$$r_{\text{coverage}}(s_t, a_t) = \frac{\sum_{i=1}^N \sum_{j=1}^6 (F_t^{i,j} - F_{t-1}^{i,j}) \cdot M_{\text{col}}}{N \cdot 6} \cdot 0.3 \quad (9)$$

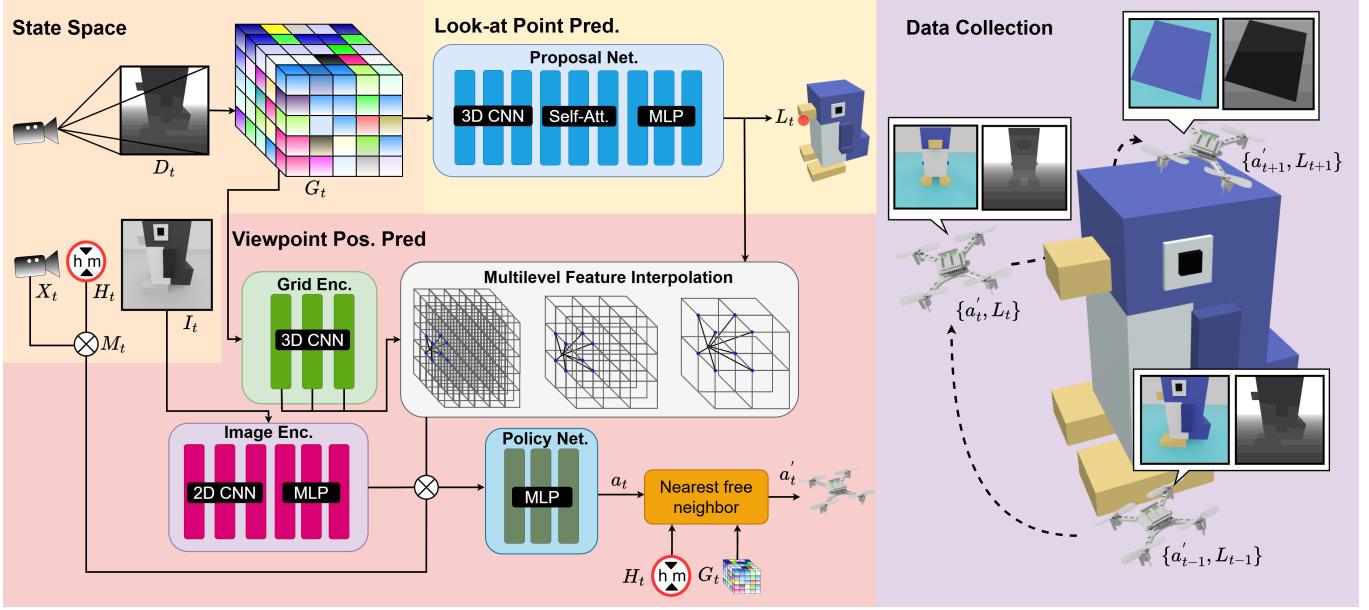


Fig. 3. Hierarchical structure of Hestia. Hestia first predicts the camera’s look-at point L_t using a proposal neural network that takes grid information G_t processed from the depth image D_t and the camera pose as input. Next, Hestia employs a grid encoder to encode the grid information G_t and performs trilinear interpolation to extract corresponding features from the encoded grid at different layers based on L_t . These multilevel interpolated features are then concatenated with the vector information M_t which includes the camera pose X_t and the maximum flyable height, H_t as well as the encoded image features. The image features are extracted using an image encoder, which takes the grayscale image I_t as input. Finally, this combined feature representation is fed into the RL policy model to predict the camera’s position a_t . Note that Hestia adopts a'_t , the nearest collision-free point to a_t , as the final camera position to ensure a collision-free viewpoint. Hence, the next-best viewpoint $\{a'_t, L_t\}$ is used for data collection.

