# OccuFly: A 3D Vision Benchmark for
# Semantic Scene Completion from the Aerial Perspective

Markus Gross[1,2,3,*]     Sai B. Matha [1]     Aya Fahmy[1]

Rui Song[4]     Daniel Cremers[2,3]     Henri Meeß[1]

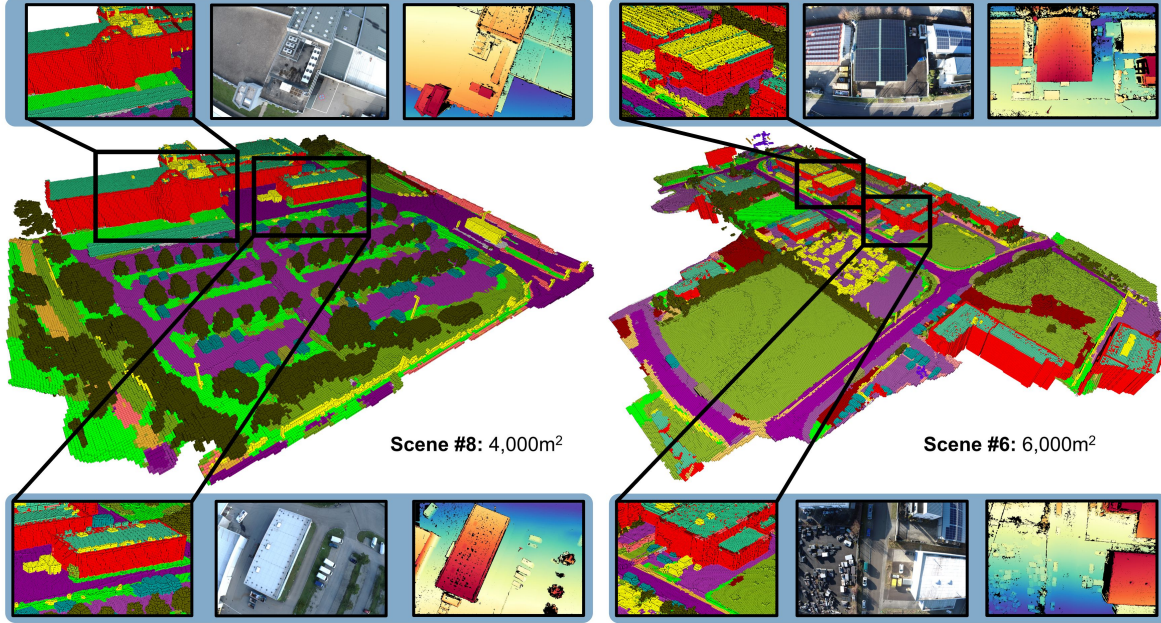[1]Fraunhofer IVI     [2]TU Munich     [3]MCML     [4]UCLA

Figure 1. OccuFly introduces the first real-world, aerial 3D SSC benchmark dataset, consisting of 9 scenes that provide over 20,000 samples of RGB images, semantic occupancy grids, and metric depth maps, including 22 semantic classes. OccuFly covers almost $200,000\,\mathrm{m}^2$ at $50\,\mathrm{m}$, $40\,\mathrm{m}$, and $30\,\mathrm{m}$ altitude in urban, industrial, and rural scenarios during spring, summer, fall, and winter. Zoom in for best view.

## Abstract

*Semantic Scene Completion (SSC) is crucial for 3D perception in mobile robotics, as it enables holistic scene understanding by jointly estimating dense volumetric occupancy and per-voxel semantics. Although SSC has been widely studied in terrestrial domains such as autonomous driving, aerial scenarios like autonomous flying remain largely unexplored, thereby limiting progress on downstream applications. Furthermore, LiDAR sensors represent the primary modality for SSC data generation, which poses challenges for most uncrewed aerial vehicles (UAVs) due to flight regulations, mass and energy constraints, and the sparsity of LiDAR-based point clouds from elevated viewpoints. To address these limitations, we introduce OccuFly, the first real-world, camera-based aerial SSC benchmark, captured at altitudes of $50\,\mathrm{m}$, $40\,\mathrm{m}$, and $30\,\mathrm{m}$ during spring, summer, fall, and winter. OccuFly covers urban, industrial, and rural scenarios, provides 22 semantic classes, and the data format adheres to established conventions to facilitate seamless integration with existing research. Crucially, we propose a LiDAR-free data generation framework based on camera modality, which is ubiquitous on modern UAVs. By utilizing traditional 3D reconstruction, our framework automates label transfer by lifting a subset of annotated 2D masks into the reconstructed point cloud, thereby substantially minimizing manual 3D annotation effort. Finally, we benchmark the state-of-the-art on OccuFly and highlight challenges specific to elevated viewpoints, yielding a comprehensive vision benchmark for holistic aerial 3D scene understanding. Code available at https://github.com/markus-42/occufly.*

# 1. Introduction

Modern approaches in 3D computer vision enable holistic scene understanding for downstream applications such as autonomous navigation, surveillance, and augmented reality [19]. To this end, one essential approach is Semantic Scene Completion (SSC) [10], which jointly infers the complete geometry of a 3D scene from a sparse observation, such as a camera image or LiDAR scan, while simultaneously assigning semantic labels to each element, typically represented in a voxelized 3D occupancy grid [67]. Recently, SSC has been extended to Panoptic Scene Completion by incorporating instance-level awareness [28]. While strong industry funding has led the research community to focus on terrestrial scenarios such as autonomous driving, low-altitude aerial scene understanding for autonomous flying of uncrewed aerial vehicles (UAVs) is largely restricted to either 2D datasets [2, 9, 12, 13, 26, 34, 39, 54, 55, 57, 59, 68, 71, 76, 78, 81, 87], or 3D mesh and point cloud datasets [4, 15, 44, 49, 52, 70, 72, 83]. Notably, the SSC objective remains largely unexplored from the aerial perspective, as no dedicated real or synthetic datasets exist, thereby confining related downstream applications to terrestrial scenarios.

Technically, SSC datasets are typically generated by (i) fusing multiple sparse LiDAR sweeps with registered poses to capture occluded regions as a dense point cloud, where (ii) each point is manually annotated with semantic labels and (iii) subsequently voxelized to produce SSC ground-truth [86]. While effective for ground vehicles, such LiDAR-based data generation becomes challenging in aerial scenarios. First, multi-modal sensor setups for UAVs are not as widespread or advanced as in autonomous driving, since UAV platforms are subject to strict flight regulations, such as US [24] or EU [22] regulations, and they must adhere to stringent mass and energy constraints, which conflict with the heavier and more power-demanding nature of LiDARs compared to cameras. Second, LiDAR sparsity persists and may even worsen from an elevated vantage point, leaving significant areas unobserved and unlabeled, which in turn would yield incomplete or low-quality ground-truth.

To address these limitations, we introduce OccuFly, the first real-world, low-altitude 3D vision benchmark for aerial Semantic Scene Completion. Crucially, we propose a data generation framework that is based on camera modality, which is considered to be ubiquitous on modern UAVs. Our dataset provides over 20,000 samples, resulting in $5\times$ the number of samples and $6\times$ the number of voxels compared to SemanticKITTI [5], which introduced the first SSC dataset for autonomous driving. Furthermore, we evaluate the state-of-the-art on OccuFly, yielding a comprehensive 3D aerial vision benchmark.

Our **contributions** can be summarized as follows:

- We present OccuFly, a real-world aerial vision benchmark consisting of 9 scenes that provide over 20,000 samples of nadir and oblique perspective images with corresponding 3D semantic voxel grids, including 22 semantic classes. OccuFly covers almost $200,000\,\mathrm{m}^2$ at $50\,\mathrm{m}$, $40\,\mathrm{m}$, and $30\,\mathrm{m}$ altitude in urban, industrial, and rural scenarios during spring, summer, fall, and winter.
- In addition to the SSC samples, OccuFly offers more than 20,000 per-frame metric depth maps. Additionally, we train and release the Depth-Anything-V2 [82] depth estimator on these depth maps, enabling state-of-the-art SSC.
- We propose a novel and scalable data generation framework to construct SSC ground-truth, thereby (i) relying on camera modality to avoid LiDAR-based point cloud sparsity, (ii) avoiding LiDAR hardware to adhere to mass and energy constraints of most UAVs, and (iii) reducing manual semantic labeling from tedious 3D annotation to efficient 2D annotation.
- Upon acceptance, we will release the OccuFly dataset, facilitating reproducibility and further research on aerial 3D Semantic Scene Completion.

# 2. Related Work

**Datasets and Benchmarks.** Apart from indoor SSC [65], the first outdoor SSC dataset was proposed by SemanticKITTI [5] for autonomous driving. Despite its impact, the limited scale and diversity of SemanticKITTI impeded the development of generalizable SSC models and their comprehensive evaluation [46]. To address these limitations, multiple impactful datasets followed, such as nuScenes [8], Waymo [66], and KITTI-360 [48], all relying on LiDAR modality for data generation. As elaborated in Sec. 1, LiDAR-based SSC ground-truth generation faces fundamental challenges stemming from sparse point clouds, as many regions remain unobserved and consequently unlabeled. Occlusions, misalignment from aggregating multiple sweeps, and dynamic objects further exacerbate these gaps, creating inaccuracies in the resulting volumetric labels. Additionally, the process of manually annotating such sparse 3D point clouds is both time-consuming and error-prone, undermining scalability for large-scale 3D reconstruction tasks.

