# Learning Reconstructability for Drone Aerial Path Planning

YILIN LIU, Shenzhen University, China
LIQIANG LIN, Shenzhen University, China
YUE HU, Shenzhen University, China
KE XIE, Shenzhen University, China
CHI-WING FU, The Chinese University of Hong Kong, China
HAO ZHANG, Simon Fraser University, Canada
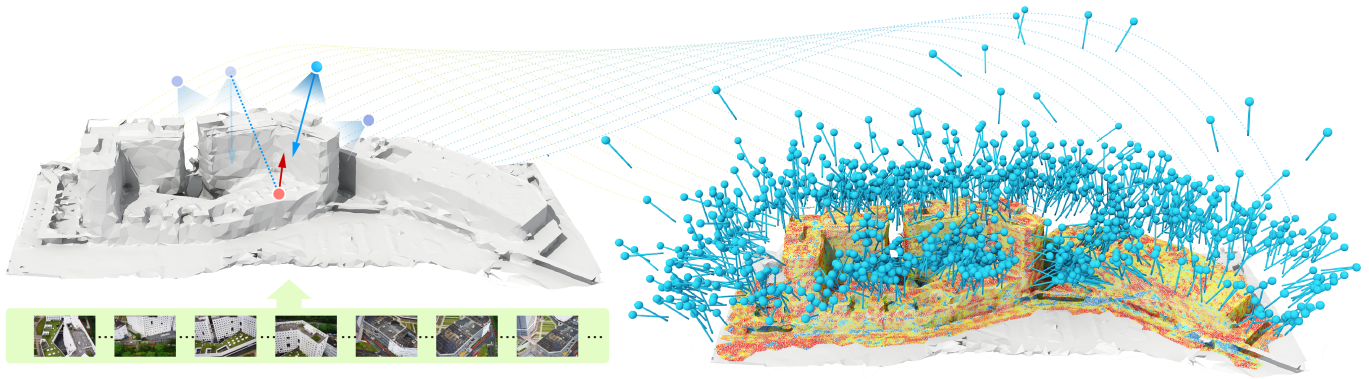HUI HUANG*, Shenzhen University, China

Fig. 1. We train a neural network to predict reconstructability for drone path planning during 3D urban scene acquisition. The prediction is based on a rough scene proxy, a set of viewpoints (bluish dots on the top), and optionally a series of images captured during the pre-flight pass, as shown on the left. The network learns both image features and viewpoint features from the perspective of sample points (red dot) on the proxy, while the predicted reconstructability (color) map, shown on the right, guides our iterative view planner to execute the onsite drone view acquisition for 3D reconstruction.

We introduce the first *learning-based reconstructability predictor* to improve view and path planning for large-scale 3D urban scene acquisition using unmanned drones. In contrast to previous heuristic approaches, our method learns a model that *explicitly* predicts how well a 3D urban scene will be reconstructed from a set of viewpoints. To make such a model trainable and simultaneously applicable to drone path planning, we simulate the proxy-based 3D scene reconstruction during training to set up the prediction. Specifically, the neural network we design is trained to predict the scene reconstructability as a function of the *proxy geometry*, a set of viewpoints, and optionally a series of scene images acquired in flight. To reconstruct a new urban scene, we first build the 3D scene proxy, then rely on the predicted reconstruction quality and uncertainty measures by our network, based off of the proxy geometry, to guide the drone path planning. We demonstrate that our data-driven reconstructability predictions are more closely correlated to the true reconstruction quality than prior heuristic measures. Further, our learned predictor can be easily integrated into existing path planners to yield improvements. Finally, we devise a new iterative view planning framework, based on the learned reconstructability, and show superior performance of the new planner when reconstructing both synthetic and real scenes.

*Corresponding author: Hui Huang (hhzhiyan@gmail.com)

Authors' addresses: Yilin Liu, whatsevenlyl@gmail.com, Shenzhen University, China; Liqiang Lin, liniquie@gmail.com, Shenzhen University, China; Yue Hu, hytraveler2000@gmail.com, Shenzhen University, China; Ke Xie, ke.xie.siat@gmail.com, Shenzhen University, China; Chi-Wing Fu, philip.chiwing.fu@gmail.com, The Chinese University of Hong Kong, China; Hao Zhang, hao.r.zhang@gmail.com, Simon Fraser University, Canada; Hui Huang, hhzhiyan@gmail.com, College of Computer Science & Software Engineering, Shenzhen University, China.

## 1 INTRODUCTION

Aerial path planning for large-scale 3D urban scene acquisition using unmanned drones has gained much attention recently [Hepp et al. 2018b; Liu et al. 2021a; Roberts et al. 2017; Smith et al. 2018; Zhang et al. 2021; Zhou et al. 2020]. The ultimate goal of the planning problem is to best reconstruct the 3D scenes in terms of completeness and accuracy while respecting physical constraints imposed by the

drones' flight speed and battery life. As the scene complexity scales up, it becomes more difficult to model the spatial relations between a scene and its viewpoints, and more importantly, the *uncertainty* surrounding scene reconstruction quality which forms the foremost criterion for path planning. Since the ground-truth 3D scene is unknown during planning, the reconstruction quality, or specifically, scene *reconstructability* with respect to a set of viewpoints of the drones, must be estimated using imperfect input data.

Current approaches to estimating reconstructability all rely on heuristics in one form or another. Some account for scene coverage [Giang et al. 2021; Schmid et al. 2020] or viewpoint correlations [Furukawa et al. 2010; Koch et al. 2019] without scene reconstruction. Others optimize path planning based on a coarse *proxy reconstruction* [Hepp et al. 2018b; Roberts et al. 2017; Smith et al. 2018; Zhang et al. 2021] obtained by an initial drone fly through. However, scene coverage and view correlation represent measures that are only *relevant to* reconstruction quality; they do not explicitly model it or strictly validate heuristic designs against it, even when the data is complete. On the other hand, the pre-constructed proxy geometries are typically coarse and inaccurate and these inaccuracies are easily propagated to the coverage and correlation estimates so as to misguide the ensuing path planning.

We formulate reconstructability estimation as a *predictive* task and introduce the first data-driven approach to *learn* reconstructability for drone aerial path planning. More formally, reconstructability measures how well the local area around a sample point in the input 3D scene can be reconstructed from information captured at a set of viewpoints. In contrast to previous heuristic approaches, our learned model *explicitly* predicts how well a 3D urban scene will be reconstructed. To make such a model trainable and simultaneously applicable to drone path planning, we simulate the *proxy-based* 3D scene reconstruction during training to set up the prediction. Specifically, our training data contains ground-truth reconstructability values per sample point on proxy geometry subject to a given viewpoint set, so that we can train the network to predict final "reconstruction quality" at the sample points, as a function of the *proxy geometry* and the viewpoints. Note that technically, the reconstruction quality is measured against the ground-truth scene, but it is unknown during inference; only the scene proxy is available.

As an extension to our learning framework, we leverage additional image inputs to refine the reconstructability prediction, since high-quality images can often be acquired during the drone preflight. These images can better capture scene details than would be possible via proxy construction, thus helping to predict *uncertainty-aware* inaccuracies over the proxy geometry. To reconstruct a new scene, we first build the 3D proxy, then rely on the network predictions and uncertainty measures, based off of the proxy geometry, to guide the view planner to find a set of viewpoints to maximize the predicted reconstruction quality for 3D scene reconstruction.

Our learning model consists of an attention-based view fusion network for reconstructability prediction. As such, the influence of the different factors related to a single viewpoint on scene reconstruction, such as viewing distances and viewing angles (with respect to surface normals over the scene geometry), together with the correlations between multiple viewpoints, such as their scale differences and baselines (i.e., distances), are all adaptively adjusted

by the learned parameters. Further, we develop another attention-based *image* fusion network to implicitly model the uncertainty of the scene geometry with respect to the acquired image observations and refine the *spatial* reconstructability that is learned from spatial relations between the viewpoints and scene proxy; see e.g., Fig. 1.

We demonstrate the effectiveness of our learning framework through extensive experiments, both quantitatively and qualitatively. We verify that our learned reconstructability more closely correlates to the true reconstruction quality than prior heuristic measures. Important for immediate practical impact, our reconstructability predictor can be easily integrated into state-of-the-art path planners [Smith et al. 2018; Zhou et al. 2020], leading to improved quality for large-scale 3D urban scene acquisition.

Finally, we complete the loop by devising a new *iterative viewpoint optimization* framework, based on the learned reconstructability, to further improve path planning. Specifically, we adjust the current viewpoints along with the path planning process to attain better reconstructability, where the adjustments include viewpoint insertion near under-reconstructed regions, deletion of redundant views, and altering the position and orientation of existing views. We show that our new adaptive scheme, built on a more accurate reconstructability prediction, can help escape local minima during path planning, a reoccurring issue which has challenged existing planners. The new planner exhibits superior reconstruction performance over existing methods on both synthetic and real scenes.

## 2 RELATED WORK

Unmanned drones have been widely employed for urban scene acquisition due to their maneuverability and large fields of view. During an acquisition, a drone usually flies along a pre-computed path, which is generated to optimize a quality measurement. Predominantly, such a measurement is related to the completeness and accuracy of the 3D urban scene to be reconstructed, i.e., reconstructability, with respect to a set of viewpoints and paths selected. However, since the ground-truth geometry is unknown, all methods must estimate the reconstructability measure.

### 2.1 Estimates of reconstructability

*View coverage and uncertainty.* One line of approaches to path optimization is based on view coverage [Giang et al. 2021; Schmid et al. 2020] by a depth sensor. Schmid et al. [2020] proposed a spatial uncertainty measure based on viewing distance. A rapid random search tree was developed to maintain the measure, facilitating the search for an optimal scanning path. Song et al. [2020] divide the acquisition process into two steps: global planning and local inspection. During global planning, they also rely on a scene uncertainty measure to generate a rough initial path. This is followed by solving a set cover problem to optimize the local viewpoints so as to capture more geometric details. Note that when planning paths to better cover the target urban scenes, these methods all consider view coverage with respect to individual viewpoints. On the other hand, methods based on multi-view stereo (MVS), e.g., [Furukawa et al. 2010; Peng and Isler 2019; Smith et al. 2018], often need to account for spatial relations between the viewpoints, since the errors

stemming from triangulation and feature matching all depend on the relative positions of the viewpoints and the scene geometry.

*View correlation.* To this end, several measures have been proposed to characterize correlations between a *pair* of viewpoints in order to model reconstruction quality. Furukawa et al. [2010] assumed that the quality measure follows a piecewise Gaussian distribution, which depends on viewpoint baselines and the pixel densities. Smith et al. [2018] decomposed this measure into two components, which respectively account for feature matching and triangulation in the context of MVS. Furthermore, their work defines reconstructability as an accumulative product of Gaussian functions, which are defined in terms of viewpoint baselines, view distances, and viewing angles. In addition, Peng and Isler [2019] also considered the impact of different view sampling rates when performing feature matching and dense reconstruction. Finally, Koch et al. [2019] factored in viewpoint overlaps during optimization.

While the above methods all consider view pairs when measuring reconstruction errors, in real MVS reconstruction, multiple viewpoints visible to a surface point would contribute to its nearby reconstruction. Roberts et al. [2017] extended the correlation model to encompass a *set* of viewpoints. Specifically, they proposed a measure based on spherical integration, which considers the impact of all visible viewpoints. The integral function is related to viewing distances, viewpoint baselines, and viewing angles. However, like other methods, which also consider these unary and relational viewpoint attributes, various assumptions have to be made, resulting in a variety of parameters that are difficult to tune in practice.