where $F_t^{i,j}$ and $F_{t-1}^{i,j}$ represent the visibility status of the j -th face of the i -th voxel at time t and $t - 1$, respectively. Here, $M_{\text{col}} \in \{0, 1\}$ is a collision indicator, set to 0 in the event of a collision, thereby preventing any positive reward for invalid actions. The term $r_{\text{constraint}}(s_t, a_t)$ penalizes unsafe and invalid actions and is defined as

$$r_{\text{constraint}}(s_t, a_t) = \begin{cases} -0.01, & \text{if } r_{\text{coverage}}(s_t, a_t) = 0, \\ & \text{or } a_t[2] > H_t, \\ & \text{or } a_t \in \text{non-free voxels}, \\ 0, & \text{otherwise.} \end{cases} \quad (10)$$

Our reward is based on the face coverage ratio rather than the point coverage ratio to ensure more comprehensive capture (see Fig. 2). Furthermore, to prevent spurious correlations, the reward design aligns with a greedy-like objective, which differs significantly from prior works [7, 54] that provide a large goal reward when the coverage ratio reaches a predefined target.

3.2 Next-Best-View Hierarchical Network

The goal of the next-best-view task is to predict a 5-DoF viewpoint for data collection. Modeling the 5-DoF viewpoint directly in the RL continuous action space is challenging due to the high-dimensional search space and complexity. To address this, Hestia introduces a hierarchical structure to simplify the problem.

Look-at point prediction. Hestia first predicts the camera’s look-at point using a proposal network (see Fig. 3), which takes the grid information G_t as input. The proposal network is a shallow 3D convolutional neural network that employs a self-attention layer to expand the receptive field. The output is then passed through linear layers to decode the look-at point L_t . To model the look-at point as a probability distribution, the reparameterization trick is used, treating it as a sample from a normal distribution.

Viewpoint position prediction. To predict the remaining 3-DoF viewpoint position, the grid information is encoded into a multilevel feature grid using a shallow 3D CNN. The look-at point L_t is then used to perform trilinear interpolation on the multilevel features from the grid. These interpolated features are concatenated with the image embedding, which is extracted by an image encoder, a shallow convolutional neural network that takes the grayscale image I_t as input. Additionally, the features are concatenated with vector information M_t , which includes the camera pose X_t and the maximum flyable height H_t . The combined features are fed into the RL policy model to predict the action a_t . While the reward function helps constrain a_t to avoid collisions, an additional constraint is applied to ensure a collision-free viewpoint. Specifically, a_t is shifted to its nearest collision-free point a'_t determined using G_t and H_t . This adjusted action a'_t serves as the final viewpoint position for data capture. Thus, Hestia’s next-best-view is represented as $\{a'_t, L_t\}$.

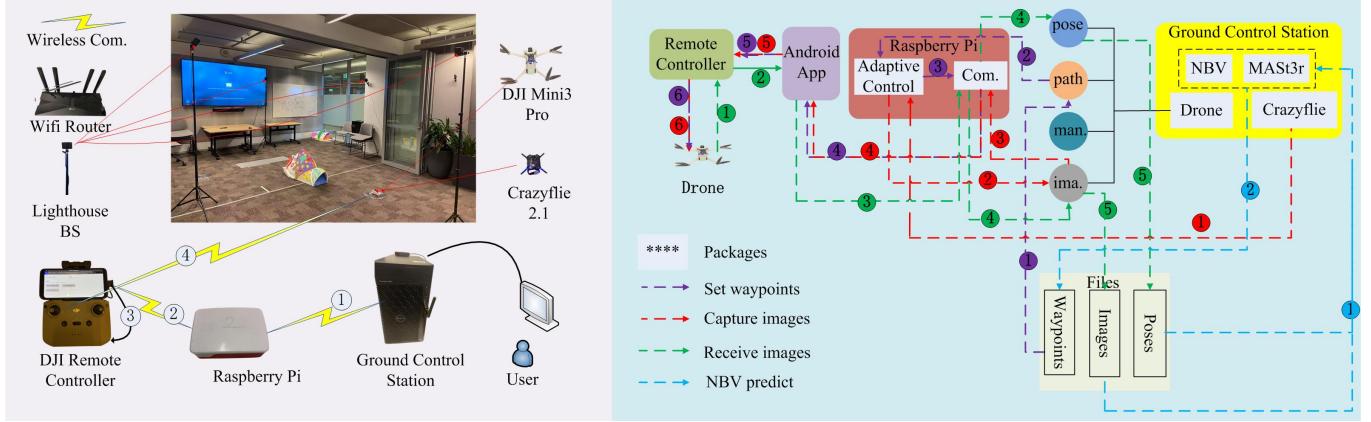


Fig. 4. The proposed data collection system and flowchart. Hestia goes beyond prior works by introducing a prototype of an intelligent data collection system utilizing a drone equipped with an RGB camera for data capture, HTC Lighthouse base stations and Crazyflie 2.1 for localization, and a WiFi router for wireless communication. Although tested indoors, this system highlights the potential of drones with next-best-view planners for large-scale data collection.