To mitigate these challenges, several benchmarks have emerged. Specifically, Occ3D [69] complements multi-sweep point cloud densification with mesh reconstruction and camera-based filtering to reduce occlusions and mislabeled voxels. Meanwhile, OpenOccupancy [73] adopts an Augmenting and Purifying pipeline, combining pseudo-label generation with extensive human annotation to double the density of occupancy labels and refine boundaries for 360° semantic coverage. In parallel, OCFBench [51] synchronizes dynamic objects via bounding box labels to address spatiotemporal occlusions, and excludes unknown voxels through ray-casting, thereby mitigating LiDAR spar-
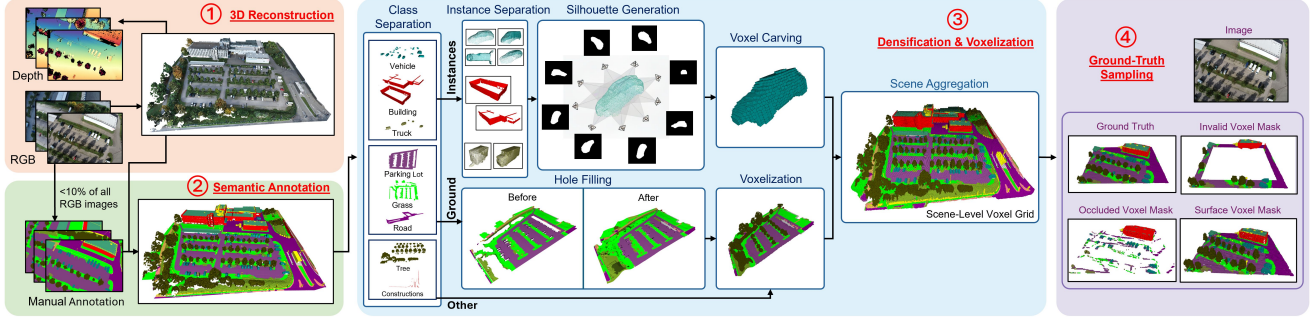
Figure 2. **Proposed image-based data generation framework.** An overview is provided in Sec. 3.1. Zoom in for best view.

sity for 4D occupancy completion and forecasting. Similarly, SSCBench [46] excludes unobserved voxels from training and evaluation, and integrates data from multiple sources to enhance geographic diversity. Additionally, Open-Scene [14] leverages voxel densification, scene completion, and flow integration to fill unobserved LiDAR regions. Finally, UniOcc [75] unifies real and synthetic (CARLA [63], OpenCOOD [80]) datasets, and provides label-free evaluation metrics to address suboptimal LiDAR coverage.

In contrast, our proposed data generation framework reconstructs metric 3D point clouds from geo-referenced imagery. We project manually annotated 2D semantic masks into these reconstructions via 2D–3D correspondences, yielding accurate per-point semantics, which substantially minimizing 3D annotation effort. We further refine and densify the semantic point cloud, voxelize it, and finally extract consistent per-frame ground-truth samples. Sec. 3.1 provides a chronological overview.

**Methods.** Apart from recurrent [74] and multi-view SSC methods [11, 29, 63, 77], single-view SSC was pioneered by MonoScene [10], which bridges 2D and 3D representations through an optics-inspired feature projection and a 3D context relation prior. Recent advancements have explored diverse feature representations to bridge 2D image cues with volumetric representations, broadly categorized into planar and voxel-based approaches [79].

In particular, planar representations, such as bird's-eye-view (BEV) [47] or tri-perspective views (TPV), enable compact feature aggregation and long-range context modeling. TPVFormer [32] introduces the TPV framework with perpendicular planes and a transformer encoder for lifting multi-view features into semantic occupancy grids, while DISC [50] extends BEV with discriminative queries and dual-attention decoders to disentangle instance-scene contexts for category-specific interactions. Building on hybrid designs, CGFormer [84] fuses voxel and TPV spaces via context-aware queries and 3D deformable cross-attention, enhancing depth-aware aggregation for superior fidelity.

In contrast, voxel-based methods directly operate on 3D grids for fine-grained occupancy and semantics. Two-stage

transformers like VoxFormer [45] and VisHall3D [53] employ sparse queries from depth priors followed by densification or visibility-aware decoupling to handle visible and occluded regions, respectively. Symphonies [36] refines voxel queries with class-centric instance propagation for dynamic 2D-3D reasoning, whereas SOAP [41] incorporates adaptive decoders drawing from semantic repositories.

Specialized enhancements include ScanSSC's [3] tri-axis voxel scanning for distant geometry. Furthermore, other methods utilize State Space Models [42, 43], as well as implicit scenes via self-supervised learning [31, 35] and Gaussian representations [33]. We refer the interested reader to the surveys of [79, 86] for a comprehensive discussion.

## 3. Data Generation Framework

### 3.1. Overview

Our proposed data generation framework consists of four modules, shown in Fig. 2. We utilize geo-referenced images to apply traditional multi-view reconstruction, generating a metric 3D point cloud (Sec. 3.3.1). This approach additionally yields 2D–3D correspondences, allowing image pixels to be associated with reconstructed 3D points, effectively streamlining the creation of 3D semantic annotations. Technically, we enable highly efficient label transfer by manually annotating only a small subset of the camera images ($<10\%$ on average) and lifting the semantic pixels into the reconstructed point cloud (Sec. 3.3.2). This reduces costly 3D annotation to efficient 2D image labeling, substantially lowering annotation effort. Subsequently, we densify individual objects with our novel densification pipeline and eventually voxelize the semantic point cloud (Sec. 3.3.3). As all previous steps are performed on a global scene level, we finally retrieve per-frame ground-truth grids by frustum-culling the scene voxel grid using geo-referenced camera poses and intrinsics, resulting in one fixed-size semantic voxel grid per camera frame. Similar to SemanticKITTI [6], we additionally construct binary masks to distinguish surface, occluded, and invalid voxels (Sec. 3.3.4).

## 3.2. Problem Formulation

Given a set of $N$ geo-referenced, calibrated RGB images $\mathcal{I} = \{\mathbf{I}_n \in \mathbb{R}^{H \times W \times 3}\}_{n=1}^N$, acquired by a pinhole camera with per-frame intrinsics $\mathcal{K} = \{\mathbf{K}_n \in \mathbb{R}^{3 \times 3}\}_{n=1}^N$, and world-to-camera poses $\mathcal{T} = \{\mathbf{T}_{c \leftarrow w}^n \in \mathrm{SE}(3)\}_{n=1}^N$, we adopt a fixed-size metric voxel-grid specification $(X, Y, Z, r)$, where $X, Y, Z$ denote the number of voxels along cartesian axes, with voxel edge length $r > 0$. Moreover, we define a semantic label set $\mathcal{C} = \{1, 2, \ldots, C\}$. Each image $\mathbf{I}_n$ is coupled to a ground-truth semantic voxel grid $\mathbf{Y}_n \in \mathcal{C}^{X \times Y \times Z}$, yielding the dataset as a set of image–grid samples $\{(\mathbf{I}_n, \mathbf{Y}_n)\}_{n=1}^N$, organized into multiple scenes with scene-dependent sample counts.

## 3.3. Method

### 3.3.1. 3D Reconstruction

From the calibrated camera intrinsics $\mathcal{K}$, and the geo-referenced images $\mathcal{I}$ and their poses $\mathcal{T}$, we obtain a metric scene reconstruction via Structure-from-Motion (SfM) [60] and Multi-View Stereo (MVS) [61], abstracted as

$$(\mathcal{P}, \mathcal{D}, \mathcal{A}) = \Psi_{\mathrm{SfM+MVS}}(\mathcal{I}, \mathcal{K}, \mathcal{T}), \tag{1}$$

where the set $\mathcal{P} = \{\mathbf{x}_m \in \mathbb{R}^3\}_{m=1}^M$ represents a dense point cloud with $M = |\mathcal{P}|$ number of points in world coordinates, and $\mathcal{D} = \{\mathbf{D}_n \in \mathbb{R}_{\geq 0}^{H \times W}\}_{n=1}^N$ are per-image metric depth maps. Moreover, $\mathcal{A}_n$ denotes per-image 2D–3D correspondences , which satisfy the projection of a 3D point $\mathbf{x}$ onto a 2D image pixel $(u, v)$, formulated as

$$\mathcal{A}_n = \left\{ ((u, v), \mathbf{x}) \,\big|\, (u, v, 1)^\top \sim \mathbf{K}_n[\mathbf{R}_n \,|\, \mathbf{t}_n][\mathbf{x}^\top \, 1]^\top \right\}, \tag{2}$$

where $\sim$ denotes equality up to a non-zero scalar, and $\mathbf{R}_n$ and $\mathbf{t}_n$ are the rotational and translational components of $\mathbf{T}_{c \leftarrow w}^n$, respectively. All correspondences of a scene are given by $\mathcal{A} = \bigcup_{n=1}^N \mathcal{A}_n$.

### 3.3.2. Semantic Annotation

We define a finite, non-empty set of semantic classes $\mathcal{C} = \{1, \ldots, C\}$. To minimize manual annotation effort, we annotate only a small subset $\mathcal{J} \subset \{1, \ldots, N\}$ of images, exploiting the fact that each 3D point is observed by multiple cameras. Note that $\mathcal{J}$ denotes the index set of images, not the images themselves. We select $\mathcal{J}$ by spatially stratified sampling over the scene area. More specifically, we partition the ground plane into a regular grid with square cells of $25\,\mathrm{m}$ edge length. For each cell center, we select the image whose pose is closest to that location. Additionally, for cells at the scene border, we select the pose closest to the border. Using the 2D–3D correspondences $\mathcal{A}$ from Sec. 3.3.1, we quantify the coverage of reconstructed points by

$$\rho(\mathcal{J}) = \frac{\left| \{\, \mathbf{x} \in \mathcal{P} \,|\, \exists n \in \mathcal{J}, \, \exists (u, v) \text{ s.t. } ((u, v), \mathbf{x}) \in \mathcal{A}_n \,\} \right|}{|\mathcal{P}|} \tag{3}$$

Empirically, $\rho(\mathcal{L}) > 0.99$, while $|\mathcal{J}|/N < 0.1$ on average, implying that annotating less than $10\%$ of all images cover more than $99\%$ of all 3D points, further detailed in in Sec. 4.3.

After manually annotating all images in $\mathcal{J}$, we lift semantic labels from images to points via back-projection. As most of the points are observed by multiple cameras, we fuse multi-view evidence to assign robust per-point labels. Technically, this is accomplished by unweighted majority voting, where ties are broken by a fixed class-prior order derived from class frequencies. Furthermore, to annotate unlabeled points, we apply k-nearest-neighbor (kNN) with inverse-distance weights within a fixed neighborhood. As a denoising step, we apply a second iteration of kNN to all points, effectively relabeling every point to the dominant class in its neighborhood, similar to [69]. The resulting semantic point cloud $\mathcal{P}_{\mathcal{S}} = \{(\mathbf{x}_m, c_m)\}_{m=1}^M$ consists of 3D points $\mathbf{x}_m$ and their corresponding semantic class labels $c_m \in \mathcal{C}$.