*Scene proxy.* During path planning, most methods up to now obtain the various measures with respect to a rough scene proxy obtained either via a rapid pre-fly and rough reconstruction [Hepp et al. 2018b; Roberts et al. 2017; Smith et al. 2018; Zhang et al. 2021], i.e., the scene proxy, or by an extraction from geological features [Zhou et al. 2020]. The inaccuracies or uncertainties of the proxy scene geometry, especially pertaining to surface normal estimation, would significantly impact the view planning. Peng and Isler [2019] developed a three-step scene reconstruction method, by iteratively finding reconstructed regions that have the lowest confidence and conducting path planning for them, to improve reconstruction quality. However, this method needs several drone flights to obtain a satisfactory reconstruction, leading to high acquisition costs.

*Data-driven methods.* Recently, the rapid proliferation of 3D scene datasets [Chang et al. 2015; Huang et al. 2019; Knapitsch et al. 2017; Lin et al. 2022; Liu et al. 2021b] have enabled data-driven methods to model correlations between viewpoints and scene geometry. Genova et al. [2017] proposed such a method for view set selection, whereby a set of views are generated for a synthetic dataset to match the content statistics of a set of example images. Sun et al. [2021] designed a neural network to model the visibility and quality of viewpoints, turning the traditional discrete viewpoint optimization problem into a continuous one. At last, Hepp et al. [2018a] utilized voxel maps for encoding the correlation between viewpoints and scene geometry, allowing one to predict the quality of the next viewpoint. However, these methods only model visibility or the quality of a single view, and they are still limited to depth-based reconstruction.

On the other hand, transformer [Vaswani et al. 2017] is an effective means for extracting the correlation between data. It has been extensively employed in machine translation [Vaswani et al. 2017], as well as stereo depth estimation [Li et al. 2021] and multi-view reconstruction [Ding et al. 2022]. In our work, we adopt transformers to fuse the geometric relations between multiple viewpoints and the image information in the reconstructability measurement.

## 2.2 Path planning

Based on the various reconstructability estimates, different path planners have been proposed to optimize viewpoint configuration for urban scene acquisition, where most of them [Roberts et al. 2017; Smith et al. 2018; Zhou et al. 2020] uniformly initialize a set of viewpoints as optimization candidates. Roberts et al. [2017] first coarsen the view optimization by determining the optimal direction for each viewpoint and finding the additive approximation of it. Then a standard integer linear program solver is employed to solve such an *orienteering* problem to obtain the optimal trajectory. Similarly, Smith et al. [2018] generate an initial trajectory at a fixed height with nadir view orientation. Then they use their reconstructability measurement as the objective function to iteratively identify whether a new position and orientation for each viewpoint are better. To this end, they resort to the Nelder-Mead method to find the global minimum of their objective function. Hepp et al. [2018b] voxelized the 3D safe space and define their objective function as the information gain towards an unknown environment. More recently, Zhou et al. [2020] leverage a dense initialization of the viewpoints and assume the initialization to be perfect but redundant for the reconstruction. Then they define the view contribution of each viewpoint according to their reconstructability measure and iteratively reduce redundant viewpoints to obtain the optimal subset of viewpoints.

On the other hand, some path planners directly generate a trajectory without any viewpoint initialization. Zhang et al. [2021] maintain and expand a rapidly-exploring random tree of the scene to directly obtain a more efficient trajectory with sufficient reconstructability. Liu et al. [2021a] generate and update the image acquisition path in real time through certain pre-defined trajectory patterns on a coarse scene proxy. However, the binary visibility function from a viewpoint to a surface point and the correlation between different viewpoints make this problem non-convex, hence hard to optimize in practice. Minimizing the objective function in an iterative way is prone to be stuck in local minima.

In our work, we develop the first *learned* reconstructability predictor, and along with an associated view optimization scheme, we can improve the performance of existing path planners.

## 3 OVERVIEW

Our learning-based framework consists of two phases: (i) *training phase* and (ii) *inference phase.* In the training phase, we prepare training data using the UrbanScene3D dataset [Lin et al. 2022] and train our neural network model to predict scene reconstructability on the proxy geometry against ground-truth information extracted from the data. Then, in the inference phase, we can integrate our trained network model, as a *reconstructability predictor*, into existing view planners for calculating the scene reconstructability. Further,
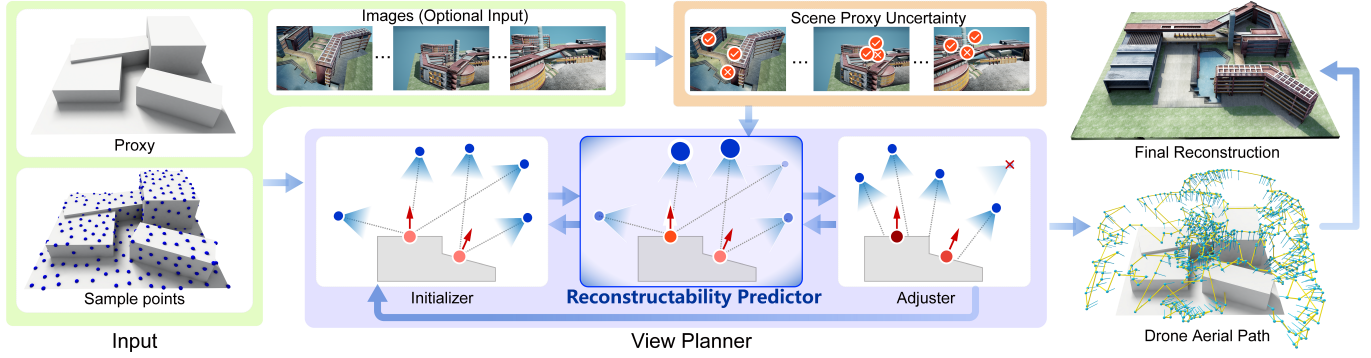
Fig. 2. *Left:* our approach takes the proxy geometry of the target scene, sample points (in red) on the proxy, images captured by a pre-flight, and the camera poses as input. *Middle:* our approach consists of a *reconstructability predictor* and a *view planner*. The *reconstructability predictor* first extracts spatial features between the sample points and the viewpoints (in blue). If the images and camera poses in the pre-flight are available, our *reconstructability predictor* can extract the scene proxy uncertainty near this sample point from the images and predict the uncertainty-aware spatial reconstructability for the view planner. The *view planner* takes the output of the *reconstructability predictor* and iteratively optimizes the number and poses of each viewpoint to maximize their ability to reconstruct the target scene. *Right:* the trajectory obtained by our method and the corresponding reconstruction result.

to address the limitations of the existing planners, we formulate a new view planning framework that combines the strengths of the existing ones to produce view configurations.

*Network inputs.* To train our neural network, we first prepare network inputs based on a given set of viewpoints in the free space and a given set of sample points on the scene proxy; see Fig. 2 (left). On the one hand, we explore every visible (sample) point-view pair, considering their locations and orientations, and encode their spatial relation geometrically as a 5D point-view feature vector. These spatial relations provide hints on how well the image captured at each viewpoint would contribute to reconstructing the local geometry surrounding the sample point. On the other hand, we extract *image features* from each pre-captured image upon its availability.

*Network predictions.* From the network inputs, we design our neural network to first extract point-view spatial features through MLPs (Multi Layer Perceptions); then, our network adopts a transformer encoder to explore the correlations across sample points and viewpoints, enabling us to better fuse features from different sample points for predicting the spatial reconstructability at each sample point on the scene proxy (Sec. 4.1). On the other hand, upon the availability of pre-captured images, our network also extracts image features and fuses these features with the point-view features to enable us to predict uncertainty-aware spatial reconstructability (Sec. 4.2). By this means, we can better account for the *inaccuracy* in scene proxy in the reconstructability prediction.

*View planning.* Last, we integrate our trained network into existing view planners as a measure for scene reconstructability (Sec. 5.1). By doing so, we found limitations of two state-of-the-art planners, [Smith et al. 2018] and [Zhou et al. 2020], on optimizing viewpoints for scene acquisition. Hence, we further formulate a new view planner (Sec. 5.2), collectively combining their complementary strengths by iteratively initializing, eliminating, and adjusting

viewpoints, as guided by our trained network, to obtain a view configuration with maximized reconstructability; see Fig. 2.

In the end, we evaluate our reconstructability predictor and view planner in both unseen synthetic and real scenes. Results presented in Sec. 6 show that the proposed reconstructability predictor can better reflect the final reconstruction quality, while the view planner can produce drone acquisition trajectories that lead to better reconstruction results compared to the previous methods [Liu et al. 2021a; Smith et al. 2018; Zhang et al. 2021; Zhou et al. 2020].

## 4 RECONSTRUCTABILITY

Reconstructability is an essential measure of how well a set of viewpoints reconstructs the target scene. Both our proposed and the existing planners [Smith et al. 2018; Zhang et al. 2021; Zhou et al. 2020] rely on it to optimize the view planning. Yet, unlike existing works, we define the reconstructability measure by formulating a learning approach and considering the ultimate goal of reconstructability, which is to enhance the quality of the final scene reconstruction. So, we adopt the reconstruction error metrics, accuracy and completeness, from Smith et al. [2018], and define the reconstructability term in our framework to be inversely proportional to the reconstruction error (Sec. 4.3). Hence, when we train our framework, minimizing the training loss would then drive our framework to learn to predict high (low) reconstructability for scene regions with low (high) reconstruction error. As a result, we can employ our framework to predict learned reconstructability in view planners to better estimate the final reconstruction quality.

The key to formulating the learning approach is to find the relation between the viewpoints and the scene geometry. Given $N$ viewpoints $\{v_i\}_{i=1}^N$ and sample point $p_j$ on proxy geometry, we want to learn function $G_s : (\mathbb{R}^{6 \times N}, \mathbb{R}^6) \rightarrow \mathbb{R}^1$ that predicts

$$R_j = G_s(\{v_i\}, p_j), \tag{1}$$

where each view $v_i$ consists of a position and an orientation; each sample point $p_j$ consists of a position and a normal vector; and $R_j$
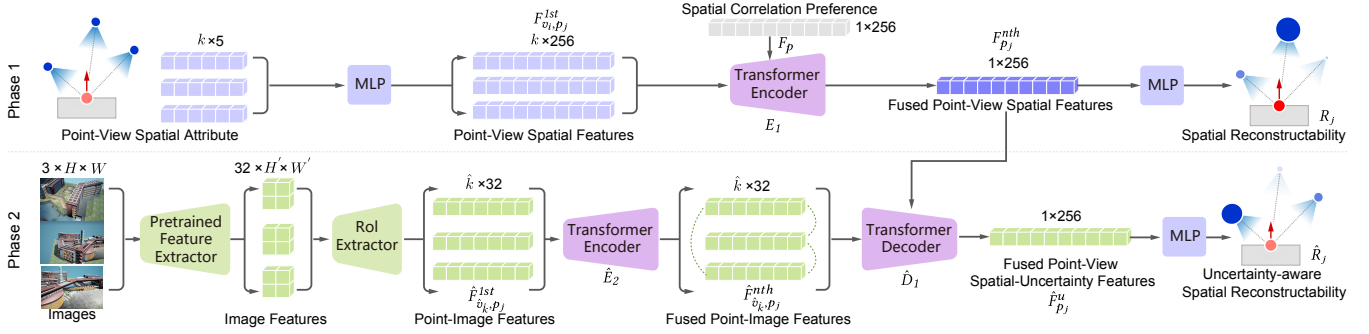
Fig. 3. The input to our network is a set of point-view spatial attributes for each sample point, and optionally images from a pre-flight. Based on the geometric characteristics between the sample points and viewpoints, our network automatically extracts the contribution of each viewpoint and predicts the *spatial reconstructability* of each sample point. If images from the pre-flight are available, our network can further extract the uncertainty of the proxy geometry near each sample point and predict the *uncertainty-aware reconstructability* to enhance the subsequent view planning process.

is the spatial reconstructability of point $p_j$ that measures how well $p_j$ can be reconstructed by the viewpoints $\{v_i\}$.