Training loss functions. Fully joint learning of two continuous action spaces through RL is complex and unnecessary for our task. The look-at point prediction network can instead be trained using supervised learning with a ground truth target. The ground truth look-at point L_t^{gt} is computed as the weighted average position of the ground truth uncaptured surface

$$L_t^{\text{gt}} = \frac{\sum_{v_i \in U} w_{v_i} p_{v_i}}{\sum_{v_i \in U} w_{v_i}} \quad (11)$$

where U represents the set of voxels containing ground truth uncaptured faces, and w_{v_i} is defined as the total number of ground truth uncaptured faces within voxel v_i

$$w_{v_i} = \sum_{f \in F_{v_i}^{\text{gt}}} 1 \quad (12)$$

where $F_{v_i}^{\text{gt}}$ is the set of ground truth uncaptured faces associated with voxel v_i . Thus, the loss function for the proposal network is formulated as

$$\mathcal{L}_{\text{proposal}} = \|L_t - L_t^{\text{gt}}\|^2 \quad (13)$$

The loss of the viewpoint prediction network is the same as the regular RL loss \mathcal{L}_{RL} which depends on the RL method used, combined with an auxiliary loss

$$\mathcal{L}_{\text{aux}} = \|a_t - a'_t\|^2 \quad (14)$$

to encourage the predicted action a_t to align with the collision-free action a'_t . Hence, the overall loss function for Hestia is

$$\mathcal{L}_{\text{all}} = \mathcal{L}_{\text{RL}} + 0.5 \cdot \mathcal{L}_{\text{aux}} + \mathcal{L}_{\text{proposal}} \quad (15)$$

4 Real-world Data Collection System

To demonstrate our work, we propose a real-world system (see the left part of Fig. 4), where a drone equipped with an RGB camera moves to the next-best viewpoint predicted by Hestia to capture images of an object. The system uses four HTC Lighthouse base stations and a Crazyflie 2.1 for localization and transmits images to the ground control station via wireless communication. MAS3R [11]

is integrated to convert RGB images into pointmaps (e.g., depth images), and three initial viewpoints are set for real-world and virtual-world synchronization [64]. The process (see Supp. Alg. 1) comprises four key processes, shown in the right part of Fig. 4 with four colored arrows and described in Supp. Alg. 2 to Supp. Alg. 5.

Specifically, the drone captures images at three initial viewpoints (see Supp. Alg. 1), where, for each viewpoint, process "set waypoints" (see Supp. Alg. 2) navigates the drone to a set of waypoints. Then, process "capture image" (see Supp. Alg. 3) starts the image capturing once the drone reaches the waypoint, and process "receive image" (see Supp. Alg. 4) transmits the captured image to the ground station. After receiving all the images captured at the initial viewpoints, process "nbv prediction" (see Supp. Alg. 5) predicts the next-best viewpoint based on the collected data. The image capturing and next-best viewpoint prediction loop will be repeated until sufficient images have been collected (see lines 8-13 in Supp. Alg. 1). For more details, please refer to the Supp. Sec. 4.

5 Experiments

5.1 Experimental Setup

Dataset. We use Objaverse [9, 10] as the training dataset to maximize the diversity of shapes encountered during training (see Fig. S1 and Fig. S2). For each shape, we generate the ground truth occupancy grid, voxel face visibility, and point cloud (see Supp. Sec. 2 for more details). Our processed training set is two orders of magnitude larger than those used in prior works [7, 23, 54] and includes at least eighteen more categories [7, 54].