### 3.3.3. Class-Aware Densification and Voxelization

**Preliminaries.** Given the semantic point cloud $\mathcal{P}_{\mathcal{S}} = \{(\mathbf{x}_m, c_m)\}_{m=1}^M$ from Sec. 3.3.2, we first partition the semantic classes into three disjoint groups

$$\mathcal{C}_{\mathrm{inst}} \cup \mathcal{C}_{\mathrm{gnd}} \cup \mathcal{C}_{\mathrm{oth}} = \mathcal{C}, \qquad \mathcal{C}_{\mathrm{inst}} \cap \mathcal{C}_{\mathrm{gnd}} \cap \mathcal{C}_{\mathrm{oth}} = \varnothing, \tag{4}$$

corresponding to instance classes $\mathcal{C}_{\mathrm{inst}}$ that are to be densified object-wise (*e.g.*, vehicles), ground classes $\mathcal{C}_{\mathrm{gnd}}$ that are to be surface-reconstructed (*e.g.*, road), and other classes $\mathcal{C}_{\mathrm{oth}}$ that are directly voxelized (*e.g.*, constructions). Based on this class partition, we apply **group separation** to retrieve group-specific point cloud subsets

$$\mathcal{P}_{\mathrm{inst}} = \{(\mathbf{x}, c) \in \mathcal{P}_{\mathcal{S}} \,|\, c \in \mathcal{C}_{\mathrm{inst}}\} \tag{5}$$

$$\mathcal{P}_{\mathrm{gnd}} = \{(\mathbf{x}, c) \in \mathcal{P}_{\mathcal{S}} \,|\, c \in \mathcal{C}_{\mathrm{gnd}}\} \tag{6}$$

$$\mathcal{P}_{\mathrm{oth}} = \{(\mathbf{x}, c) \in \mathcal{P}_{\mathcal{S}} \,|\, c \in \mathcal{C}_{\mathrm{oth}}\} \tag{7}$$

Furthermore, let $\mathcal{G} = [X] \times [Y] \times [Z]$ be a voxel grid. For any point cloud subset $\mathcal{Q} \subset \mathcal{P}$ and a target voxel resolution $r > 0$, we denote by $\mathrm{Vox}_r(\mathcal{Q}) \subset \mathcal{G}$ the set of occupied voxels obtained by standard binning (point rasterization) [23] at resolution $r$. For triangle meshes, we use the same notation to indicate triangle-to-voxel scan conversion [37].

**Instance classes.** We apply **instance separation** by decomposing $\mathcal{P}_{\mathrm{inst}}$ into object instances using Euclidean clustering via DBSCAN [21] with class-specific parameters $(\varepsilon_c, \mathrm{minPts}_c)$. This process yields $J$ object instances

$$\mathbb{S} = \mathrm{DBSCAN}(\mathcal{P}_{\mathrm{inst}}, \varepsilon_c, \mathrm{minPts}_c) = \{\mathcal{S}_j \subset \mathcal{P}_{\mathrm{inst}}\}_{j=1}^J, \tag{8}$$

where $\mathcal{S}_j$ represents an instance point cloud.

To create the voxelized visual hull from an instance point cloud $\mathcal{S}$, we perform two major steps: Silhouette extraction

and silhouette-based voxel carving. For each instance $\mathcal{S} \in \mathbb{S}$, let $\mathrm{pos}(\mathcal{S}) = \{\mathbf{x} \mid (\mathbf{x}, c) \in \mathcal{S}\}$ denote its 3D positions.

To **extract silhouettes**, we (i) place $K$ virtual cameras $\mathbb{V} = \{(\mathbf{K}_k, \mathbf{T}_{c \leftarrow w}^k)\}_{k=1}^K$ quasi-uniformly distributed on the viewing sphere around $\mathrm{pos}(\mathcal{S})$ (in practice we use $K{=}24$); (ii) project $\mathrm{pos}(\mathcal{S})$ to each view to obtain 2D point sets $\mathcal{U}_k = \{\pi_k(\mathbf{x}) \mid \mathbf{x} \in \mathrm{pos}(\mathcal{S})\}$, where $\pi_k$ is the pinhole projection induced by $(\mathbf{K}_k, \mathbf{T}_{c \leftarrow w}^k)$; and (iii) compute a binary silhouette $\Omega_k \subset \mathbb{R}^2$ via the $\alpha$-shape boundary [20] of $\mathcal{U}_k$, where $\Omega_k$ is the set of pixels that belong to the instance.

For **voxel carving**, we apply multi-view silhouette carving [40]. To this end, we back-project each silhouette to a generalized cone

$$\mathcal{R}_k = \{\mathbf{x} \in \mathbb{R}^3 : \pi_k(\mathbf{x}) \in \Omega_k\}. \qquad (9)$$

The continuous instance hull is $\mathcal{H}(\mathcal{S}) = \bigcap_{k=1}^K \mathcal{R}_k$. To carve within a finite space around the instance hull, let $\mathcal{B}(\mathcal{S})$ be a tight axis-aligned bounding box of $\mathrm{pos}(\mathcal{S})$, dilated by a small margin. We discretize $\mathcal{B}(\mathcal{S})$ into a 3D grid of voxels indexed by $\mathbf{v} = (i, j, k)$. The carved occupancy set is

$$\mathcal{O}_{\mathrm{inst}}(\mathcal{S}) = \{\mathbf{v} \in \mathcal{G} \mid \mathrm{center}(\mathbf{v}) \in \mathcal{B}(\mathcal{S}) \cap \mathcal{H}(\mathcal{S})\}. \quad (10)$$

All voxels $\mathbf{v} \in \mathcal{O}_{\mathrm{inst}}(\mathcal{S})$ receive the semantic label of its original semantic instance point cloud $\mathcal{P}_{\mathrm{inst}}$. Finally, we aggregate all instances via $\mathcal{O}_{\mathrm{inst}} = \bigcup_{\mathcal{S} \in \mathbb{S}} \mathcal{O}_{\mathrm{inst}}(\mathcal{S})$.

**Ground classes.** We densify $\mathcal{P}_{\mathrm{gnd}}$ via Poisson surface reconstruction [38] to obtain a watertight triangle mesh $\mathcal{M}_{\mathrm{gnd}} = \Psi_{\mathrm{Poisson}}(\mathcal{P}_{\mathrm{gnd}})$, which **fills holes** and enforces surface continuity, similar to [77]. We then voxelize the mesh:

$$\mathcal{O}_{\mathrm{gnd}} = \mathrm{Vox}_r(\mathcal{M}_{\mathrm{gnd}}), \qquad (11)$$

assigning to each occupied voxel the majority ground class of contributing mesh samples in its cell.

**Other classes.** For $\mathcal{P}_{\mathrm{oth}}$, we apply direct voxelization:

$$\mathcal{O}_{\mathrm{oth}} = \mathrm{Vox}_r(\mathcal{P}_{\mathrm{oth}}), \qquad (12)$$

with per-voxel semantics determined by majority voting of points falling into the voxel.

**Aggregation.** We construct a scene-level semantic voxel grid by combining all groups with a fixed precedence order $\mathrm{inst} \succ \mathrm{oth} \succ \mathrm{gnd}$ to resolve label conflicts:

$$\mathbf{Y}(\mathbf{v}) = \begin{cases} \text{label from } \mathcal{O}_{\mathrm{inst}}, & \mathbf{v} \in \mathcal{O}_{\mathrm{inst}}, \\ \text{label from } \mathcal{O}_{\mathrm{oth}}, & \mathbf{v} \in \mathcal{O}_{\mathrm{oth}} \setminus \mathcal{O}_{\mathrm{inst}}, \\ \text{label from } \mathcal{O}_{\mathrm{gnd}}, & \mathbf{v} \in \mathcal{O}_{\mathrm{gnd}} \setminus (\mathcal{O}_{\mathrm{inst}} \cup \mathcal{O}_{\mathrm{oth}}), \\ 0, & \text{otherwise,} \end{cases}$$
$$(13)$$

where $0$ denotes empty. The resulting $\mathbf{Y}$ constitutes the semantic voxel grid for the whole scene.

### 3.3.4. Ground-Truth Sampling

**Frustum Culling.** Given the scene-level semantic voxel grid $\mathbf{Y} \in (\{0\} \cup \mathcal{C})^{X \times Y \times Z}$ and per-frame camera parameters $(\mathbf{K}_n, \mathbf{T}_{c \leftarrow w}^n)$, we construct per-frame ground-truth by frustum-culling and rasterization at a fixed metric specification $(X, Y, Z, r)$. Let $\pi_n$ denote the pinhole projection induced by $(\mathbf{K}_n, \mathbf{T}_{c \leftarrow w}^n)$, and let $[d_{\min}, d_{\max}]$ be near/far clipping distances, respectively. Define the truncated frustum

$$\mathcal{F}_n = \{\mathbf{x} \in \mathbb{R}^3 : \pi_n(\mathbf{x}) \in [0, W] \times [0, H]\},$$

with $d_{\min} \leq d_n(\mathbf{x}) \leq d_{\max}$, where $d_n(\mathbf{x})$ is the camera-centric depth of $\mathbf{x}$. We obtain the per-frame grid $\mathbf{Y}_n \in (\{0\} \cup \mathcal{C})^{X \times Y \times Z}$ by discretizing $\mathcal{F}_n$ at resolution $r$ and sampling $\mathbf{Y}$ at voxel centers:

$$\mathbf{Y}_n(\mathbf{v}) = \begin{cases} \mathbf{Y}(\mathbf{v}), & \text{if center}(\mathbf{v}) \in \mathcal{F}_n, \\ 0, & \text{otherwise,} \end{cases} \quad \mathbf{v} \in \mathcal{G}.$$

**Binary Masks.** In addition, we construct three binary masks, similar to SemanticKITTI [6]: invalid $\mathbf{M}_n^{\mathrm{inv}}$, surface $\mathbf{M}_n^{\mathrm{surf}}$, and occluded $\mathbf{M}_n^{\mathrm{occ}}$, all in $\{0, 1\}^{X \times Y \times Z}$.

*Invalid Mask.* The invalid mask represents voxels outside the field of view:

$$\mathbf{M}_n^{\mathrm{inv}}(\mathbf{v}) = \mathbf{1}[\mathrm{center}(\mathbf{v}) \notin \mathcal{F}_n], \qquad (14)$$

where $\mathbf{1}[\cdot]$ denotes the indicator (Iverson) function.