Also, we want to consider the uncertainty of the given scene geometry when predicting the reconstructability, *if* some images of the target scene are given. Compared to the spatial relations function $G_s$ we modeled above, images provide rich texture information, which can help improve both the *reconstructability* prediction and the subsequent path planning. Specifically, given $L$ existing RGB images $\{\hat{I}_l\}_{l=1}^{L}$ and their poses $\{\hat{v}_l\}_{l=1}^{L}$, we want to learn another function $\hat{G}_s : (\mathbb{R}^{L \times 3 \times H \times W}, \mathbb{R}^{6 \times L}, \mathbb{R}^{6 \times N}, \mathbb{R}^6) \rightarrow \mathbb{R}^1$ that predicts
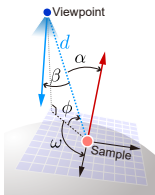
$$\hat{R}_j = \hat{G}_s(\{\hat{I}_l\}, \{\hat{v}_l\}, \{v_i\}, p_j), \quad (2)$$

where $H$ and $W$ are the size of images and $\hat{R}_j$ denotes the refined *reconstructability*, considering the potentially inaccurate geometry nearby. We show our network structure in Fig. 3. In the following section, we present the geometric representation of $R_s$ in Sec. 4.1 and how to compute the refined measurement $\hat{R}_s$ in Sec. 4.2. Also, we give the training details and a more specific calculation for our reconstructability measure in Sec. 4.3.

## 4.1 Spatial Reconstructability

When measuring the quality of a set of viewpoints, we want to formulate $G_s$ to extract the relation among the viewpoints, sample points, and *reconstructability*. Unlike the previous method [Smith et al. 2018], which assumes that $G_s$ follows a Gaussian distribution associated with some manually-defined parameters, we use a data-driven method to learn this function.

*Individual view-point feature extraction.* For each sample point $p_j$, we first locate all visible viewpoints at $p_j$ and calculate their spatial attributes with respect to $p_j$. As shown in the inset figure below, we consider the following five elements as the spatial features between point $p_j$ and view $v_i$ to predict the spatial reconstructability: the local spherical coordinates $(\omega, \phi, d)$ of the viewpoint with respect to the sample point, angle $\alpha$ between the sample point normal and the direction from the sample point to the viewpoint, and angle $\beta$ between the

viewpoint's viewing direction and the direction from the viewpoint to sample point. In practice, angles $\alpha$ and $\phi$ are complementary, so we only need to calculate one of them.

To better encode the individual influence of viewpoint $v_i$ on sample point $p_j$, we use an MLP (Multi Layer Perception) module to map the 5D spatial features to a 256D latent vector $F_{v_i,p_j}^{1st}$, *which* represents the influence of *individual* viewpoint on the associated local scene geometry during the reconstruction.

*Viewpoint feature fusion.* In the multi-view stereo (MVS) pipeline, viewpoints are highly coupled. The change in relative position and orientation of the viewpoints has a big impact on the final reconstruction [Furukawa et al. 2010; Smith et al. 2018]. So, we further extract a higher-order correlation between the viewpoints. Unlike the previous method [Smith et al. 2018], which exhaustively computes the feature of every view pair to approximate this correlation, we adopt a transformer encoder to learn the correlations among viewpoints and sample points, extract the contribution of each viewpoint, and then transform the individual feature $F_{v_i,p_j}^{1st}$ of $K$ visible viewpoints to a fused point-view spatial feature $F_{p_j}^{nth}$. Specifically, we train the transformer encoder, $E_1 : \mathbb{R}^{(K+1) \times 256} \rightarrow \mathbb{R}^{1 \times 256}$, with

$$F_{p_j}^{nth} = E_1(F_p, \{F_{v_i,p_j}^{1st}\}_{i=1}^{K}) \quad (3)$$

where the query, key, and value are all from $K$ individual features $\{F_{v_i,p_j}^{1st}\}_{i=1}^{K}$. Similar to DETR [Carion et al. 2020], we also use trainable parameter $F_p : \mathbb{R}^{1 \times 256}$ to represent the spatial correlation preference. We stack $F_p$ on the input features to form a $K + 1$ tensor at the beginning and extract the fused point-view spatial feature $F_{p_j}^{nth}$ from $F_p$ after the fusion.

*Spatial reconstructability.* Last, we use a standard MLP to learn to determine the spatial reconstructability $R_j$ from the fused point-view spatial feature $F_{p_j}^{nth}$.

## 4.2 Uncertainty-aware Spatial Reconstructability

The quality of the final reconstructed model is not only related to the geometric relation between the viewpoints and the scene geometry but also influenced by the surface appearance. More importantly,

the scene geometry we used in the above computation is usually reconstructed by a quick pre-flight pass [Hepp et al. 2018a; Smith et al. 2018; Zhou et al. 2020], which provides only coarse or even inaccurate geometry information. Hence, we train another function $\hat{G}_s$, which leverages the images captured from the pre-flight to further refine the reconstructability $R_j$ into the *uncertainty-aware spatial reconstructability* $\hat{R}_j$.

*Individual point-image feature extraction.* For each pre-captured image, we first extract their 32D latent code using a pre-trained convolutional neural network [Gu et al. 2020]. In order to obtain the image feature of each sample point $p_j$ on the scene geometry, we project point $p_j$ on the image and use a feature interpolation operator [He et al. 2017] to extract the individual feature $\hat{F}^{1st}_{\hat{v}_{\hat{k}}, p_j}$ : $\mathbb{R}^{32\times 1}$ of sample point $p_j$ in $\hat{k}th$ viewpoint $\hat{v}_{\hat{k}}$.

*Point-image feature fusion.* As the sample point can be visible at multiple viewpoints, we adopt another transformer encoder, $\hat{E}_2$ : $\mathbb{R}^{\hat{K}\times 32} \rightarrow \mathbb{R}^{\hat{K}\times 32}$, to correlate and fuse the $\hat{K}$ individual point-image feature $\hat{F}^{1st}_{\hat{v}_{\hat{k}}, p_j}$ to produce the fused image feature $\hat{F}^{nth}_{\hat{v}_{\hat{k}}, p_j}$ over all the visible viewpoints $\hat{V}_k$.

*Uncertainty-aware spatial reconstructability fusion.* Then, we can use the fused point-image feature of point $p_j$ to refine the spatial reconstructability $R_j$ that we predicted before. We adopt transformer decoder, $\hat{D}_1 : (\mathbb{R}^{\hat{K}\times 32}, \mathbb{R}^{1\times 256}) \rightarrow \mathbb{R}^{1\times 256}$, to extract the importance of the fused point-image feature $\hat{F}^{nth}_{\hat{v}_{\hat{k}}, p_j}$ to the spatial feature $F^{nth}_{p_j}$ and output the fused feature $\hat{F}^u_{p_j} : \mathbb{R}^{1\times 256}$ of point $p_j$ with

$$\hat{F}^u_{p_j} = \hat{D}_1(F^{nth}_{p_j}, \hat{F}^{nth}_{\hat{v}_{\hat{k}}, p_j}) \qquad (4)$$

where we use spatial feature $F^{nth}_{p_j}$ we predicted before as the query tensor to represent the pure spatial feature around this point and the fused image feature $\hat{F}^{nth}_{\hat{v}_{\hat{k}}, p_j}$ as the key and value to refine the feature, injecting semantics around point $p_j$ into the prediction.

*Uncertainty-aware spatial reconstructability.* Finally, we use an MLP module to predict the final uncertainty-aware spatial reconstructability $\hat{R}_j$ for point $p_j$ from the fused feature $\hat{F}^u_{p_j}$.

### 4.3 Training

While the reconstructability measures in existing works are well-defined and easy to obtain, taking them as target to train our framework will only drive our framework to predict reconstructability that mimics the existing measures. Our learning approach goes beyond existing works by considering the ultimate goal of reconstructability, i.e., to enhance the final reconstruction quality. Hence, we explicitly supervise the training of our network by simulating proxy-based acquisitions for various scenes and also the reconstruction process with planned viewpoints, such that during the training, we can define the target value of the network-predicted reconstructability based on the ground-truth reconstruction errors.

Specifically, we prepare our training data using UrbanScene3D [Lin et al. 2022], which consists of different scenes and different levels of the proxy [Smith et al. 2018; Zhou et al. 2020], as well as trajectories and associated reconstruction results from different path
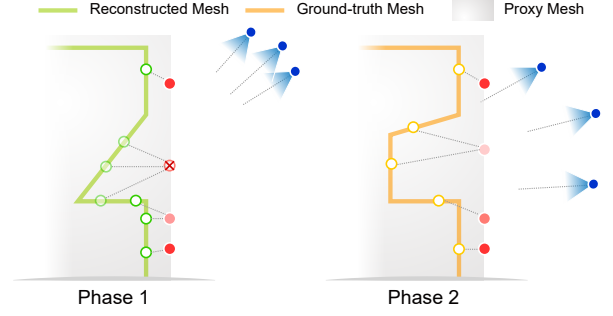
Fig. 4. We use different projection mechanisms for generating ground-truth labels (targets) in different training phases. Blue dots are viewpoints and red dots are sample points on proxy, where deeper red means the sample point has higher reconstructability. In *Phase 1*, we use the reconstruction accuracy between the reconstructed model and proxy to generate the ground-truth *reconstructability* for training our network, enabling it to learn the spatial correlation between the sample points and viewpoints. In *Phase 2*, we use the reconstruction completeness between the ground-truth surface and proxy surface to generate the target for our learned reconstructability, encouraging the network model to encode the uncertainty measurement in the reconstructability prediction.

planners [Smith et al. 2018; Zhang et al. 2021; Zhou et al. 2020]. As described in Sec. 3, our network input is a set of 5D relative information between the sample point and viewpoints, the visible pre-captured images, and their poses. We can easily obtain such input data from UrbanScene3D. However, training our framework is still hard, since it involves multiple data sources and configurations. Also, the image data is optional, as it may not be available during path planning. So, we split our training process into two phases and adopt different strategies to prepare training data in each phase. In particular, we use the *fine* and *inter* levels of proxy to train our framework in phase 1, since they have smaller difference from the ground-truth model. Then in phase 2, we use all four proxy levels to train the uncertainty-aware reconstructability predictor.