Scene setup. We use NVIDIA IsaacLab [41] to create scenes for training, running 256 scenes in parallel. Each sample is scaled to a maximum side length of 8 meters and placed within a 20m x 20m x 20m scene. During training, the object is randomly translated within the scene. For testing, the object is positioned at its original location, and the four corners for benchmarking. An RGB-D camera is mounted on the Crazyflie model, positioned 0.1 meters in front of

- [78] Huangying Zhan, Jiyang Zheng, Yi Xu, Ian Reid, and Hamid Rezatofighi. 2022. Activemap: Radiance field for active mapping and planning. *arXiv preprint arXiv:2211.12656* (2022).
- [79] Kai Zhang, Sai Bi, Hao Tan, Yuanbo Xiangli, Nanxuan Zhao, Kalyan Sunkavalli, and Zexiang Xu. 2024. GS-LRM: Large Reconstruction Model for 3D Gaussian Splatting. *European Conference on Computer Vision* (2024).
- [80] Shunyuan Zheng, Boyan Zhou, Ruizh Shao, Boning Liu, Shengping Zhang, Liqiang Nie, and Yebin Liu. 2024. GPS-Gaussian: Generalizable Pixel-wise 3D Gaussian Splatting for Real-time Human Novel View Synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [81] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. 2018. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817* (2018).
- [82] Chen Ziwen, Hao Tan, Kai Zhang, Sai Bi, Fujun Luan, Yicong Hong, Li Fuxin, and Zexiang Xu. 2024. Long-LRM: Long-sequence Large Reconstruction Model for Wide-coverage Gaussian Splats. *arXiv preprint 2410.12781* (2024).
- [83] Zi-Xin Zou, Zhipeng Yu, Yuan-Chen Guo, Yangguang Li, Ding Liang, Yan-Pei Cao, and Song-Hai Zhang. 2024. Triplane meets gaussian splatting: Fast and generalizable single-view 3d reconstruction with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10324–10335.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

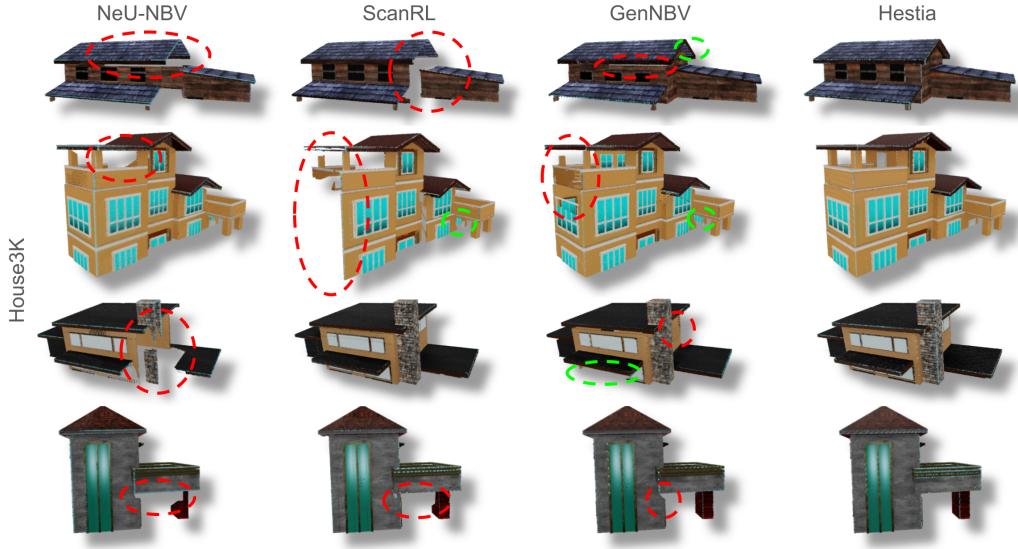


Fig. 5. **Qualitative comparison on the House3K [54].** Compared to the baselines, the point cloud reconstructed from the depth maps collected using our approach captures more fine-grained details, such as the roof soffit, particularly in self-occluded areas.