*Surface Mask (view-independent).* Let $\mathcal{N}_6(\mathbf{v})$ be the 6-neighborhood in the grid, and let $\mathcal{E}(\mathbf{v}) = \mathbf{1}[\mathbf{Y}_n(\mathbf{v}) \neq 0]$ denote occupancy. We mark geometric boundary voxels by

$$\mathbf{M}_n^{\mathrm{surf}}(\mathbf{v}) = \mathbf{1}[\mathcal{E}(\mathbf{v}) = 1 \ \wedge \ \exists \mathbf{u} \in \mathcal{N}_6(\mathbf{v}) \mid \mathcal{E}(\mathbf{u}) = 0].$$

*Occluded Mask (view-dependent).* For each pixel $(u, v) \in [0, W] \times [0, H]$, consider the ordered set of frustum voxels $\mathcal{V}_n(u, v) = \langle \mathbf{v}_1, \mathbf{v}_2, \dots \rangle$ traversed by the ray $\rho_n(u, v)$ from near to far. Let $i^\star = \min\{i : \mathcal{E}(\mathbf{v}_i) = 1\}$ if such $i$ exists. Then we assign

$$\mathbf{M}_n^{\mathrm{occ}}(\mathbf{v}_j) = \begin{cases} 0, & j = i^\star, \\ 1, & j > i^\star \text{ and } \mathcal{E}(\mathbf{v}_j) = 1, \\ 0, & \text{otherwise,} \end{cases}$$

for all $\mathbf{v}_j \in \mathcal{V}_n(u, v)$. We set $\mathbf{M}_n^{\mathrm{occ}}(\mathbf{v}) = 0$ if $\mathbf{v}$ lies outside all rays or is invalid.

Finally, the per-frame ground-truth sample consists of $(\mathbf{I}_n, \mathbf{Y}_n, \mathbf{M}_n^{\mathrm{inv}}, \mathbf{M}_n^{\mathrm{surf}}, \mathbf{M}_n^{\mathrm{occ}})$, where invalid voxels are excluded from evaluation, surface voxels represent view-independent geometric boundaries, and occluded voxels capture view-dependent occupied regions behind the first visible surface along camera rays.

Table 1. OccuFly dataset statistics, discussed in Sec. 4.2.

| Scene | Season | Scenario | Area [m²] | Number of Samples | | | |
|---|---|---|---|---|---|---|---|
| | | | | 50 m | 40 m | 30 m | Total |
| **Training** | | | | | | | **14,321** |
| 01 | Winter | Rural | 41,234 | 421 | 469 | 512 | 1402 |
| 02 | Winter | Urban | 8,529 | 294 | 474 | 132 | 900 |
| 03 | Spring | Urban | 25,077 | 1,039 | 1,475 | 1,531 | 4,045 |
| 04 | Spring | Industrial | 55,579 | 1,157 | 1,312 | 1,571 | 4,040 |
| 05 | Summer | Rural | 24,810 | 997 | 1,682 | 1,255 | 3,934 |
| **Validation** | | | | | | | **1,997** |
| 06 | Spring | Urban | 5,428 | 283 | 385 | 343 | 1,011 |
| 07 | Spring | Industrial | 5,802 | 375 | 272 | 339 | 986 |
| **Test** | | | | | | | **3,842** |
| 08 | Fall | Industrial | 4,314 | 188 | 316 | 389 | 893 |
| 09 | Spring | Urban | 23,165 | 1,292 | 1,416 | 241 | 2,949 |
| **Total** | | | **193.938** | **6,046** | **7,801** | **6,313** | **20,160** |

Table 2. Comparison of terrestrial and aerial vision-based SSC benchmarks, detailed in Sec. 4.2.

| | Camera Views | # of Samples | Depth Maps | Semantic Classes | Grid resolution | |
|---|---|---|---|---|---|---|
| | | | | | $X \times Y \times Z$ | $r$ |
| **Terrestrial Benchmarks** | | | | | | |
| SemanticKITTI [6] | single | 4,649 | ✗ | 19 | $256 \times 256 \times 32$ | 0.20 |
| OpenOccupancy [73] | multi | 34,149 | ✗ | 16 | $512 \times 512 \times 40$ | 0.20 |
| SSCBench-Waymo [46] | multi | 19,985 | ✗ | 14 | $256 \times 256 \times 32$ | 0.20 |
| SSCBench-nuScenes [46] | multi | 34,078 | ✗ | 16 | $256 \times 256 \times 32$ | 0.20 |
| SSCBench-KITTI-360 [46] | multi | 12,865 | ✗ | 19 | $256 \times 256 \times 32$ | 0.20 |
| Occ3D-Waymo [69] | multi | 200,000 | ✗ | 14 | $3200 \times 3200 \times 128$ | 0.05 |
| Occ3D-nuScenes [69] | multi | 40,000 | ✗ | 16 | $200 \times 200 \times 16$ | 0.40 |
| **Aerial Benchmarks** | | | | | | |
| OccuFly (ours) | single | 20,160 | ✓ | 22 | $192 \times 128 \times 128$ | 0.50 |

Table 3. Semantic class comparison with 2D aerial image datasets.

| | UDD [13] | VDD [9] | UAVid [54] | AeroScapes [57] | ICG [34] | SkyScapes [2] | OccuFly (ours) |
|---|---|---|---|---|---|---|---|
| # Classes | 4 | 7 | 8 | 11 | 20 | 31 (20) | 22 |

# 4. OccuFly Dataset

## 4.1. Data Collection

We utilize two DJI UAV platforms for photogrammetric data acquisition [27]: The Phantom 4 RTK (P4) [16] and the DJI Mavic 3 Enterprise Series (M3-ES) [17], capturing images at 5472×3648 and 4000×3000 pixels, respectively. Geo-referenced camera poses and orientations are recorded by each flight controller via onboard sensor fusion (GNSS, IMU, compass, and magnetometer) [18]. We collect data from nine scenes in urban, industrial, and rural scenarios within a single geographic region, spanning spring, summer, fall, and winter. Note that the locations of data collection are withheld during the review process but will be disclosed upon publication. For data acquisition, we executed automated double-grid flight patterns at altitudes of 50 m, 40 m, and 30 m. These missions yielded ground sampling distances (GSD) of 1.4, 1.1, and 0.8 cm/pixel with the P4 platform, and 6.7, 5.4, and 4.1 cm/pixel with the M3-ES platform. The double-grid pattern provided, on average, 67 % side and 74 % forward image overlap. For oblique acquisitions, camera tilt angles were set to −75° at 50 m and 40 m, and −70° at 30 m. Additionally, using the M3-ES UAV, we collected nadir (0° tilt) imagery at all altitudes.

## 4.2. Dataset Statistics

We summarize OccuFly in Tab. 1. The dataset comprises 9 scenes and more than 20,000 annotated samples, each including (i) an RGB image, (ii) a semantic occupancy grid, and (iii) a metric depth map. Voxel grids are annotated with 22 semantic classes, and per-class frequencies are reported in Fig. 4. OccuFly spans approximately 193,938 m² across altitudes of 50 m, 40 m, and 30 m, covering urban, industrial, and rural environments in spring, summer, fall, and winter. 3D space is discretized into voxel grids of resolution $192 \times 128 \times 128$ with voxel size $r = 0.5$ m. Moreover, the dataset follows the SSCbench [46] data organization structure. We report the depth distribution in the supplementary material.

Ground-truth at 40 m and 30 m is generated via frus-

tum culling (Sec. 3.3.4) of the 50 m scene-level semantic voxel grid. Furthermore, we provide group assignments (Sec. 3.3.3) for each semantic class in Fig. 3.

## 4.3. Dataset Evaluation

We assess OccuFly against established vision-based terrestrial SSC datasets (Tab. 2). Similar to SemanticKITTI [6], which introduced real-world SSC to autonomous driving, OccuFly introduces real-world SSC to the aerial domain, but at a substantially larger scale: the number of samples is more than five times higher, and the total number of labeled voxels is over six times larger than SemanticKITTI. OccuFly further provides the largest class taxonomy (22 classes) among the compared SSC datasets while adhering to SSCBench-style data structuring [46] for seamless integration.

*3D Reconstruction Quality.* We quantify geometric consistency of the SfM+MVS pipeline (Sec. 3.3.1) using the standard RMSE reprojection error. Across all scenes, the average reprojection RMSE of 1.24 pixels demonstrates strong geometric fidelity [30], reflecting precise SfM/MVS alignment under high-resolution inputs. Scene-wise reprojection errors and qualitative reconstructions for all scenes are provided in the supplementary, while representative results are shown in Fig. 2.

*Volumetric Completeness.* Although classical multi-view reconstruction may leave holes in texture-poor regions, our class-aware densification and voxelization (Sec. 3.3.3) remove such artifacts in the final volumetric ground-truth. Additional reconstructions and resulting scene-level grids are included in the supplementary.

*2D-to-3D Label Transfer Efficiency.* Semantic annotation leverages multi-view correspondences to lift a small subset of manually annotated images to 3D (Sec. 3.3.2). In practice, annotating fewer than 10% of the images per scene suffices to automatically label over 99% of reconstructed points, as
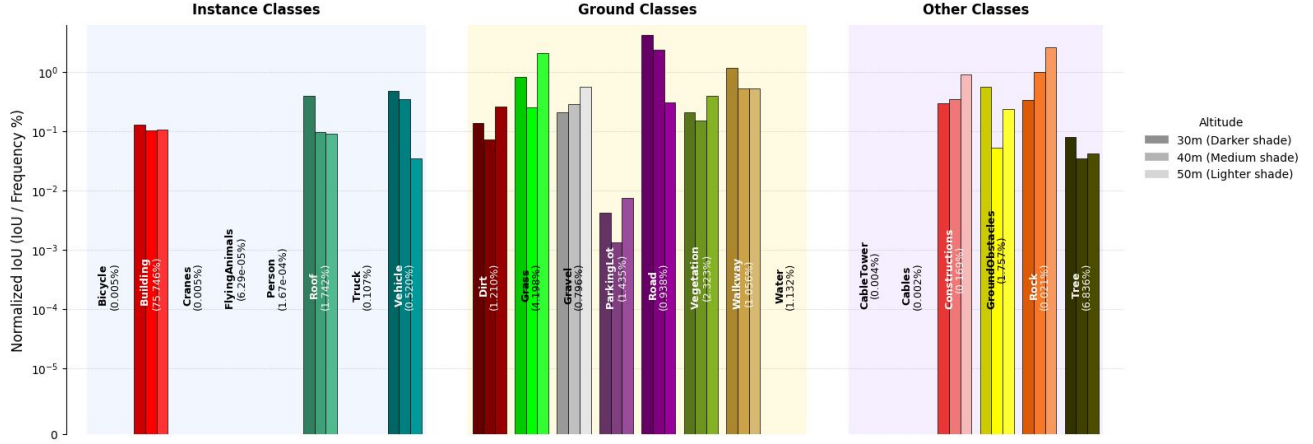
Figure 3. **Quantitative results on OccuFly test set.** We evaluate CGFormer [84] (i) class-wise, (ii) group-wise, and (iii) altitude-wise, effectively investigating implications on key characteristics that are natural to image-based aerial data generation and learning for SSC.

determined by Eq. (3). We report per-scene annotation ratios and coverage statistics in the supplementary. Semantic point clouds in Fig. 2 further illustrate the strong labeling fidelity.