*Phase 1 training.* The goal of phase 1 is to model the pure geometric function $G_s$, which associates only with the relative information geometrically between the sample point on proxy and the viewpoints. From UrbanScene3D, we can obtain the *reconstruction accuracy* [Smith et al. 2018], which is defined on the sample points on the reconstructed model and measured by the shortest distance to the ground-truth surface. Then, we can project the reconstruction accuracy to the sample points on the proxy geometry, such that we can estimate the reconstruction accuracy on the proxy. In detail, for each sample point $p_j$ on proxy, we find a set of nearest sample points $\{p_q\}$ on the reconstructed mesh within distance threshold $\tau$. Then, by averaging accuracy $acc_q$ of each point $p_q$, we can implicitly encode the potential error near each sample point $p_j$ on proxy:

$$R^g_j = \frac{|\{p_q\}|}{\sum_{q\in\{p_q\}} acc_q}, \qquad (5)$$

where $R^g_j$ is used as the target of our reconstructability measure when training our framework. Also, we discard sample points on proxy with distances to the reconstructed model larger than $\tau$. This

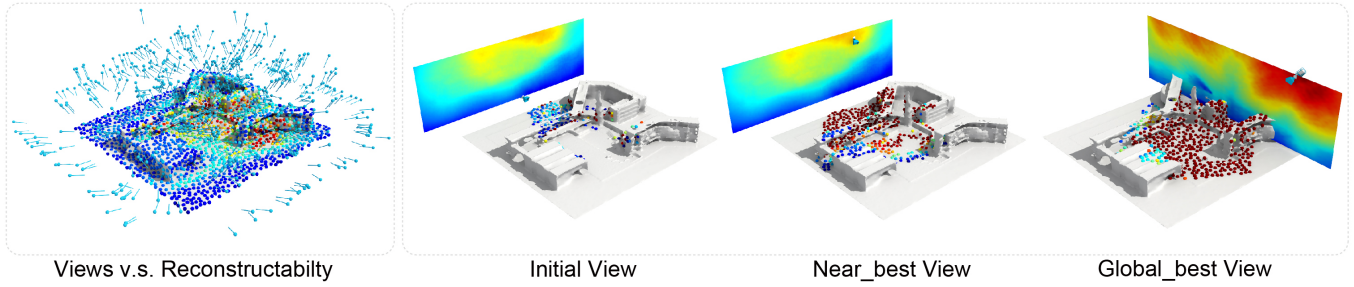| Views v.s. Reconstructabilty | Initial View | Near_best View | Global_best View |

Fig. 5. Here we show a common phenomenon: *local minima* may easily occur in the optimization process of previous view planners. The left figure shows all viewpoints and sample points with their associated reconstructability [Smith et al. 2018] in the current configuration. As shown in the second figure, we select a specific viewpoint as our analysis target. We change its position in the plane and calculate the reconstructability *increment* at that position. For each position, we randomly sample 128 directions and find the best one with the highest reconstructability to the sample points. To better visualize the calculated value nearby, we only show a section of the reconstructability field. The "best" place, or the local minimum, near this viewpoint can be found in the third figure, where the reconstructability of all sample points can be maximized. However, the right figure shows that the actual best place, or the global minimum, appears far away from the initial view location. In practice, it is very difficult to find the local minimum during the planning, since previous planners [Smith et al. 2018; Zhou et al. 2020] only leverage local information to optimize viewpoints. Our proposed view planning framework samples more viewpoints near poorly-reconstructed regions and reduces the number of viewpoints near well-reconstructed regions to help the planner escape from the local minima.

helps avoid taking proxy inaccuracy into our framework, as these sample points are likely located on inaccurate parts of the proxy.

*Phase 2 training.* Then, phase 2 further considers the pre-captured images [Smith et al. 2018], enabling us to predict uncertainty-aware reconstructability. Specifically, we define the training target directly using the *reconstruction completeness* [Smith et al. 2018], which is defined on the ground-truth model and measured by the shortest distance to the reconstructed model. Note that the proxy may differ significantly from the ground-truth model, due to insufficient data when preparing the proxy; see Fig. 4. Hence, phase 2 further encodes such uncertainty into the reconstructability prediction.

*Training details.* As Fig. 3 shows, we set all hidden dimensions in our network to 256, except for the dimension of the image features (i.e., 32), which is from the pre-trained model [Gu et al. 2020]. The input images are resized to $800 \times 600$ with color values normalized to $[0, 1]$. We use the standard RoI-Align operator [He et al. 2017] to extract the feature of each specific point in the given image. The MLP module we adopted consists of a linear layer and an ReLU activation. For data generation, we use $\tau = 20cm$ in all our experiments. Also, we use an *L1* loss to train the network. The whole training process takes about 10 hours on an RTX 3090Ti GPU.

## 5 VIEW PLANNER

In this section, we first discuss the integration of our reconstructability predictor with existing view planners in Sec. 5.1. Due to the non-linearity of the optimization process, existing planners may easily be stuck at local minima. So, in Sec. 5.2, we further present a new view planning framework to overcome this problem.

### 5.1 Integration with Existing Planners

Our reconstructability measure can be easily integrated into existing planners [Smith et al. 2018; Zhou et al. 2020]. Smith et al. [2018] can use our reconstructability to find an optimal viewpoint set that maximizes the reconstructability of sample points. We can readily

replace the reconstructability calculation using our method. More specifically, we can use our predicted reconstructability to identify if a new viewpoint configuration is better than the original one when executing the downhill simplex method in their method. So, we can use the same view planner to minimize the same objective.

Also, we can integrate our method into the view planner in [Zhou et al. 2020]. We use our predicted reconstructability to calculate the view redundancy of a given view configuration and perform the subsequent min-max view reduction. In each iteration, we follow [Zhou et al. 2020] to select the viewpoint with the highest redundancy and remove it temporally. Then, we can use our reconstructability predictor to test if any sample point receives a reconstructability lower than the threshold. If the test passes all associated sample points, we can remove the viewpoint permanently.

During the view planning, the objective function usually involves the relative position between the sample points and viewpoints, as well as between different viewpoints. Thus, the objective function is highly non-linear. Since the existing planners mostly optimize the view configuration in an iterative manner, they may easily fall into local minima. We further show this phenomenon in Fig. 5.

Previous planners often fail to escape from local minima. Smith et al. [2018] only find view candidates near the current position, preventing it from moving redundant views near poorly-reconstructed regions. As for Zhou et al. [2020], local minimum occurs when the view initialization is not perfect. However, it is hard to obtain a perfect view initialization for scenes with complex structures and occlusions. Also, having more initialized viewpoints will increase the computational burden on the visibility test and reconstructability calculation, as the complexity of reconstructability calculation is $O(|V|^2)$, where $|V|$ is the number of viewpoints.

Yet, we find that these two methods can complement each other. The local view adjustment from Smith et al. [2018] can help increase the reconstructability even with a poor view initialization from Zhou et al. [2020] by finding a better position and orientation near each viewpoint. On the other hand, the view initialization and elimination
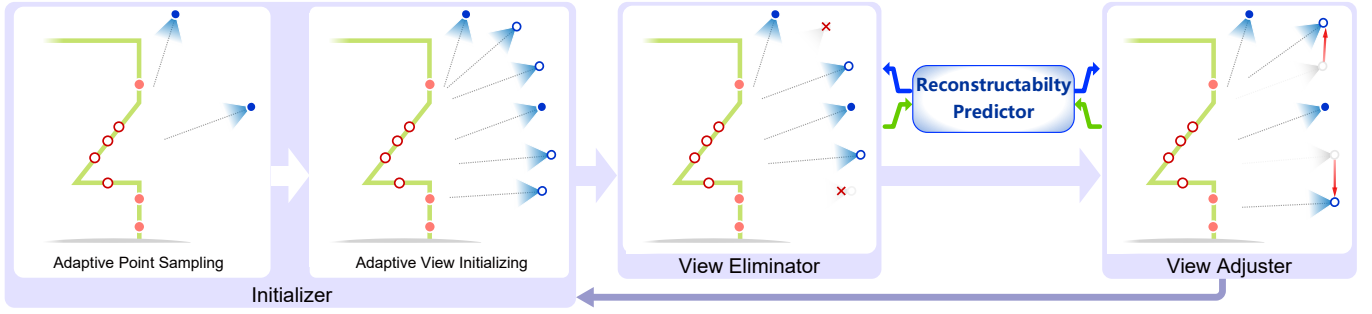
Fig. 6. Our proposed view-planning framework. Different from the existing ones, which are based on either view elimination [Zhou et al. 2020] or view adjustment [Smith et al. 2018], our *initializer* iteratively finds sample points with low reconstructability and try to allocate more views around these regions. The whole planner runs in an iterative way, helping it to escape from local minima during the optimization.
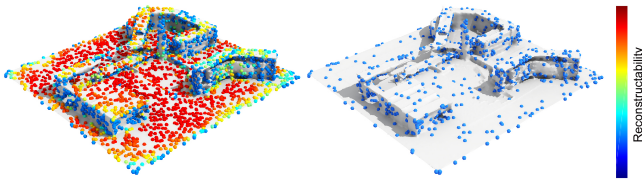


Fig. 7. Our sampling mechanism that generates the optimized target at each iteration. The left figure shows the current reconstructability of each sample point in the scene. The probability of each point to be sampled in the current iteration is calculated according to its reconstructability. The right figure shows the sampling result. These points are used to initialize, eliminate, and adjust viewpoints in the current iteration. This strategy helps the planner focus on regions that are hard to be reconstructed.

from Zhou et al. [2020] can help the local view adjustment from Smith et al. [2018] escape from local minima. Particularly, it helps avoid redundant views near well-reconstructed regions and allocate more views around regions with poor reconstructability.

## 5.2 New Planner

To this end, we develop a new view planning framework, which iteratively initializes, eliminates, and adjusts viewpoints to obtain a view configuration with maximum reconstructability. Our view planner consists of an initializer, a view eliminator, and a view adjuster, as shown in Fig. 6. Compared with the previous planners, our planner optimizes viewpoints in an iterative manner and can better escape from local minima during the optimization through adaptive point sampling and view initializing.

*Adaptive points sampling.* Zhou et al. [2020] first collect a dense set of viewpoints as the initialization and assume a perfect initialization. So, the planner only needs to reduce the number of viewpoints. However, it is hard to obtain a perfect viewpoint initialization, as the proxy geometry is usually coarse and inaccurate.

Specifically, we select sample points in an adaptive manner. The probability $Prob_{p_j}$ of selecting sample point $p_j$ is calculated based
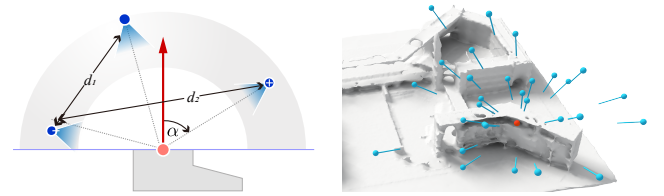


Fig. 8. We initialize viewpoints according to the normal vector of each sample point and the existing visible viewpoints toward this point. Specifically, we want the initialized viewpoint closer to the normal direction of the sample point (smaller $\alpha$), while being far away from the existing viewpoints (larger $d_1, d_2$). The left shows an example of our view initialization.

on its current reconstructability $R_j$ (or $\hat{R}_j$, if images are provided):

$$Prob_{p_j} = \frac{\sum_{q \in P_n} \frac{1}{R_q} e^{-d_q}}{|P_n|}, \qquad (6)$$

where $P_n$ is the set of sample points nearest to $p_j$ and $d_q$ is the distance from nearest sample point $q$ to point $p_j$. By this distance-weighted average, we can find regions that are not well-reconstructed and sample more points in them; see Fig. 7 for an illustration.

Compared with previous methods, which use uniform sample points as the optimization target, our proposed adaptive sampling can enable us to obtain better resolutions for regions that are previously hard to be reconstructed.

In each iteration, we sample $N$ points on the proxy surface as the optimization target according to the above probability. Note that the following view initialization elimination and adjustment will only be performed on the selected sample points.

*Adaptive view initialization.* For each sample point on the proxy surface, we create a set of viewpoints as a local initialization around the sample point. Zhou et al. [2020] directly extend the normal vector of each sample point to a specific view distance and place a viewpoint towards the sample point. However, complex geometric structures and occlusions on the proxy geometry will simply break this initialization, as shown in Fig. 8. Also, the viewing direction of the initialized viewpoint in Zhou et al. [2020] always points to the associated sample point. Such a setting could be optimal for
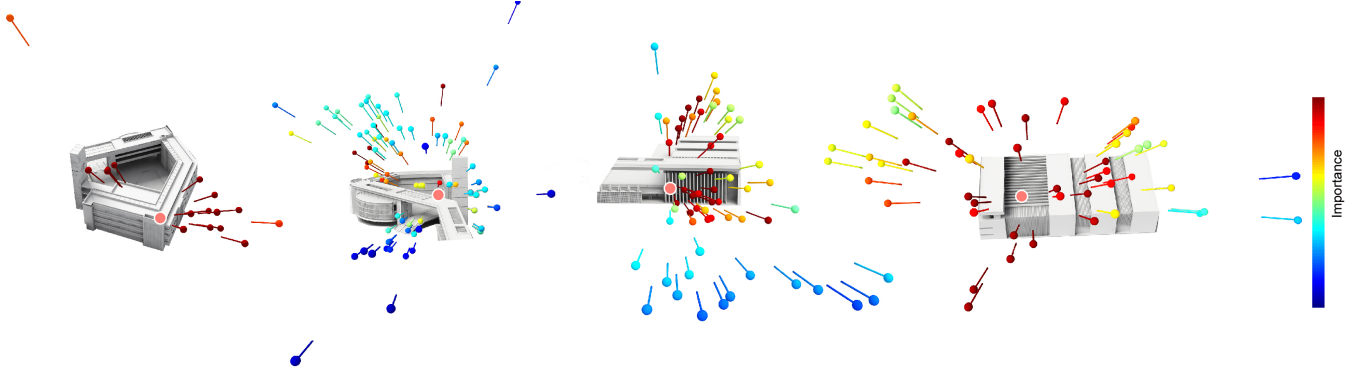
Fig. 9. The contribution of each viewpoint when predicting the reconstructability. Viewpoints with high contribution are marked red. Note that we only use the final reconstructability as the supervised signal during training. The contribution automatically extracted from the network indicates that viewpoints from far away have less impact on results, and viewpoints with appropriate baselines have higher weights when computing reconstructability.