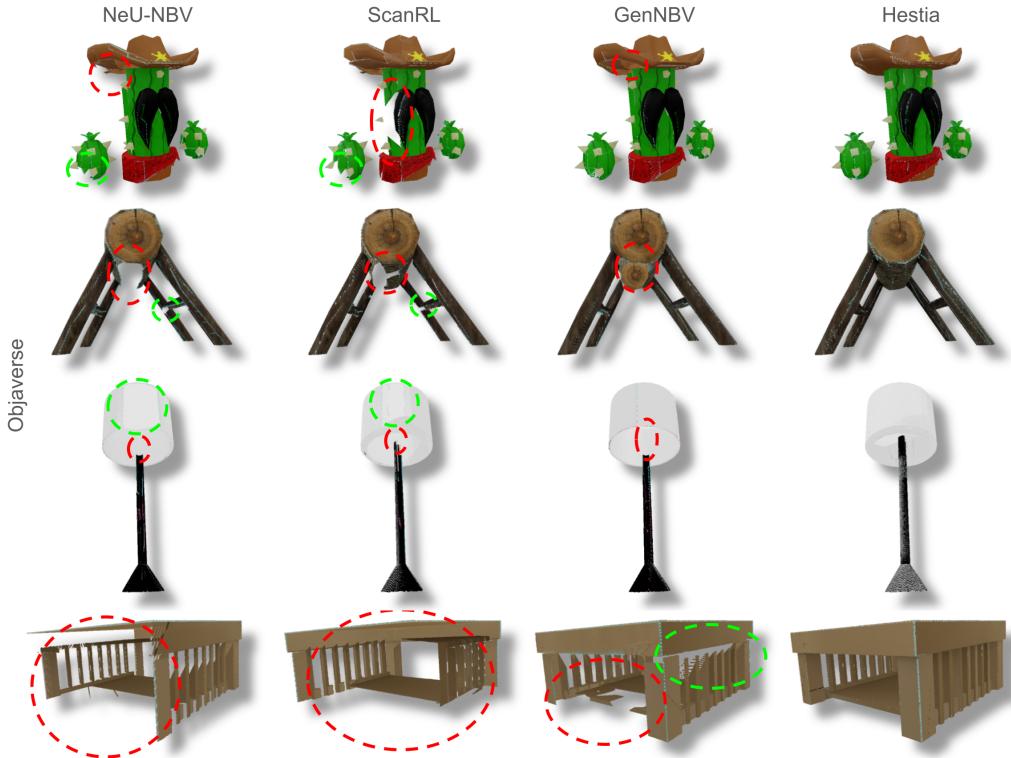


Fig. 6. **Qualitative comparison on the Objaverse [9, 10].** Compared to the baselines, the point cloud reconstructed from the depth maps collected using our approach captures finer details from bottom-up views.