*Semantic Taxonomy.* Beyond SSC, Tab. 3 positions OccuFly among established 2D aerial semantic segmentation datasets, where its 22-class taxonomy ranks second. While Skysapes [2] ranks first, 12 of its 31 classes are lane-markings, effectively reducing its distinct class count to 20. Consequently, OccuFly provides one of the most detailed aerial taxonomies to date, strengthening fine-grained semantic evaluation and enabling seamless comparability with established 2D benchmarks.

# 5. Benchmark Experiments

## 5.1. Experimental Setup

**Aerial Semantic Scene Completion.** We benchmark CGFormer [84], a state-of-the-art and established SSC method, and use official implementations and evaluation protocols to ensure scientific rigor. Since CGFormer utilizes depth maps to back-project geometric priors, we replace its MobileStereNetV2 [62] with our DAv2-OccuFly (see next paragraph). Regarding evaluation metrics, geometry is assessed by voxel-level Intersection-over-Union (IoU), and semantics by mean IoU (mIoU) over non-empty classes, in line with prior research [32, 45]. Specifically, we evaluate CGFormer [84] (i) class-wise, (ii) group-wise, and (iii) altitude-wise, presented in Sec. 5.2, effectively investigating implications on key characteristics that are natural to image-based aerial data generation and learning for SSC.

**Metric Monocular Depth Estimation.** This task is crucial for vision-based SSC, as state-of-the-art methods, such as CGFormer [84] and others [36, 45, 85], initialize 3D geometric priors via back-projecting metric depth maps (see Sec. 2). Notably, no established metric mono depth models exist for

Table 4. Altitude-wise metric monocular depth estimation performance with Depth Anything V2 (DAv2) [82] on OccuFly test set.

| Altitude | Method | Higher is better ↑ | | | Lower is better ↓ | | | |
|---|---|---|---|---|---|---|---|---|
| | | $\delta_1$ | $\delta_2$ | $\delta_3$ | AbsRel | RMSE | MAE | SILog |
| 30m | DAv2-metric | 0.005 | 0.029 | 0.050 | 0.737 | 22.961 | 22.493 | 0.147 |
| | DAv2-OccuFly (ours) | **0.919** | **0.982** | **0.998** | **0.103** | **3.108** | **2.666** | **0.086** |
| 40m | DAv2-metric | 0.002 | 0.024 | 0.111 | 0.693 | 25.031 | 23.918 | 0.208 |
| | DAv2-OccuFly (ours) | **0.795** | **0.957** | **0.998** | **0.148** | **4.486** | **3.823** | **0.119** |
| 50m | DAv2-metric | 0.000 | 0.000 | 0.003 | 0.767 | 34.576 | 33.515 | 0.191 |
| | DAv2-OccuFly (ours) | **0.844** | **0.997** | **1.000** | **0.129** | **5.985** | **5.261** | **0.114** |

the aerial domain. To this end, we evaluate the potential of OccuFly's metric depth maps by benchmarking Depth Anything V2 ViT-Small Metric (DAv2-metric) [82], a state-of-the-art metric monocular model. Technically, we follow the DAv2 metric adaptation protocol and fine-tune the affine-invariant model on OccuFly's training split with metric depth supervision, referred to as DAv2-OccuFly. We evaluate the performance with established metrics following [82].

## 5.2. Quantitative Results

**Aerial Semantic Scene Completion.** Our analysis in Fig. 3 reveals a strong correlation between performance and semantic class frequency, consistent across (i) classes, (ii) altitudes, and (iii) semantic groups (Sec. 3.3.3). This indicates that CGFormer is well-suited to aerial SSC under near-uniform class distributions. The insight underscores the need for balanced data and validates our scalable data-generation framework, enabling advances in aerial SSC from both methodological and data-generation perspectives.

**Metric Monocular Depth Estimation.** Altitude-wise results in Tab. 4 show that our DAv2-OccuFly consistently and substantially outperforms DAv2-metric across all metrics and altitudes. Notably, normalized error measures (AbsRel, SILog) remain relatively stable with altitude, indicating that the model exhibits robust scale-invariant behavior. In con-
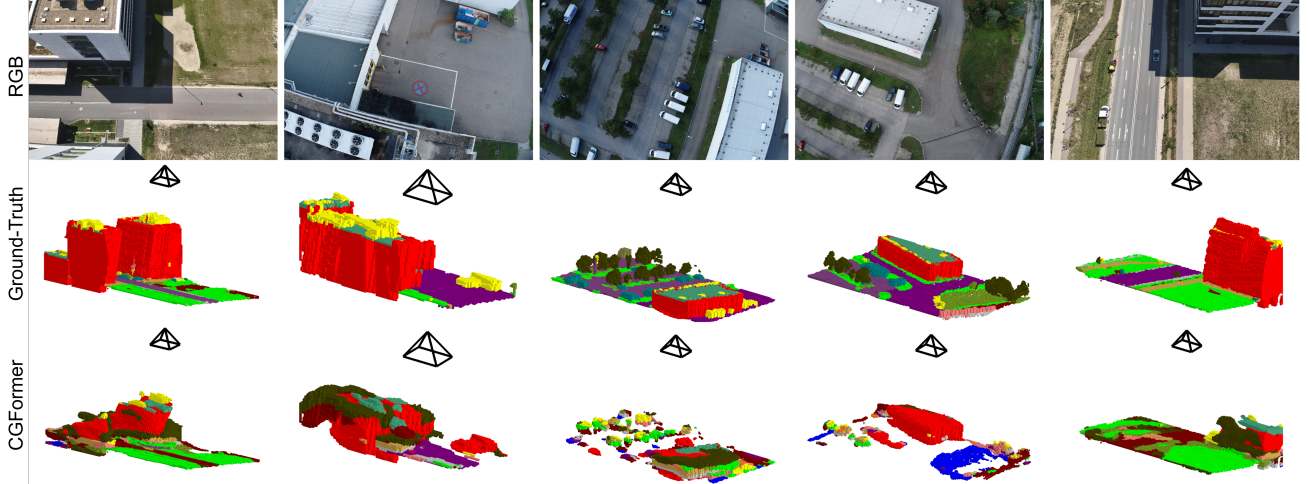
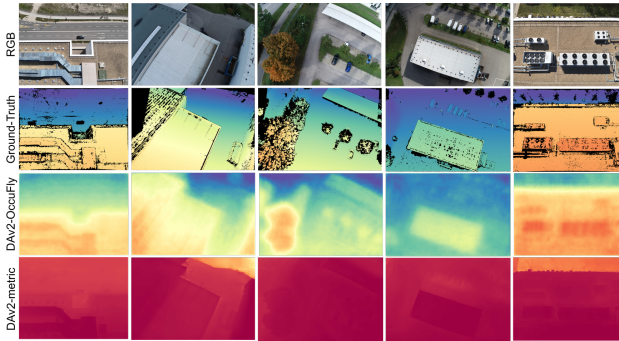Figure 4. Qualitative visualization results on the OccuFly test set (Sec. 4.2).



Figure 5. Qualitative evaluation results of DAv2-OccuFly (Sec. 5.1) on the OccuFly test set (see Sec. 4.2). Zoom in for best view.

trast, absolute errors (RMSE, MAE) increase with altitude, suggesting a positive correlation between viewpoint height and metric error. Taken together, these findings confirm that altitude materially affects metric depth estimation in aerial imagery and motivate altitude-aware training and evaluation as a promising direction for future work.

### 5.3. Qualitative Results

**Aerial Semantic Scene Completion.** Our qualitative analysis on the aerial OccuFly dataset shows that, although coarse geometry is captured, semantic consistency suffers significantly and performs inferiorly compared to terrestrial settings. This gap exposes the domain-specific challenges of aerial imagery and reveals that existing SSC models fall short in this domain. OccuFly therefore serves as a rigorous testbed to propel progress on aerial image-based 3D scene understanding.

**Metric Monocular Depth Estimation.** Qualitative depth visualizations in Fig. 5 utilize a single, metrically consistent colormap across ground truth, DAv2-OccuFly (ours), and

DAv2-metric, allowing for a direct comparison of absolute ranges. DAv2-OccuFly visually reconstitutes the ground-truth topology with coherent depth gradients and realistic color distributions, indicating well-calibrated metric estimates. In contrast, DAv2-metric often shows sharper object boundaries but is dominated by saturated red hues, evidencing a systematic overestimation of distance and poor absolute scaling. This visual mismatch mirrors the quantitative results. Overall, the qualitative evidence underscores that viewpoint altitude and in-domain fine-tuning are key to achieving accurate metric calibration.

## 6. Conclusion

OccuFly introduces the first real-world aerial 3D SSC benchmark, comprising 9 scenes and over 20,000 samples with RGB images, semantic occupancy grids, and per-frame metric depth maps across 22 semantic classes. Our LiDAR-free, image-based data data generation framework is highly scalable, requiring minimal manual annotation.

To this end, our camera-centric data generation framework faces challenges that offer further research opportunities. (1) It assumes static scenes, thus suppressing truly dynamic objects (unlike LiDAR sweep aggregation that may retain motion traces [6]). Dynamic-capable reconstruction methods, such as Dynamic NeRF [58] and 4D Gaussian Splatting [64], are a promising remedy. (2) Data aquisition may incur temporal inconsistencies across cross-altitude capture, when per-frame frustum culling uses images taken later. (3) The labeling process can be made fully automated by replacing manually annotated masks with robust 2D pseudo-labels from 2D semantic segmentation models, further scaling data generation.

Finally, our proposed dataset and underlying data generation framework foster holistic aerial 3D scene understanding.