Fig. 10. The training set for our reconstructability predictor. We use *Bridge, Castle, Town* from UrbanScene3D dataset [Lin et al. 2022], which contains 72 trajectories and the corresponding reconstruction results. For testing, we use a completely new scene: *School* to evaluate our predictor. Also, we use another dataset [Smith et al. 2018] and three real scenes to test the integration of our reconstructability predictor with various path planners [Smith et al. 2018; Zhou et al. 2020].

observing the associated sample point, but there should be better choices by considering multiple adjacent sample points. Instead, we randomly sample $M_m$ viewpoints in a hemisphere around each sample point. We then filter the visible viewpoints according to the proxy and add the best $M_b$ viewpoints to our viewpoint set. The weight of each sample viewpoint $v_m$ is calculated by

$$Score_{v_m} = dot(v_m - x_j, n_j) * \min_{v_v \in V_v} dot(v_m - x_j, v_v - x_j), \quad (7)$$

where $x_j$, $n_j$ are the position and normal vector of sample point $p_j$; and $V_v$ is the existing set of visible viewpoints at point $p_j$. The calculation encourages the viewpoint to be closer to the normal vector of point $p_j$, while further away from the existing viewpoints.

*View elimination and adjustment.* Based on the initialized viewpoints, we use Zhou et al. [2020] to compute the redundancy of each viewpoint and remove the redundant ones. Similar to Smith et al. [2018], we also adjust the viewpoints to further increase the reconstructability of the sample points after the view elimination.

## 6   RESULTS AND EVALUATION

We start with an analysis of the proposed reconstructability predictor in Sec. 6.1, by reporting the correlation factor between the

reconstructability predicted by different methods and the final reconstruction quality. Next, we integrate our reconstructability predictor into several existing view planners to demonstrate improved final reconstruction quality in Sec. 6.2. This is followed by experimenting with and evaluating the new view planner we propose, in Sec. 6.3. Finally, Sec. 6.4 presents results from our full view planning and reconstruction pipeline on three real scenes.

### 6.1   Reconstructability

We compare our reconstructability predictor with the heuristic estimate from Smith et al. [2018] using the *Spearman's rank correlation coefficient*, which is a measurement to quantify correlations. Indeed, a quality reconstructability estimate or prediction should output values that best reflect the final reconstruction quality. Specifically, the highest *Spearman correlation factor* is obtained when the order of the predicted reconstructability is the same as the order of the reconstruction quality in a scene.

Our experiments have been conducted using the trajectories and corresponding reconstruction results from the UrbanScene3D dataset [Lin et al. 2022]. Specifically, we train our network on three scenes: *Castle, Village, Bridge*, as shown in Fig. 10, with testing done on a new scene, *School*. We report the Spearman correlation of *Inter, Fine* proxy for phase 1 and all four levels of proxy *Box, Coarse, Inter, Fine* for phase 2. The whole test set contains 24 trajectories and their corresponding reconstruction results. These trajectories share different characteristics, e.g., in terms of flight patterns and view density. We use the reconstructability calculated by Smith et al. [2018] and the number of visible views as baselines for comparison.

*Reconstructability predicted by Smith et al. [2018].* For each sample point, we extract the visible viewpoints and calculate the corresponding reconstructability. We use the default parameter $k1 = 32$, $k3 = 8$, $alpha1 = \pi/16$, $alpha3 = \pi/4$, as in their paper.

*Number of visible views.* Alternatively, with the intuition that sample points with more visible viewpoints tend to lead to higher

Table 1. Quantitative comparison between different reconstructability estimates on the test scene: *School*, *without* image inputs. Higher Spearman correlation factor indicates better prediction. *Visible number* denotes the correlation factor between the reconstruction quality and the number of visible viewpoints at each sample point. Compared with the two baselines, the reconstructability predicted by our method better matches the final reconstruction quality.

| Planner | Overlap | Proxy | Image (#) | Ours | Smith et al. | | Visible Number | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | Spearman | Spearman | Inc. | Spearman | Inc. |
| Smith et al. | 70 | Inter | 559 | **32.42%** | 17.83% | 81.83% | 22.83% | 42.01% |
| | | Fine | 559 | **43.35%** | 36.66% | 18.25% | 36.15% | 19.92% |
| Smith et al. | 90 | Inter | 559 | **14.79%** | 10.31% | 43.45% | 12.32% | 20.05% |
| | | Fine | 559 | **34.94%** | 30.86% | 13.22% | 26.68% | 30.96% |
| Zhou et al. | 70 | Inter | 342 | **25.49%** | 17.62% | 44.67% | 14.37% | 77.38% |
| | | Fine | 518 | **44.51%** | 41.48% | 7.30% | 37.84% | 17.63% |
| Zhou et al. | 90 | Inter | 595 | **23.02%** | 11.48% | 100.52% | 12.02% | 91.51% |
| | | Fine | 1243 | **38.79%** | 32.07% | 20.95% | 29.82% | 30.08% |
| Zhang et al. | 70 | Inter | 330 | **30.53%** | 27.83% | 9.70% | 27.23% | 12.12% |
| | | Fine | 518 | **36.35%** | 30.58% | 18.87% | 34.37% | 5.76% |
| Zhang et al. | 90 | Inter | 570 | **32.88%** | 24.63% | 33.50% | 31.46% | 4.51% |
| | | Fine | 1043 | **49.53%** | 44.57% | 11.13% | 43.39% | 14.15% |
| Average Inc. | | | | | | **33.62%** | | **30.51%** |

Table 2. Quantitative comparison between different reconstructability estimates on the *School* scene, *with* image inputs, where the image features were computed for images captured during drone pre-flight. Again, our predictor performs better than the baselines even with inaccurate proxies (*Coarse, Inter*).

| Planner | Overlap | Proxy | Image (#) | Ours | Smith et al. | | Visible Number | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | Spearman | Spearman | Inc. | Spearman | Inc. |
| Smith et al. | 70 | Box | 559 | **26.89%** | 7.30% | 268.36% | 7.69% | 249.67% |
| | | Coarse | 559 | **27.07%** | 9.16% | 195.52% | 8.89% | 204.50% |
| | | Inter | 559 | **36.02%** | 16.88% | 113.39% | 23.50% | 53.28% |
| | | Fine | 559 | **55.51%** | 53.12% | 4.50% | 52.19% | 6.36% |
| Smith et al. | 90 | Box | 559 | **15.69%** | 8.66% | 81.18% | 11.55% | 35.84% |
| | | Coarse | 559 | **7.65%** | 1.13% | 576.99% | 1.08% | 608.33% |
| | | Inter | 559 | **29.85%** | 11.62% | 156.88% | 20.27% | 47.26% |
| | | Fine | 559 | **45.12%** | 40.06% | 12.63% | 37.34% | 20.84% |
| Zhou et al. | 70 | Box | 416 | 23.23% | 21.54% | 7.85% | **24.00%** | -3.21% |
| | | Coarse | 330 | **25.50%** | 17.78% | 43.42% | 19.99% | 27.56% |
| | | Inter | 342 | **41.30%** | 19.38% | 113.11% | 18.32% | 125.44% |
| | | Fine | 518 | **56.35%** | 55.23% | 2.03% | 49.40% | 14.07% |
| Zhou et al. | 90 | Box | 614 | **16.58%** | 11.23% | 47.64% | 11.05% | 50.05% |
| | | Coarse | 570 | **22.25%** | -2.02% | 1201.49% | 4.89% | 355.01% |
| | | Inter | 595 | **37.12%** | 6.16% | 502.60% | 12.40% | 199.35% |
| | | Fine | 1243 | **47.01%** | 40.41% | 16.33% | 39.90% | 17.82% |
| Zhang et al. | 70 | Box | 518 | **29.79%** | 16.00% | 86.19% | 14.30% | 108.32% |
| | | Coarse | 330 | **23.57%** | 14.32% | 64.59% | 14.68% | 60.56% |
| | | Inter | 330 | **44.12%** | 33.76% | 30.69% | 35.25% | 25.16% |
| | | Fine | 518 | 57.04% | 52.36% | 8.94% | **57.88%** | -1.45% |
| Zhang et al. | 90 | Box | 614 | **23.78%** | 14.60% | 62.88% | 11.10% | 114.23% |
| | | Coarse | 570 | 9.92% | 6.90% | 43.77% | **12.13%** | -18.22% |
| | | Inter | 570 | **45.42%** | 33.03% | 37.51% | 37.82% | 20.10% |
| | | Fine | 1043 | **60.82%** | 59.55% | 2.13% | 60.45% | 0.61% |
| Average Inc. | | | | | | **153.36%** | | **96.73%** |

reconstruction quality, we also calculate the number of visible viewpoints for each sample point and compute its Spearman correlation factor with respect to the final reconstruction quality.

As shown in Table 1, the values predicted by our reconstructability predictor have a higher Spearman correlation factor than those

Table 3. Quantitative evaluation, on reconstruction quality using F-score, Precision, and Recall, of integrating our reconstructability predictor into two different view planners [Smith et al. 2018; Zhou et al. 2020], either using image inputs or not. The method being compared to, marked by "Smith", employed reconstructability estimates by the method in Smith et al. [2018] to guide the view planning for 3D scene reconstruction.

| Proxy | Recon. | Smith et al. Planner | | | | Zhou et al. Planner | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Image (#) | F-score↑ | Precision↑ | Recall↑ | Image (#) | F-score↑ | Precision↑ | Recall↑ |
| Box | Smith | 1,770 | 31.8491 | 44.8348 | 24.6963 | 806 | 25.6819 | 42.4220 | 18.4151 |
| | Ours (w/o) | 1,717 | 32.5645 | 45.6230 | 25.3178 | 699 | 27.6641 | 44.2545 | **20.1210** |
| | Ours | 1,610 | **34.5483** | **52.4389** | **25.7598** | 796 | **28.1596** | **48.5530** | 19.8303 |
| Inter | Smith | 714 | 34.3846 | 50.9298 | **25.9533** | 825 | 29.2754 | 46.7762 | 21.3045 |
| | Ours (w/o) | 747 | **34.6040** | **52.3987** | 25.8315 | 747 | 29.9230 | 47.4468 | 21.8522 |
| | Ours | 696 | 34.4410 | 52.1732 | 25.7047 | 696 | **32.4972** | **53.8454** | **23.2709** |

Table 4. Quantitative comparison between the different planners, on reconstruction quality measured using *accuracy*, as explained in Section 6.2.