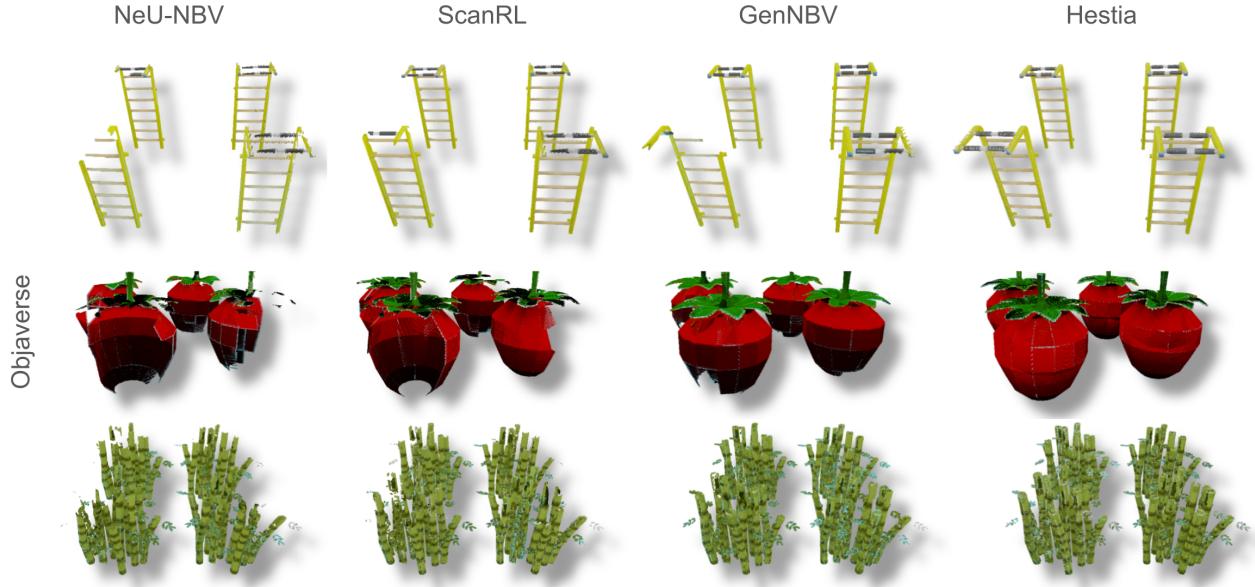


Fig. 7. Qualitative comparison of objects at the four corners from Objaverse [9, 10]. Compared to the baselines, the point clouds reconstructed from the depth maps collected using our approach are more robust and consistent across different object centers.

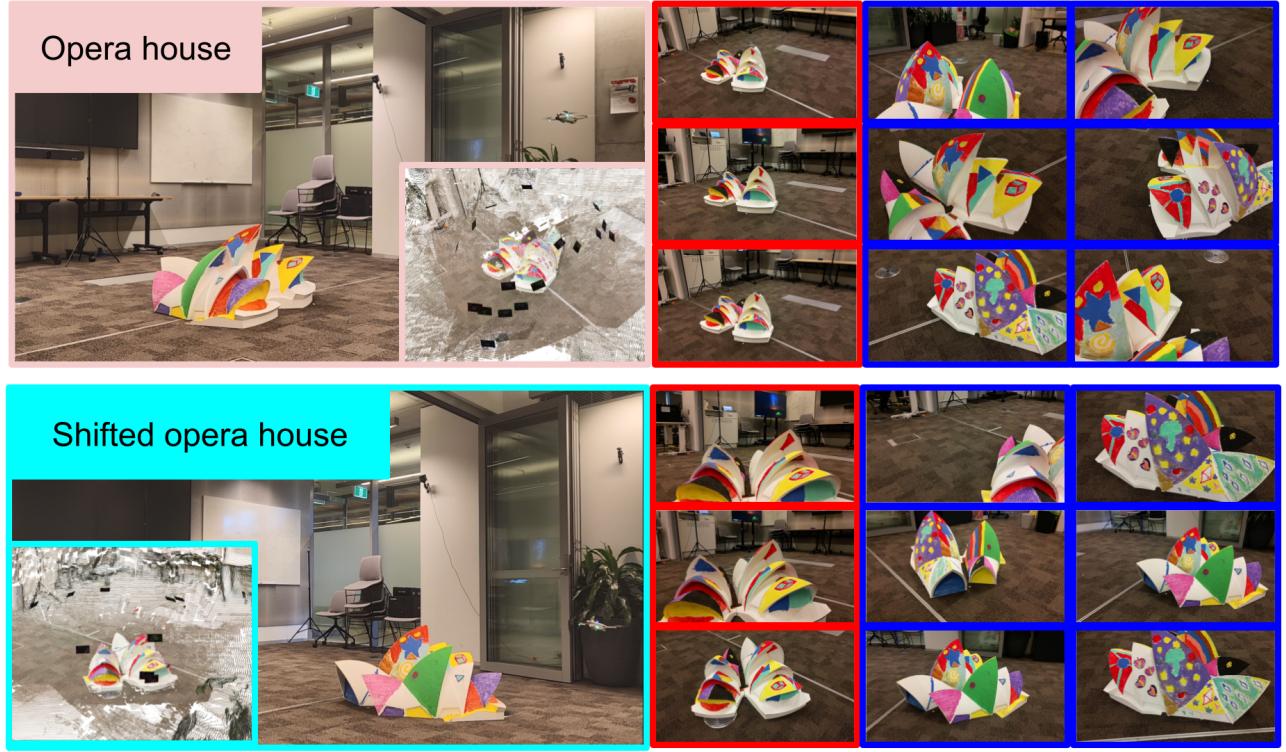


Fig. 8. Real-world demonstration. Hestia can operate in real-world scenarios with three initial viewpoints (red rectangles) and predict the next-best viewpoints for capture (blue rectangles for the first six views), even when the depth camera is unavailable. Point cloud reconstruction results are shown on the left, with black rectangles representing camera poses.