# References

[1] Agisoft LLC. *Agisoft Metashape Professional, Version 2.2*. Agisoft LLC, 2025. Photogrammetric processing software. 1

[2] Seyed Majid Azimi, Corentin Henry, Lars Sommer, Arne Schumann, and Eleonora Vig. Skyscapes fine-grained semantic understanding of aerial scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7393–7403, 2019. 2, 6, 7

[3] Jongseong Bae, Junwoo Ha, and Ha Young Kim. Three cars approaching within 100m! enhancing distant geometry by tri-axis voxel scanning for camera-based semantic scene completion, 2025. 3

[4] Radu Beche and Sergiu Nedevschi. Claravid: A holistic scene reconstruction benchmark from aerial perspective with delentropy-based complexity profiling, 2025. 2

[5] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall. SemanticKITTI: A dataset for semantic scene understanding of LiDAR sequences. In *Proc. of the IEEE/CVF International Conf. on Computer Vision (ICCV)*, 2019. 2

[6] Jens Behley, Martin Garbade, Andres Milioto, Jan Quenzel, Sven Behnke, Cyrill Stachniss, and Jurgen Gall. Semantickitti: A dataset for semantic scene understanding of lidar sequences. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9297–9307, 2019. 3, 5, 6, 8, 1

[7] Yohann Cabon, Naila Murray, and Martin Humenberger. Virtual kitti 2, 2020. 3

[8] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11621–11631, 2020. 2

[9] Wenxiao Cai, Ke Jin, Jinyan Hou, Cong Guo, Letian Wu, and Wankou Yang. Vdd: Varied drone dataset for semantic segmentation. *Journal of Visual Communication and Image Representation*, 109:104429, 2025. 2, 6

[10] Anh-Quan Cao and Raoul de Charette. MonoScene: Monocular 3D semantic scene completion. In *CVPR*, 2022. 2, 3

[11] Dubing Chen, Huan Zheng, Jin Fang, Xingping Dong, Xianfei Li, Wenlong Liao, Tao He, Pai Peng, and Jianbing Shen. Rethinking temporal fusion with a unified gradient descent view for 3d semantic occupancy prediction. In *2025 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1505–1515, 2025. 3

[12] Lyujie Chen, Feng Liu, Yan Zhao, Wufan Wang, Xiaming Yuan, and Jihong Zhu. Valid: A comprehensive virtual aerial image dataset. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2009–2016, 2020. 2

[13] Yu Chen, Yao Wang, Peng Lu, Yisong Chen, and Guoping Wang. Large-scale structure from motion with semantic constraints of aerial images. In *Chinese Conference on Pattern Recognition and Computer Vision (PRCV)*, pages 347–359. Springer, 2018. 2, 6

[14] OpenScene Contributors. Openscene: The largest up-to-date 3d occupancy prediction benchmark in autonomous driving, 2023. 3

[15] Oussema Dhaouadi, Johannes Meier, Luca Wahl, Jacques Kaiser, Luca Scalerandi, Nick Wandelburg, Zhuolun Zhuo, Nijanthan Berinpanathan, Holger Banzhaf, and Daniel Cremers. Highly accurate and diverse traffic data: The deepscenario open 3d dataset. In *2025 IEEE Intelligent Vehicles Symposium*. IEEE, 2025. 2

[16] DJI. Phantom 4 rtk. `https://www.dji.com/phantom-4-rtk/info`, 2016. Accessed: 2025-11-08. 6, 2

[17] DJI. Mavic 3 enterprise series. `https://enterprise.dji.com/mavic-3-enterprise`, 2022. Accessed: 2025-11-08. 6, 2

[18] DJI. Onboard SDK documentation. `https://developer.dji.com/onboard-api-reference/group__telem.html`, 2025. Accessed: 2025-11-08. 6

[19] Aloisio Dourado, Teofilo E. De Campos, Hansung Kim, and Adrian Hilton. Edgenet: Semantic scene completion from a single rgb- d image. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pages 503–510, 2021. 2

[20] Herbert Edelsbrunner, David Kirkpatrick, and Raimund Seidel. On the shape of a set of points in the plane. *IEEE Transactions on Information Theory*, 29(4):551–559, 1983. 5, 1

[21] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, page 226–231. AAAI Press, 1996. 4, 1

[22] European Commission. Commission implementing regulation (eu) 2019/947 of 24 may 2019 on the rules and procedures for the operation of unmanned aircraft, 2019. `https://eur-lex.europa.eu/eli/reg_impl/2019/947/oj`. 2

[23] Olivier Faugeras. *Three-Dimensional Computer Vision: A Geometric Viewpoint*. MIT Press, 1993. Section 2.3.2: Voxel Grids. 4

[24] Federal Aviation Administration. Small unmanned aircraft systems, 2016. `https://www.ecfr.gov/current/title-14/chapter-I/subchapter-F/part-107`. 2

[25] Horatiu Florea, Vlad-Cristian Miclea, and Sergiu Nedevschi. Wilduav: Monocular uav dataset for depth estimation tasks. *2021 IEEE 17th International Conference on Intelligent Computer Communication and Processing (ICCP)*, 2021. 2, 3

[26] Michael Fonder and Marc Van Droogenbroeck. Mid-air: A multi-modal dataset for extremely low altitude drone flights. In *Conference on Computer Vision and Pattern Recognition Workshop (CVPRW)*, 2019. 2

[27] Markus Gerke and Heinz-J"urgen Przybilla. Accuracy analysis of photogrammetric uav image blocks: Influence of onboard rtk-gnss and cross flight patterns. *Photogrammetrie - Fernerkundung - Geoinformation*, 2016(1):17–30, 2016. 6

[28] Markus Gross, Aya Fahmy, Danit Niwattananan, Dominik Muhle, Rui Song, Daniel Cremers, and Henri Meeß. IP-Former: Visual 3d panoptic scene completion with context-

adaptive instance proposals. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. 2

[29] Xiyue Guo, Jiarui Hu, Junjie Hu, Hujun Bao, and Guofeng Zhang. Sgformer: Satellite-ground fusion for 3d semantic scene completion. *arXiv preprint arXiv:2503.16825*, 2025. 3

[30] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003. 6

[31] Adrian Hayler, Felix Wimbauer, Dominik Muhle, Christian Rupprecht, and Daniel Cremers. S4c: Self-supervised semantic scene completion with neural fields. In *2024 International Conference on 3D Vision (3DV)*, pages 409–420. IEEE, 2024. 3

[32] Yuanhui Huang, Wenzhao Zheng, Yunpeng Zhang, Jie Zhou, and Jiwen Lu. Tri-perspective view for vision-based 3d semantic occupancy prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9223–9232, 2023. 3, 7

[33] Yuanhui Huang, Amonnut Thammatadatrakoon, Wenzhao Zheng, Yunpeng Zhang, Dalong Du, and Jiwen Lu. Gaussianformer-2: Probabilistic gaussian superposition for efficient 3d occupancy prediction. In *2025 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 27477–27486, 2025. 3

[34] Institute of Computer Graphics and Vision, Graz University of Technology. Semantic drone dataset, 2019. 2, 6

[35] Aleksandar Jevtić, Christoph Reich, Felix Wimbauer, Oliver Hahn, Christian Rupprecht, Stefan Roth, and Daniel Cremers. Feed-forward SceneDINO for unsupervised semantic scene completion. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2025. 3

[36] Haoyi Jiang, Tianheng Cheng, Naiyu Gao, Haoyang Zhang, Tianwei Lin, Wenyu Liu, and Xinggang Wang. Symphonize 3d semantic scene completion with contextual instance queries. *CVPR*, 2024. 3, 7

[37] Arie Kaufman and Eyal Shimony. 3d scan-conversion algorithms for voxel-based graphics. In *Proceedings of the 1986 Workshop on Interactive 3D Graphics*, pages 45–75, Chapel Hill, NC, USA, 1986. 4

[38] Michael M. Kazhdan, Matthew Bolitho, and Hugues Hoppe. Poisson surface reconstruction. In *Proceedings of the Fourth Eurographics Symposium on Geometry Processing*, pages 61–70, Aire-la-Ville, Switzerland, Switzerland, 2006. Eurographics Association. 5, 1

[39] Benedikt Kolbeinsson and Krystian Mikolajczyk. DDOS: The drone depth and obstacle segmentation dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2024. 2

[40] A. Laurentini. The visual hull concept for silhouette-based image understanding. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(2):150–162, 1994. 5

[41] Hyo-Jun Lee, Yeong Jun Koh, Hanul Kim, Hyunseop Kim, Yonguk Lee, and Jinu Lee. Soap: Vision-centric 3d semantic scene completion with scene-adaptive decoder and occluded region-aware view projection. In *2025 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 17145–17154, 2025. 3

[42] Heng Li, Yuenan Hou, Xiaohan Xing, Yuexin Ma, Xiao Sun, and Yanyong Zhang. Occmamba: Semantic occupancy prediction with state space models. In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR)*, pages 11949–11959, 2025. 3

[43] Shijie Li, Zhongyao Cheng, Rong Li, Shuai Li, Juergen Gall, Xun Xu, and Xulei Yang. Global-aware monocular semantic scene completion with state space models, 2025. 3

[44] Yixuan Li, Lihan Jiang, Linning Xu, Yuanbo Xiangli, Zhenzhi Wang, Dahua Lin, and Bo Dai. Matrixcity: A large-scale city dataset for city-scale neural rendering and beyond. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3205–3215, 2023. 2

[45] Yiming Li, Zhiding Yu, Christopher Choy, Chaowei Xiao, Jose M Alvarez, Sanja Fidler, Chen Feng, and Anima Anandkumar. Voxformer: Sparse voxel transformer for camera-based 3d semantic scene completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9087–9098, 2023. 3, 7

[46] Yiming Li, Sihang Li, Xinhao Liu, Moonjun Gong, Kenan Li, Nuo Chen, Zijun Wang, Zhiheng Li, Tao Jiang, Fisher Yu, Yue Wang, Hang Zhao, Zhiding Yu, and Chen Feng. Sscbench: A large-scale 3d semantic scene completion benchmark for autonomous driving. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2024. 2, 3, 6, 1

[47] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Qiao Yu, and Jifeng Dai. Bevformer: Learning bird's-eye-view representation from multi-camera images via spatiotemporal transformers. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2022. 3

[48] Yiyi Liao, Jun Xie, and Andreas Geiger. Kitti-360: A novel dataset and benchmarks for urban scene understanding in 2d and 3d. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(3):3292–3310, 2023. 2

[49] Lin Liqiang, Liu Yilin, Hu Yue, Yan Xingguang, Xie Ke, and Huang Hui. Capturing, reconstructing, and simulating: the urbanscene3d dataset. In *ECCV*, pages 93–109, 2022. 2

[50] Enyu Liu, En Yu, Sijia Chen, and Wenbing Tao. Disentangling instance and scene contexts for 3d semantic scene completion, 2025. 3

[51] Xinhao Liu, Moonjun Gong, Qi Fang, Haoyu Xie, Yiming Li, Hang Zhao, and Chen Feng. Lidar-based 4d occupancy completion and forecasting. *arXiv preprint arXiv:2310.11239*, 2023. 2

[52] Rafael Lopez-Campos and Jose Martinez-Carranza. Espada: Extended synthetic and photogrammetric aerial-image dataset. *IEEE Robotics and Automation Letters*, 6(4):7981–7988, 2021. 2

[53] Haoang Lu, Yuanqi Su, Xiaoning Zhang, Longjun Gao, Yu Xue, and Le Wang. Vishall3d: Monocular semantic scene completion from reconstructing the visible regions to hallucinating the invisible regions, 2025. 3