| Proxy | Recon. | Zhou et al. Planner | | | | | Smith et al. Planner | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Image (#) | 70%↓ | 80%↓ | 90%↓ | 95%↓ | Image (#) | 70%↓ | 80%↓ | 90%↓ | 95%↓ |
| Box | Smith | 1,770 | 0.0256 | 0.0391 | 0.0695 | 0.1085 | 806 | 0.0280 | 0.0427 | 0.0769 | 0.1229 |
| | Ours (w/o) | 1,717 | 0.0246 | 0.0378 | 0.0675 | 0.1038 | 699 | 0.0253 | 0.0384 | 0.0680 | 0.1052 |
| | Ours | 1,610 | **0.0186** | **0.0289** | **0.0531** | **0.0844** | 796 | **0.0218** | **0.0346** | **0.0650** | **0.1044** |
| Inter | Smith | 714 | 0.0193 | 0.0298 | 0.0545 | 0.0860 | 825 | 0.0236 | 0.0368 | 0.0686 | 0.1110 |
| | Ours (w/o) | 747 | **0.0187** | **0.0291** | **0.0536** | **0.0848** | 747 | 0.0227 | 0.0351 | 0.0627 | 0.0987 |
| | Ours | 696 | 0.0189 | 0.0297 | 0.0543 | 0.0850 | 696 | **0.0185** | **0.0303** | **0.0573** | **0.0887** |

Table 5. Quantitative comparison between the different planners, on reconstruction quality measured using *completeness*, as explained in Section 6.2.

| Proxy | Recon. | Zhou et al. Planner | | | | | Smith et al. Planner | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Image (#) | 70%↓ | 80%↓ | 90%↓ | 95%↓ | Image (#) | 70%↓ | 80%↓ | 90%↓ | 95%↓ |
| Box | Smith | 1,770 | 0.2409 | 0.4607 | 0.9519 | 2.5048 | 806 | 0.3867 | 0.6710 | **1.1889** | **2.6265** |
| | Ours (w/o) | 1,717 | **0.2352** | **0.4203** | **0.8960** | **1.8019** | 699 | **0.3265** | **0.6027** | 1.1894 | 2.6428 |
| | Ours | 1,610 | 0.2576 | 0.4968 | 1.0011 | 2.4704 | 796 | 0.3928 | 0.6794 | 1.1870 | 2.6571 |
| Inter | Smith | 714 | **0.2471** | **0.4789** | **0.9611** | 2.3776 | 825 | 0.3083 | 0.5661 | 1.1051 | 2.6055 |
| | Ours (w/o) | 747 | 0.2549 | 0.4929 | 0.9898 | 2.4303 | 747 | **0.2856** | **0.5356** | 1.1274 | 2.6525 |
| | Ours | 696 | 0.2594 | 0.5013 | 0.9917 | **2.3767** | 696 | 0.2969 | 0.5475 | **1.0382** | **2.5033** |

Table 6. Comparing our new view planner with those from [Smith et al. 2018; Zhou et al. 2020] using F-score, Precision, and Recall. All planners employ our reconstruction predictor with the help of image features.

| Planner | Box Proxy | | | | Inter Proxy | | | |
|---|---|---|---|---|---|---|---|---|
| | Image (#) | F-score↑ | Precision↑ | Recall↑ | Image (#) | F-score↑ | Precision↑ | Recall↑ |
| Smith et al. | 735 | 28.1596 | 48.5530 | 19.8303 | 714 | 32.4972 | 53.8454 | 23.2709 |
| Zhou et al. | 1610 | 34.5483 | 52.4389 | **25.7598** | 696 | 34.4410 | 52.1732 | **25.7047** |
| Ours | 764 | **34.9627** | **58.0929** | 25.0062 | 754 | **35.0951** | **59.8070** | 24.8339 |

obtained by the other baselines. Furthermore, the margin of gains by our method is even greater when the proxy is inaccurate, as can be seen from Table 2. Note that the *Inter* level of proxy refers to the proxy generated by the rapid pre-flight using *Oblique Photography*. *Coarse* and *Box* proxies refer to those extracted by satellite images. Both of these image acquisition and scene reconstruction workflows have been common practices in the industry.

In Fig. 9, we show the weight of each viewpoint when predicting reconstructability. While the contribution of each viewpoint is automatically extracted by our network without any supervision, there are still observable patterns, e.g., larger viewing distance tends to result in lower contribution; viewpoints with an appropriate baseline and scale difference often produce higher contributions.

Table 7. Comparing our new view planner with those from [Smith et al. 2018; Zhou et al. 2020] using *accuracy*. All planners employ our reconstruction predictor and image features.

| Planner | Box Proxy | | | | | Inter Proxy | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Image (#) | 70%↓ | 80%↓ | 90%↓ | 95%↓ | Image (#) | 70%↓ | 80%↓ | 90%↓ | 95%↓ |
| Smith et al. | 735 | 0.0218 | 0.0346 | 0.0649 | 0.1044 | 714 | 0.0185 | 0.0303 | 0.0573 | 0.0888 |
| Zhou et al. | 1610 | 0.0186 | 0.0289 | 0.0531 | 0.0844 | 696 | 0.0189 | 0.0297 | 0.0544 | 0.0850 |
| Ours | 764 | **0.0153** | **0.0240** | **0.0444** | **0.0721** | 754 | **0.0147** | **0.0235** | **0.0448** | **0.0740** |

Table 8. Comparing our new view planner with those from [Smith et al. 2018; Zhou et al. 2020] using *completeness*. All planners employ our reconstruction predictor and image features.

| Planner | Box Proxy | | | | | Inter Proxy | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Image (#) | 70%↓ | 80%↓ | 90%↓ | 95%↓ | Image (#) | 70%↓ | 80%↓ | 90%↓ | 95%↓ |
| Smith et al. | 735 | 0.3928 | 0.6794 | 1.1870 | 2.6571 | 714 | 0.2969 | 0.5475 | 1.0382 | 2.5033 |
| Zhou et al. | 1610 | **0.2576** | **0.4968** | **1.0011** | **2.4704** | 696 | **0.2594** | **0.5013** | **0.9917** | **2.3767** |
| Ours | 764 | 0.3139 | 0.5971 | 1.1862 | 2.6151 | 754 | 0.3188 | 0.5959 | 1.1362 | 2.4887 |

Table 9. Quantitative comparisons on reconstruction quality, in terms of *accuracy* and *completeness*, between different methods using another reconstruction benchmark in [Smith et al. 2018]. Note that both test scenes are unseen by our data-driven reconstructability predictor. Our view planner is clearly the best performing one. In the table, we highlight the best performing numbers in each column in bold, and the second best performing numbers in italic.

| Method | | Acc.↓ 90% | Acc.↓ 95% | Comp.↑ 0.02m | Comp.↑ 0.05m | Comp.↑ 0.075m | | Acc.↓ 90% | Acc.↓ 95% | Comp.↑ 0.02m | Comp.↑ 0.05m | Comp.↑ 0.075m |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Smith et al. | | 0.053 | 0.792 | 36.010 | 44.740 | 49.470 | | **0.028** | *0.051* | **32.040** | *37.740* | *40.620* |
| Zhou et al. | NY-1 | *0.030* | 0.342 | **38.190** | *45.220* | *49.780* | UK-1 | *0.030* | 0.054 | 30.750 | 35.960 | 38.770 |
| Liu et al. | | **0.029** | **0.107** | N/A | 44.72 | 49.33 | | **0.028** | 0.052 | N/A | 36.71 | 39.58 |
| Ours | | 0.039 | *0.192* | *38.077* | **46.046** | **50.523** | | 0.031 | **0.050** | *31.755* | **37.843** | **40.795** |

Table 10. Comparing our new path planer with state-of-the-art alternatives [Smith et al. 2018; Zhou et al. 2020] on two real scenes. Each real scene has a high precision LiDAR capture as the ground truth model. We report F-score, accuracy, and completeness of the final reconstruction results.

| Scene | Method | Image (#) | F-score↑ | Precision↑ | Recall↑ | Acc.↓ 80% | Acc.↓ 90% | Acc.↓ 95% | Comp.↓ 80% | Comp.↓ 90% | Comp.↓ 95% |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Smith et al. | 678 | 9.9549 | 7.1530 | 16.3659 | 0.2106 | 0.4385 | 0.7271 | 0.3011 | 0.8440 | 2.0124 |
| Polytech | Zhou et al. | 1364 | 10.3966 | 7.3969 | 17.4888 | 0.1872 | 0.4229 | 0.7108 | 0.2775 | 0.8348 | 1.8146 |
| | Ours | 1141 | **16.5603** | **11.8909** | **27.2685** | **0.1435** | **0.3331** | **0.6368** | **0.2368** | **0.7448** | **1.7062** |
| | Smith et al. | 2900 | 8.4057 | 5.7568 | 15.5701 | 0.2338 | 0.4066 | 0.6450 | 0.5787 | 1.5299 | 2.5265 |
| ArtSci | Zhou et al. | 3286 | 11.5342 | 8.1013 | 20.0158 | **0.1965** | **0.3517** | **0.5494** | 0.5650 | 1.5680 | 2.5586 |
| | Ours | 2600 | **13.7585** | **9.1734** | **27.5069** | 0.2167 | 0.3984 | 0.6440 | **0.4632** | **1.3959** | **2.3358** |

## 6.2 Integration with Existing Planners

We simulate a complete scene reconstruction pipeline to evaluate the integration of our reconstructability predictor with existing path planners [Smith et al. 2018; Zhou et al. 2020], after slight modifications as discussed in Sec. 5. These newly adapted planners are compared to their counterparts with the reconstructability estimated by the approach in Smith et al. [2018]. Similar to Sec. 6.1, we test reconstructability prediction and path planning on the *School* scene. Also, we use two proxy levels, *Coarse* and *Inter*, to plan trajectories as they represent two common ways to obtain scene proxies in

practice: a quick reconstruction from a rapid flight and 2.5D box extraction from satellite images, respectively.

For evaluation, two different metrics, accuracy and completeness, are employed to compare quality of the reconstructed meshes, as in Smith et al. [2018]. While accuracy measures how close a reconstructed mesh is to the ground truth mesh, completeness reveals how well the ground truth mesh is "covered" by the reconstruction. In other words, accuracy accounts from distances from the reconstructed mesh to the ground truth mesh, while completeness is measured based on distances from the ground truth mesh to the reconstructed mesh. Specifically, for each sample point on the