[54] Ye Lyu, George Vosselman, Gui-Song Xia, Alper Yilmaz, and Michael Ying Yang. Uavid: A semantic segmentation dataset for uav imagery. *ISPRS journal of photogrammetry and remote sensing*, 165:108–119, 2020. 2, 6

[55] Alina Marcu, Mihai Pirvu, Dragos Costea, Emanuela Haller, Emil Slusanschi, Ahmed Nabil Belbachir, Rahul Sukthankar, and Marius Leordeanu. Self-supervised hypergraphs for learning multiple world interpretations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 983–992, 2023. 2

[56] F. Nex, E. K. Stathopoulou, F. Remondino, M. Y. Yang, L. Madhuanand, Y. Yogender, B. Alsadik, M. Weinmann, B. Jutzi, and R. Qin. Usegeo - a uav-based multi-sensor dataset for geospatial research. *ISPRS Open Journal of Photogrammetry and Remote Sensing*, 13:100070, 2024. 2, 3

[57] Ishan Nigam, Chen Huang, and Deva Ramanan. Ensemble knowledge transfer for semantic segmentation. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1499–1508. IEEE, 2018. 2, 6

[58] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10313–10322, 2021. 8

[59] Giulia Rizzoli, Francesco Barbato, Matteo Caligiuri, and Pietro Zanuttigh. Syndrone-multi-modal uav dataset for urban scenarios. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2210–2220, 2023. 2

[60] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 4

[61] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 4

[62] Faranak Shamsafar, Samuel Woerz, Rafia Rahim, and Andreas Zell. Mobilestereonet: Towards lightweight deep networks for stereo matching. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2417–2426, 2022. 7

[63] Rui Song, Chenwei Liang, Hu Cao, Zhiran Yan, Walter Zimmer, Markus Gross, Andreas Festag, and Alois Knoll. Collaborative semantic occupancy prediction with hybrid feature fusion in connected automated vehicles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 17996–18006, 2024. 3

[64] Rui Song, Chenwei Liang, Yan Xia, Walter Zimmer, Hu Cao, Holger Caesar, Andreas Festag, and Alois Knoll. Coda-4dgs: Dynamic gaussian splatting with context and deformation awareness for autonomous driving. In *IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE/CVF, 2025. 8

[65] Shuran Song, Fisher Yu, Andy Zeng, Angel X. Chang, Manolis Savva, and Thomas Funkhouser. Semantic scene completion from a single depth image. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 190–198, 2017. 2

[66] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott Ettinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[67] Sebastian Thrun, Wolfram Burgard, and Dieter Fox. *Probabilistic robotics*. MIT Press, Cambridge, Mass., 2005. 2

[68] Pengju Tian, Zhirui Wang, Peirui Cheng, Yuchao Wang, Zhechao Wang, Liangjin Zhao, Menglong Yan, Xue Yang, and Xian Sun. Ucdnet: Multi-uav collaborative 3-d object detection network by reliable feature mapping. *IEEE Transactions on Geoscience and Remote Sensing*, 63:1–16, 2025. 2

[69] Xiaoyu Tian, Tao Jiang, Longfei Yun, Yue Wang, Yilun Wang, and Hang Zhao. Occ3d: A large-scale 3d occupancy prediction benchmark for autonomous driving. *arXiv preprint arXiv:2304.14365*, 2023. 2, 4, 6

[70] Khiem Vuong, Anurag Ghosh, Deva Ramanan, Srinivasa Narasimhan, and Shubham Tulsiani. Aerialmegadepth: Learning aerial-ground reconstruction and view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2025. 2

[71] Jiahao Wang, Xiangyu Cao, Jiaru Zhong, Yuner Zhang, Haibao Yu, Lei He, and Shaobing Xu. Griffin: Aerial-ground cooperative detection and tracking dataset and benchmark, 2025. 2

[72] Wenshan Wang, Delong Zhu, Xiangwei Wang, Yaoyu Hu, Yuheng Qiu, Chen Wang, Yafei Hu, Ashish Kapoor, and Sebastian Scherer. Tartanair: A dataset to push the limits of visual slam. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020. 2

[73] Xiaofeng Wang, Zheng Zhu, Wenbo Xu, Yunpeng Zhang, Yi Wei, Xu Chi, Yun Ye, Dalong Du, Jiwen Lu, and Xingang Wang. Openoccupancy: A large scale benchmark for surrounding semantic occupancy perception. *arXiv preprint arXiv:2303.03991*, 2023. 2, 6

[74] Xuzhi Wang, Xinran Wu, Song Wang, Lingdong Kong, and Ziping Zhao. Monocular semantic scene completion via masked recurrent networks. In *Proceedings of the IEEE/CVF Conference on International Conference on Computer Vision (ICCV)*, 2025. 3

[75] Yuping Wang, Xiangyu Huang, Xiaokang Sun, Mingxuan Yan, Shuo Xing, Zhengzhong Tu, and Jiachen Li. Uniocc: A unified benchmark for occupancy forecasting and prediction in autonomous driving. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, 2025. 3

[76] Z. Wang, P. Cheng, M. Chen, P. Tian, Z. Wang, X. Li, X. Yang, and X. Sun. Drones help drones: A collaborative framework for multi-drone object trajectory prediction and beyond. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 2

[77] Yi Wei, Linqing Zhao, Wenzhao Zheng, Zheng Zhu, Jie Zhou, and Jiwen Lu. Surroundocc: Multi-camera 3d occupancy prediction for autonomous driving. *arXiv preprint arXiv:2303.09551*, 2023. 3, 5

[78] Rouwan Wu, Xiaoya Cheng, Juelin Zhu, Xuxiang Liu, Maojun Zhang, and Shen Yan. Uavd4l: A large-scale dataset for

uav 6-dof localization. In *International Conference on 3D Vision (3DV)*, 2024. 2

[79] Huaiyuan Xu, Junliang Chen, Shiyu Meng, Yi Wang, and Lap-Pui Chau. A survey on occupancy perception for autonomous driving: The information fusion perspective. *Information Fusion*, 114:102671, 2025. 3

[80] Runsheng Xu, Hao Xiang, Xin Xia, Xu Han, Jinlong Li, and Jiaqi Ma. Opv2v: An open benchmark dataset and fusion pipeline for perception with vehicle-to-vehicle communication. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 2583–2589, 2022. 3

[81] Qi Yan, Jianhao Zheng, Simon Reding, Shanci Li, and Iordan Doytchinov. Crossloc: Scalable aerial localization assisted by multimodal synthetic data. *arXiv preprint arXiv:2112.09081*, 2021. 2

[82] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *arXiv preprint arXiv:2406.09414*, 2024. 2, 7, 1

[83] Hui Ye, Raj Sunderraman, and Shihao Ji. Uav3d: A large-scale 3d perception benchmark for unmanned aerial vehicles. In *The 38th Conference on Neural Information Processing Systems (NeurIPS)*, 2024. 2

[84] Zhu Yu, Runmin Zhang, Jiacheng Ying, Junchen Yu, Xiaohai Hu, Lun Luo, Si-Yuan Cao, and Hui-Liang Shen. Context and geometry aware voxel transformer for semantic scene completion. In *Advances in Neural Information Processing Systems*, 2024. 3, 7, 1, 4

[85] Yunpeng Zhang, Zheng Zhu, and Dalong Du. Occformer: Dual-path transformer for vision-based 3d semantic occupancy prediction. *arXiv preprint arXiv:2304.05316*, 2023. 7

[86] Yanan Zhang, Jinqing Zhang, Zengran Wang, Junhao Xu, and Di Huang. Vision-based 3d occupancy prediction in autonomous driving: a review and outlook. *arXiv preprint arXiv:2405.02595*, 2024. 2, 3

[87] Pengfei Zhu, Longyin Wen, Dawei Du, Xiao Bian, Heng Fan, Qinghua Hu, and Haibin Ling. Detection and tracking meet drones challenge. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021. 2

# OccuFly: A 3D Vision Benchmark for Semantic Scene Completion from the Aerial Perspective
## Supplementary Material

## 7. Implementation Details

**Data Generation.** We generate data on an AMD Ryzen Threadripper PRO 7985WX 64-Cores (allocating 8 cores) with $120\,\mathrm{GB}$ of memory.

For **3D reconstruction** (Sec. 3.3.1), we utilize the Agisoft Metashape 2.2.0 photogrammetric reconstruction software [1]. After reconstruction, we ensure high geometric fidelity across the whole scene by removing a small, noisy margin at the border of the scene, which naturally arises from aerial SfM+MVS due to decreased image overlap (Fig. 9). Subsequently, during **semantic annotation** (Sec. 3.3.2), we partition the ground plane into a regular grid with square cells of $25\,\mathrm{m}$ edge length to determine the subset of frames for manual annotation. Moreover, we apply kNN with $k = 100$ for unlabelled point assignment, and $k = 200$ for subsequent label refinement. Furthermore, for **densification and voxelization** (Sec. 3.3.3) of *instance classes*, DBSCAN [21] clustering is performed with class-wise parameters detailed in Tab. 5, and we set $\alpha = 0.05$ for $\alpha$-Shapes [20] and use $K = 24$ camera views during silhouette extraction. For *ground classes*, we set Poisson reconstruction parameters [38] to a depth of 8 and a scale of 1.2. Finally, group assignments for all classes are reported in Tab. 6.

**Aerial Semantic Scene Completion.** Following the training protocol of CGFormer [84], we train this benchmark method for 25 epochs, using the official implementation from GitHub. We adapt the codebase to OccuFly's voxel grid of $192 \times 128 \times 128$, replacing the $256 \times 256 \times 32$ grid used in SemanticKITTI [6] and SSCBench [46], on which the original model was trained. All experiments are conducted on a single NVIDIA A100 80GB GPU with a batch size of 1. The peak GPU memory usage is $32.39\,\mathrm{GB}$.