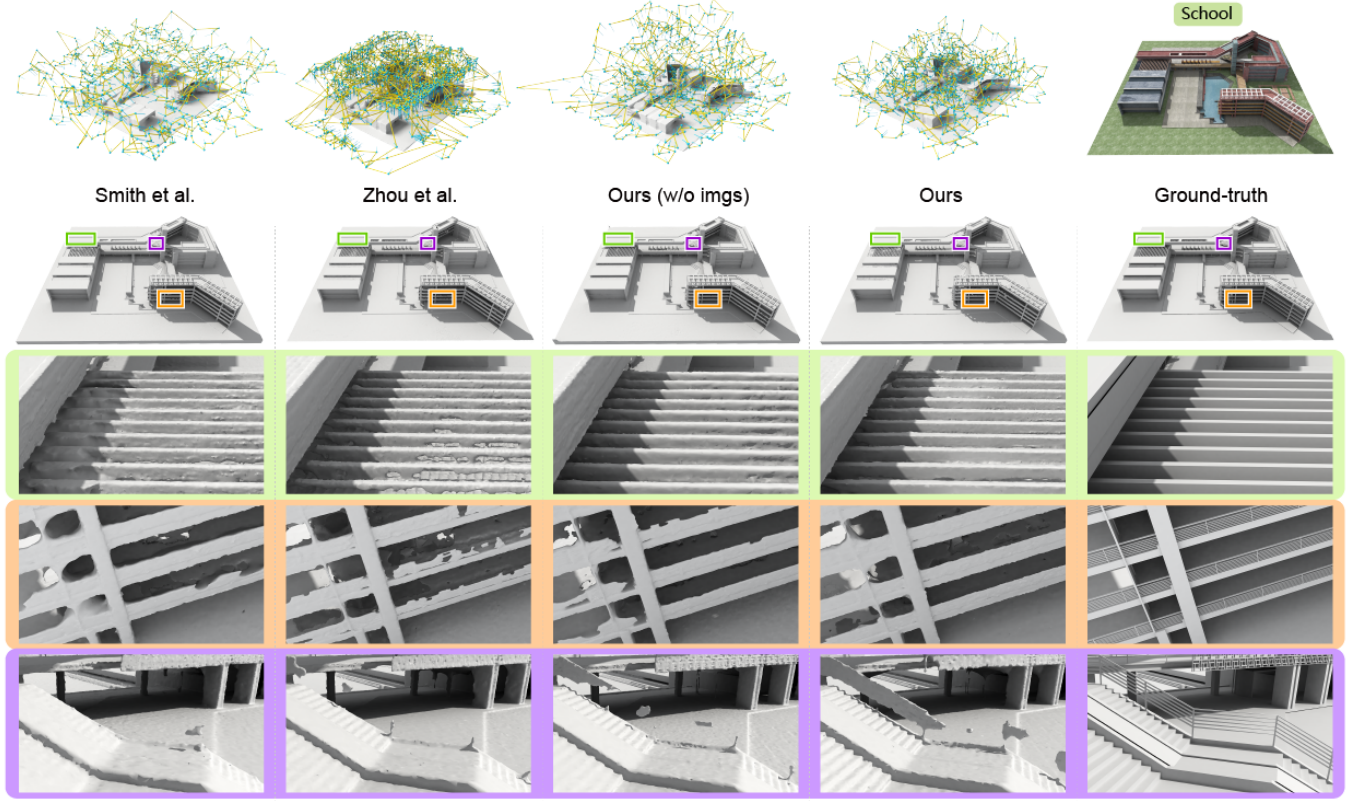| Smith et al. | Zhou et al. | Ours (w/o imgs) | Ours | Ground-truth |

Fig. 11. A visual comparison of reconstruction results produced by different methods: Smith et al. [2018], Zhou et al. [2020], and ours (with vs. without using image features). With the ground truth meshes as references, we can see that the reconstruction results obtained by our method can more faithfully recover geometric details and sharp features.

surface of a reconstructed mesh, we find the closest point on the ground truth mesh and compute their distance. We then sort these distances for all the sample points and report an accuracy number $y$, with respect to a given percentage $x$%, if exactly $x$% of the distances are less than $y$. Clearly, a lower accuracy number would corroborate with higher reconstruction quality. For completeness, we simply switch the roles of the reconstructed mesh and the ground truth mesh and perform the same computations. Additionally, we report F-scores [Knapitsch et al. 2017] as a summary metric, where a threshold of $10cm$ was applied to extract inliers and outliers.

Quantitative comparisons on reconstruction quality between the different planners are shown in Tables 3, 4, and 5. Specifically, Table 3 shows that, in terms of F-score, Precision, and Recall, the planners employing our reconstructability predictor, either with image inputs or without, generally outperform the baseline planners guided by reconstructability estimates from Smith et al. [2018]. Comparisons in terms of accuracy and completeness measures also exhibit a similar trend, as shown in Tables 4 and 5.

### 6.3 Integration with the New Planner

To evaluate the performance of our new view planner, we compare it to the planners from [Smith et al. 2018; Zhou et al. 2020], all using our reconstructability predictor with the help of image features. We

follow the same evaluation strategy as in Sec. 6.2, also testing on the *School* scene. The quantitative results in Tables 6, 7, and 8 show that our planner generally outperforms the alternatives, attaining higher precision while maintaining a comparable recall on the final reconstruction. Fig. 11 provides a visual comparison between reconstruction results produced by the different methods.

We also evaluate our method using the scene reconstruction benchmark proposed by Smith et al. [2018]. Following Zhou et al. [2020], we choose *NY-1* and *UK-1* as the test scenes, since *Bridge-1* has been used for training. Note that both of these test scenes are new to our reconstructability predictor. As the results in Table 9 show, while our new planner does not come on top in every reported measure, it is clearly the best performing one among the four methods compared. Note that to strictly follow the measures employed by the reconstruction benchmark, the reported completeness measure in Table 9 is different from those shown in the other tables: the roles of $x$ and $y$ in the prior completeness definition are switched so that larger numbers reflect higher reconstruction quality.

### 6.4 Test on Real Scenes

To further demonstrate the performance of our method, we show qualitative reconstruction results obtained by our new view planner on several challenging real 3D scenes. We chose three scenes
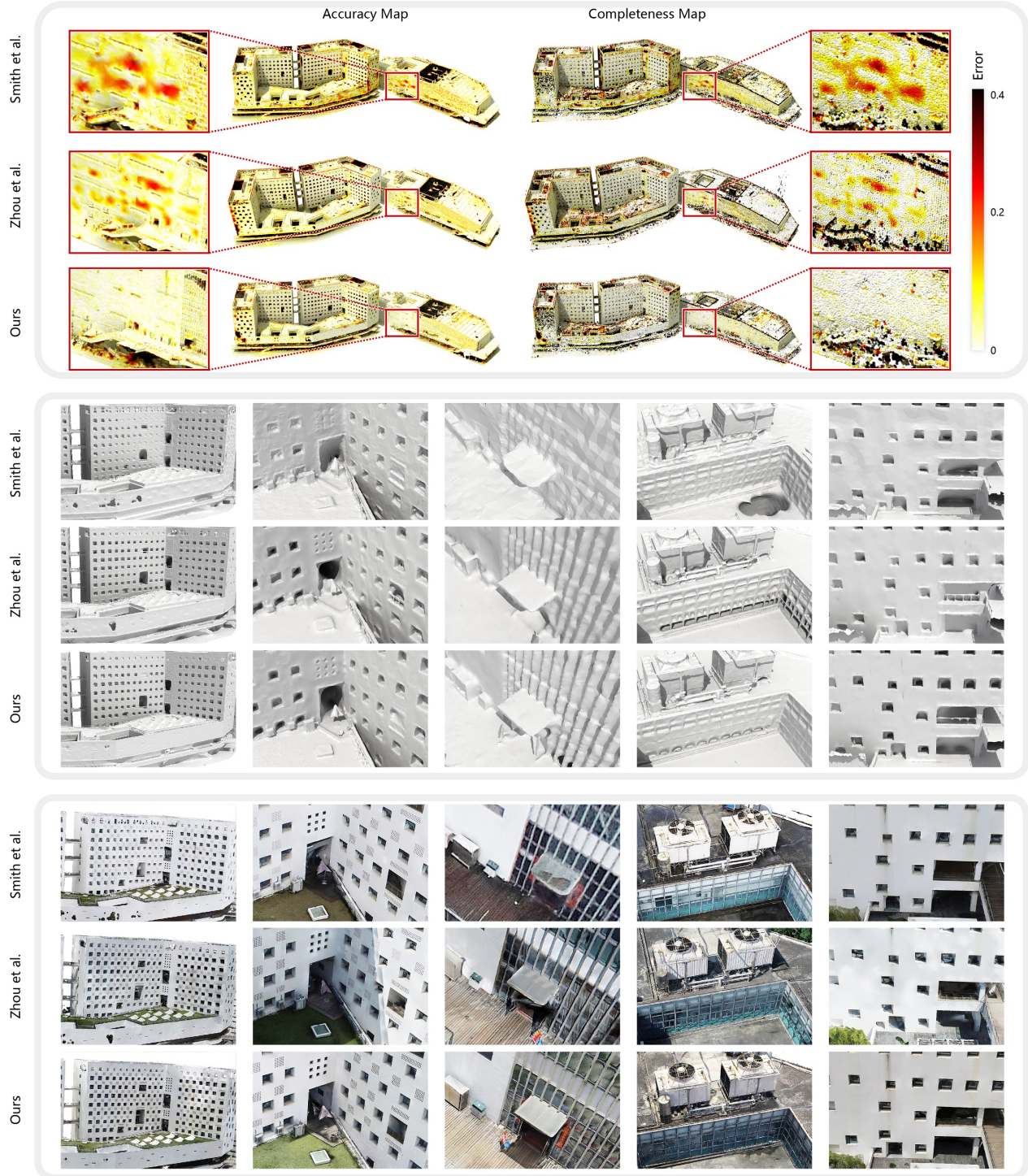
Fig. 12. Qualitative results on the *Polytech* scene, compared with previous planning and reconstruction methods [Smith et al. 2018; Zhou et al. 2020]. *Top:* Global and Local error map on the scene. The accuracy map is collected by calculating the shortest distance from the reconstructed mesh to the GT LiDAR points, while the completeness map is collected by collecting the shortest distance from the ground truth LiDAR points to the reconstructed mesh. *Middle:* Local geometry details on the reconstructed mesh. *Bottom:* Local textured details on the reconstructed mesh.

Fig. 13. Visualization of the reconstruction results on the *ArtSci* scene compared with previous methods [Smith et al. 2018; Zhou et al. 2020]. Unlike *Polytech*, *ArtSci* contains two irregular buildings with more complex geometries, leading to increased difficulty towards trajectory planning. We also show the LiDAR points that were collected by the high-resolution LiDAR scanner, which are used as the ground truth model for evaluation.
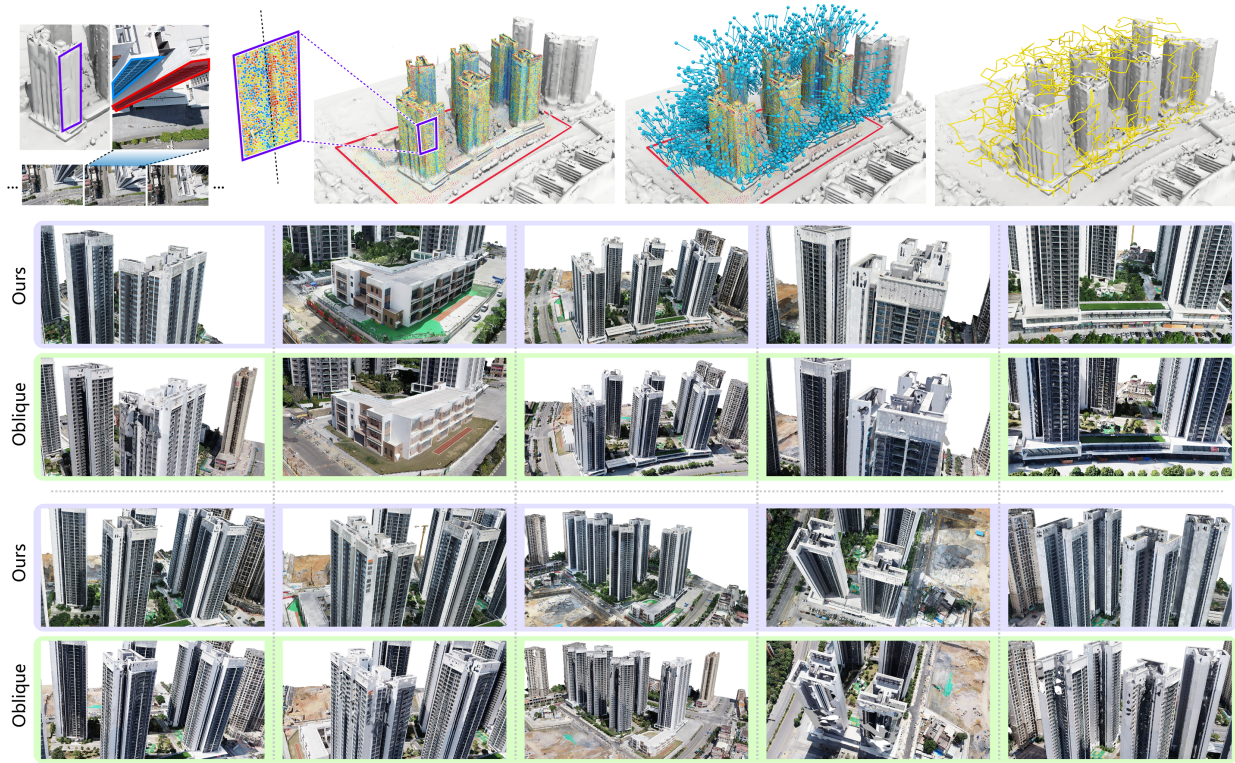
Fig. 14. Qualitative results on the *Apartment* scene, compared with *Oblique Photography*, a widely-used path planning method in the industry. *Apartment* contains six high-rise buildings (higher than 80m), making it difficult to obtain a quality reconstruction using *Oblique Photography*. Moreover, our method can extract and incorporate scene uncertainties from input images into reconstructability prediction. As shown in the top left of the image, sample points located over inaccurate regions of the proxy can be implicitly identified by our reconstructability predictor. They have lower reconstructability, which would encourage the view planner to produce viewpoints with large baselines in order to obtain a comprehensive observation over these inaccurate regions.

possessing different scales and characteristics. Specifically, *Polytech* contains two buildings with weak texture, repeated patterns, and large height differences, while *ArtSci* contains two irregular buildings with complex geometries and occlusions. Finally, *Apartment* contains six high-rise buildings covering the largest area.

We compare our reconstruction results for *Polytech* and *ArtSci* with the high resolution LiDAR point clouds from UrbanScene3D [Lin et al. 2022] as ground truth. As shown in Table 10, the scenes reconstructed by our method are generally of superior quality. Note that our model was trained using synthetic data, without any fine-tuning on real scenes. Fig. 12 and Fig. 13 show visual comparisons, demonstrating that our reconstruction results are more accurate, especially over regions with complex geometries and occlusions.

Since we do not have a ground-truth mesh for *Apartment*, we simply compare our reconstruction results qualitatively with those obtained by a widely-adopted industrial algorithm: *Oblique Photography*. As demonstrated in Fig. 14, our method achieves better reconstruction quality visually, especially over regions near the ground, where images from *Oblique Photography* can hardly be observed. Moreover, we show that our reconstructability predictor can extract the potential inaccurate geometry from the images and propagate the uncertainty to the reconstructability predictor.

## 7 CONCLUSION, LIMITATION, AND FUTURE WORK

Our work shows that reconstructability, in the context of drone path planing for urban scene acquisition, is a learnable measure. Specifically, we define reconstructability, i.e., the expected scene reconstruction quality, as a function of proxy geometry and a set of viewpoints, and introduce the first data-driven predictor trained on synthetic data from the new UrbanScene3D dataset.

While our acquisition problem falls under the general realm of "learning to reconstruct", it differs from most reconstructive tasks including all recent works on neural fields [Xie et al. 2022]. In these latter works, and under the typical setting for 3D reconstruction, the input at test time is a direct, albeit under-, representation of the target 3D scene, e.g., a shape code for IM-Net [Chen and Zhang 2019] or multi-view images for NeRF [Mildenhall et al. 2020]. In our problem setting, such inputs are not given; we must predict a view plan to acquire these inputs first, on-the-fly during test time, and then reconstruct the scene. As a result, the design of our learning framework has to handle at least two gaps: the domain gap between synthetic and real data, and the accuracy gap between the proxy and true geometry of the reconstructed scene.

Extensive experiments have been conducted to demonstrate that our learned reconstructability better correlates with the true reconstruction quality than existing heuristic estimates. Combined with an iterative view optimization scheme, our predictor can be integrated into both previous and our new path planners, leading to consistent improvements on reconstruction quality. Qualitative and quantitative results are presented for both synthetic and real scenes to demonstrate generalizability of our learning framework.

As a first attempt at a learning framework for reconstructability, our work still has several limitations. For example, while our feature learning is geometry- and uncertainty-aware, it does not explicitly account for material properties of the acquired scene. More critically, since our predictor operates on point-view and point-image features that are both defined on the scene proxy, the scale and variability of the gaps between the proxy and the true scene geometry can impact the quality of the learned model. Currently, we relate points from the proxy and the true surface via closest distances, which is a simple heuristic but not a reliable correspondence.

Furthermore, both the view eliminator from Zhou et al. [2020] and the view adjuster from Smith et al. [2018] need a scoring function to transform the calculated reconstructability from the sample points to the viewpoint, as the view planner must decide whether a viewpoint is redundant or a new choice is better. Our current scoring function is quite heuristic, e.g., Zhou et al. [2020] use the smallest reconstructability of the sample point that the viewpoint can observe as the score. Compared with the existing reconstructability definitions, one may also consider how to directly define the reconstructability on the optimized viewpoints.

Besides addressing the above limitations, it would also be interesting to explore the correlation between viewpoints in the MVS pipeline. Learning such correlations is useful for view planning, pose estimation, as well as sparse/dense urban reconstruction. Finally, we are interested in applying and adapting our learning framework for robot-assisted 3D object or indoor scene acquisition.

## ACKNOWLEDGMENTS

## REFERENCES

Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-End Object Detection with Transformers. In *Proc. Euro. Conf. on Computer Vision*. 213–229.

Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. 2015. ShapeNet: An Information-Rich 3D Model Repository. *arXiv preprint arXiv:1512.03012* (2015).

Zhiqin Chen and Hao Zhang. 2019. Learning Implicit Fields for Generative Shape Modeling. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*. 5939–5948.

Yikang Ding, Wentao Yuan, Qingtian Zhu, Haotian Zhang, Xiangyue Liu, Yuanjiang Wang, and Xiao Liu. 2022. TransMVSNet: Global Context-Aware Multi-View Stereo Network With Transformers. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*. 8585–8594.

Yasutaka Furukawa, Brian Curless, Steven M Seitz, and Richard Szeliski. 2010. Towards Internet-scale Multi-view Stereo. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*. 1434–1441.

Kyle Genova, Manolis Savva, Angel X. Chang, and Thomas A. Funkhouser. 2017. Learning Where to Look: Data-Driven Viewpoint Set Selection for 3D Scenes. *arXiv preprint arXiv:1704.02393* (2017).

Khang Truong Giang, Soohwan Song, Daekyum Kim, and Sunghee Choi. 2021. Sequential Depth Completion With Confidence Estimation for 3D Model Reconstruction. *IEEE Robotics and Automation Letters* 6 (2021), 327–334.

Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. 2020. Cascade Cost Volume for High-Resolution Multi-View Stereo and Stereo Matching. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*. 2492–2501.

Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask R-CNN. In *Proc. Int. Conf. on Computer Vision*. 2980–2988.

Benjamin Hepp, Debadeepta Dey, Sudipta N. Sinha, Ashish Kapoor, Neel Joshi, and Otmar Hilliges. 2018a. Learn-to-Score: Efficient 3D Scene Exploration by Predicting View Utility. *Proc. Euro. Conf. on Computer Vision* 11219 (2018), 455–472.

Benjamin Hepp, Matthias Nießner, and Otmar Hilliges. 2018b. Plan3D: Viewpoint and Trajectory Optimization for Aerial Multi-View Stereo Reconstruction. *ACM Trans. on Graphics* 38, 1 (2018), 4:1–4:17.

Xinyu Huang, Peng Wang, Xinjing Cheng, Dingfu Zhou, Qichuan Geng, and Ruigang Yang. 2019. The Apolloscape Open Dataset for Autonomous Driving and its Application. *IEEE Trans. Pattern Analysis & Machine Intelligence* 42, 10 (2019), 2702–2719.

Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. 2017. Tanks and Temples: Benchmarking Large-Scale Scene Reconstruction. *ACM Trans. on Graphics* 36, 4 (2017), 78:1–78:13.

Tobias Koch, Marco Körner, and Friedrich Fraundorfer. 2019. Automatic and Semantically-aware 3D UAV Flight Planning for Image-based 3D Reconstruction. *Remote Sensing* 11, 13 (2019), 1550.

Zhaoshuo Li, Xingtong Liu, Nathan Drenkow, Andy Ding, Francis X. Creighton, Russell H. Taylor, and Mathias Unberath. 2021. Revisiting Stereo Depth Estimation From a Sequence-to-Sequence Perspective With Transformers. In *Proc. Int. Conf. on Computer Vision*. 6197–6206.

Liqiang Lin, Yilin Liu, Yue Hu, Xingguang Yan, Ke Xie, and Hui Huang. 2022. Capturing, Reconstructing, and Simulating: the UrbanScene3D Dataset. In *Proc. Euro. Conf. on Computer Vision*.

Yilin Liu, Ruiqi Cui, Ke Xie, Minglun Gong, and Hui Huang. 2021a. Aerial Path Planning for Online Real-Time Exploration and Offline High-Quality Reconstruction of Large-Scale Urban Scenes. *ACM Trans. on Graphics* 40, 6 (2021), 226:1–226:16.

Yilin Liu, Ke Xie, and Hui Huang. 2021b. VGF-Net: Visual-Geometric Fusion Learning for Simultaneous Drone Navigation and Height Mapping. *Graphical Models* 116 (2021), 101108.

Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. 2020. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *Proc. Euro. Conf. on Computer Vision*. 405–421.

Cheng Peng and Volkan Isler. 2019. Adaptive View Planning for Aerial 3D Reconstruction. In *Proc. IEEE Int. Conf. on Robotics & Automation*. 2981–2987.

Mike Roberts, Debadeepta Dey, Anh Truong, Sudipta Sinha, Shital Shah, Ashish Kapoor, Pat Hanrahan, and Neel Joshi. 2017. Submodular Trajectory Optimization for Aerial 3D Scanning. In *Proc. Int. Conf. on Computer Vision*. 5324–5333.

Lukas Schmid, Michael Pantic, Raghav Khanna, Lionel Ott, Roland Siegwart, and Juan Nieto. 2020. An Efficient Sampling-Based Method for Online Informative Path Planning in Unknown Environments. *IEEE Robotics and Automation Letters* 5 (2020), 1500–1507.

Neil Smith, Nils Moehrle, Michael Goesele, and Wolfgang Heidrich. 2018. Aerial Path Planning for Urban Scene Reconstruction: A Continuous Optimization Method and Benchmark. *ACM Trans. on Graphics* 37, 6 (2018), 183:1–183:15.

Soohwan Song, Daekyum Kim, and Sungho Jo. 2020. Online Coverage and Inspection Planning for 3D Modeling. *Autonomous Robots* 44 (2020), 1431–1450.

Yifan Sun, Qixing Huang, Dun-Yu Hsiao, Li Guan, and Gang Hua. 2021. Learning View Selection for 3D Scenes. In *Proc. IEEE Conf. on Computer Vision & Pattern Recognition*. 14464–14473.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Proc. Conf. on Neural Information Processing Systems*. 5998–6008.

Yiheng Xie, Towaki Takikawa, Shunsuke Saito, Or Litany, Shiqin Yan, Numair Khan, Federico Tombari, James Tompkin, Vincent Sitzmann, and Srinath Sridhar. 2022. Neural Fields in Visual Computing and Beyond. *Computer Graphics Forum* 41, 2 (2022), 641–676.

Han Zhang, Yucong Yao, Ke Xie, Chi-Wing Fu, Hao Zhang, and Hui Huang. 2021. Continuous Aerial Path Planning for 3D Urban Scene Reconstruction. *ACM Trans. on Graphics* 40, 6 (2021), 225:1–225:15.

Xiaohui Zhou, Ke Xie, Kai Huang, Yilin Liu, Yang Zhou, Minglun Gong, and Hui Huang. 2020. Offsite Aerial Path Planning for Efficient Urban Scene Reconstruction. *ACM Trans. on Graphics* 39, 6 (2020), 192:1–192:16.