Table 5. Class-wise DBSCAN [21] parameters for instance separation, discussed in Sec. 3.3.3.

| Class | $\epsilon$ | MinPts |
|---|---|---|
| Building | 4.0 | 1000 |
| Roof | 1.0 | 1000 |
| Vehicle | 1.0 | 500 |
| Crane | 1.0 | 500 |
| Bicycle | 0.4 | 80 |
| Person | 0.3 | 10 |
| Flying Animal | 0.3 | 30 |
| Truck | 1.0 | 500 |

Table 6. Semantic class frequencies, group assignments (Sec. 3.3.3), and semantic color table of the OccuFly dataset.

| Group | Color | Name | % |
|---|---|---|---|
| Instance | ■ | Building | 75.7457 |
| | ■ | Roof | 1.7417 |
| | ■ | Vehicle | 0.5195 |
| | ■ | Crane | 0.0052 |
| | ■ | Bicycle | 0.0046 |
| | ■ | Person | 0.0002 |
| | ■ | Flying Animal | 0.0001 |
| | ■ | Truck | 0.1067 |
| Ground | ■ | Grass | 4.1978 |
| | ■ | Vegetation | 2.3234 |
| | ■ | Water | 1.1322 |
| | ■ | Walkway | 1.0560 |
| | ■ | Dirt | 1.2102 |
| | ■ | Road | 0.9377 |
| | ■ | Gravel | 0.7960 |
| | ■ | Parking Lot | 1.4349 |
| Others | ■ | Tree | 6.8357 |
| | ■ | Ground Obstacle | 1.7566 |
| | ■ | Construction | 0.1689 |
| | ■ | Cable Tower | 0.0040 |
| | ■ | Rock | 0.0210 |
| | ■ | Cable | 0.0017 |

**Metric Monocular Depth Estimation.** We employ Depth Anything V2 Vi-T-Small [82] with $24.8\,\mathrm{M}$ parameters for metric monocular depth estimation. Technically, we use the authors' official implementation hosted on Hugging Face and run training on a single NVIDIA A100 80 GB GPU with a batch size of 12. Under this configuration, the peak GPU memory footprint is $19.5\,\mathrm{GB}$.

## 8. Additional OccuFly Dataset Evaluation

**3D Reconstruction.** We provide scene-wise reprojection errors in Tab. 7. An average root mean square reprojection error of $1.24$ pixels in our geo-referenced images validates the high metric accuracy of the reconstructed point cloud. Scene-wise reconstructed point clouds are shown in Fig. 10.

**Semantic Annotation.** As detailed in Sec. 3.3.2, we manually annotate only a small subset of images and subsequently lift semantic labels to 3D. In Fig. 6, we present qualitative examples of the manual annotations, which exhibit exceptional pixel-accurate delineation. Furthermore, Tab. 8 reports

Table 7. Scene-wise root mean square (RMS) reprojection error after 3D reconstruction (Sec. 3.3.1).

| Scene | RMS Reprojection Error [px] |
|-------|------------------------------|
| 1 | 0.469 |
| 2 | 0.474 |
| 3 | 0.388 |
| 4 | 0.451 |
| 5 | 0.422 |
| 6 | 2.13 |
| 7 | 2.04 |
| 8 | 2.61 |
| 9 | 2.22 |
| Average | 1.24 |

Table 8. Scene-wise manual semantic annotation ratios for UAV platforms DJI Phantom 4 RTK (P4) [16] and DJI Mavic 3 Enterprise Series (M3-ES) [17]. Note that the number of aquired images marginally differs from the number of images finally provided in the dataset, as we remove images at the border of each reconstructed scene to ensure high geometric fidelity (see Sec. 7).

| Scene | UAV Paltform | Aquired Images | Annotated Images | Ratio [%] |
|-------|--------------|----------------|------------------|-----------|
| 1 | P4 | 421 | 73 | 17.34 |
| 2 | P4 | 338 | 48 | 14.20 |
| 3 | M3-ES | 1048 | 66 | 6.30 |
| 4 | M3-ES | 1252 | 102 | 8.15 |
| 5 | M3-ES | 1082 | 74 | 6.84 |
| 6 | P4 | 380 | 52 | 13.68 |
| 7 | P4 | 284 | 40 | 14.08 |
| 8 | P4 | 251 | 38 | 15.14 |
| 9 | M3-ES | 1337 | 93 | 6.96 |
| **Total** | | **6393** | **586** | **9.17** |

per-scene annotation ratios, achieving an average annotation ratio of $<10\%$, indicating exceptional annotation efficiency. Additionally, Fig. 9 shows scene-wise image overlap during data collection, which exceeds $>90\%$ for all scenes. This substantial overlap ensures accurate 3D reconstruction and semantic label lifting. Finally, Fig. 10 illustrates (i) the reconstructed RGB point cloud, (ii) the semantic point cloud after label lifting, and (iii) the resulting scene-level semantic voxel grid. These visualizations demonstrate the remarkable fidelity of our data generation framework and its ability to propagate sparse 2D annotations to a globally consistent, voxel-level 3D ground truth.

**Metric Depth Maps.** Table 7 reports a mean reprojection error of $1.24$ pixels, indicating high geometric consistency of the reconstruction. Since the metric depth maps are derived from these reconstructed points, a low reprojection error serves as a strong proxy for depth accuracy. Moreover, in Tab. 9, we compare OccuFly to other real-world, low-
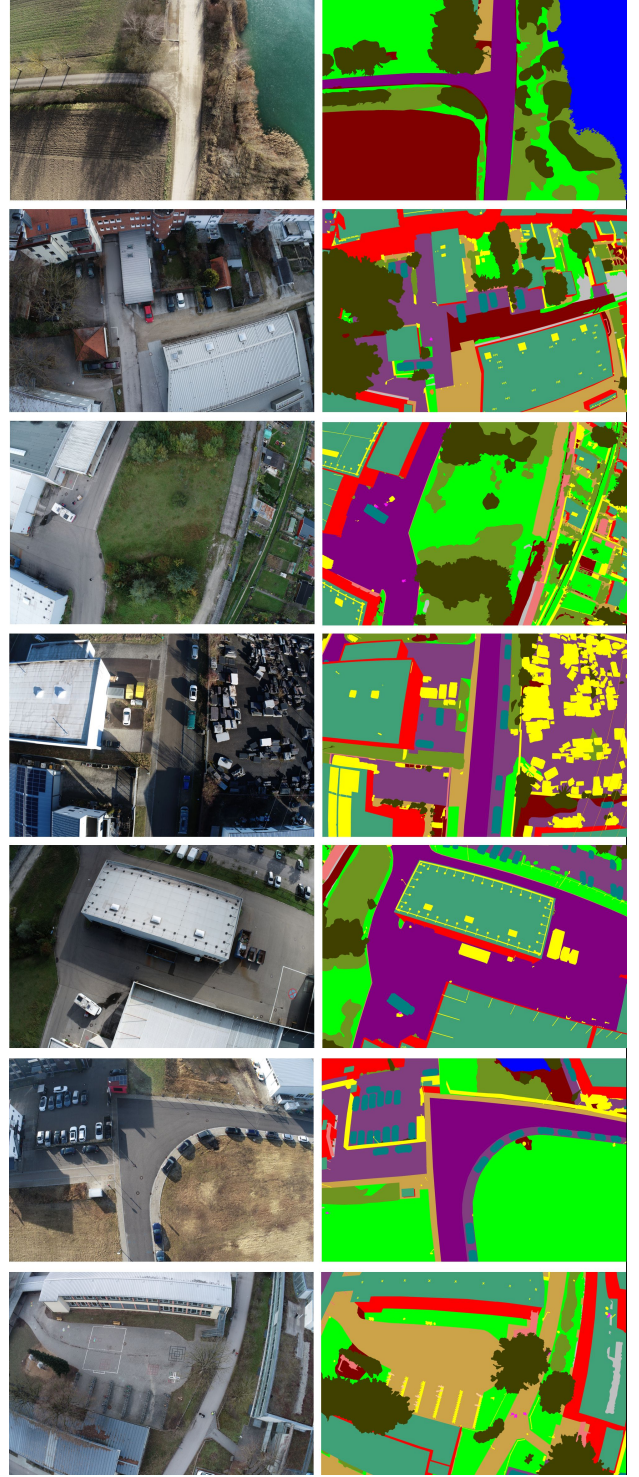


Figure 6. Manual image annotations used for 3D label lifting, showing exceptional pixel-accurate delineation. Zoom in for best view.

altitude aerial datasets that include metric depth maps. To the best of our knowledge, WildUAV [25] and UseGeo [56]

Table 9. Comparison of the OccuFly dataset to other real-world low-altitude aerial datasets containing metric depth maps.

| Dataset | # Depth Maps | Density | Scenarios | Seasons |
|---|---|---|---|---|
| WildUAV [25] | $\sim 1,500$ | Dense | Rural | Summer, Autumn |
| UseGeo [56] | 829 | Sparse | Urban | n.a. |
| OccuFly (ours) | 20,160 | Dense | Rural, Urban, Industrial | Spring, Summer, Autumn, Winter |

Table 10. Reciprocal, out-of-domain, zero-shot generalization on metric monocular depth estimation, discussed in Sec. 9.2.

| Method | Higher is better ↑ | | | Lower is better ↓ | | | |
|---|---|---|---|---|---|---|---|
| | $\delta_1$ | $\delta_2$ | $\delta_3$ | AbsRel | RMSE | MAE | SILog |
| DAv2-metric on OccuFly | 0.002 | 0.015 | 0.059 | **0.729** | 28.369 | 27.381 | **0.192** |
| DAv2-OccuFly (ours) on V-KITTI-2 [7] | **0.122** | **0.337** | **0.521** | 1.098 | **16.711** | **15.251** | 0.686 |

are the only publicly available dataset of this kind. OccuFly is substantially larger, providing more than $13\times$ and more than $24\times$ as many metric depth maps, respectively, while spanning a broader range of scenarios and seasons. This positions OccuFly as the largest and most diverse publicly available low-altitude metric depth estimation dataset to date, which further enables holistic vision-based 3D scene understanding, such as Aerial Semantic Scene Completion.

Finally, we provide per-scene depth histograms for all nine scenes to illustrate the dataset's metric depth distributions (see Fig. 7). Most scenes show peaks around 30–50 meters, reflecting the image acquisition altitudes, while certain scenes, such as Scene 09, exhibit more diverse depth ranges.

# 9. Additional Benchmark Evaluation

## 9.1. Aerial Semantic Scene Completion

We provide class-wise evaluation metrics of CGFormer [84] in Tab. 11, and additional qualitative results in Fig. 8. While overall performance is limited, the model attains its overall peak accuracy (IoU and mIoU) at $30\,\mathrm{m}$ and performs comparatively well at $50\,\mathrm{m}$. This pattern indicates that altitude-related viewpoint variations influence the performance of state-of-the-art SSC methods. In line with our previous findings, the performance on our aerial OccuFly dataset indicates that, although the network captures coarse 3D structure, its predictions remain semantically fragmented and inferior to those obtained in terrestrial scenarios. These discrepancies highlight the unique challenges introduced by aerial viewpoints and emphasize that existing SSC architectures are not yet optimized for this regime. Consequently, our dataset
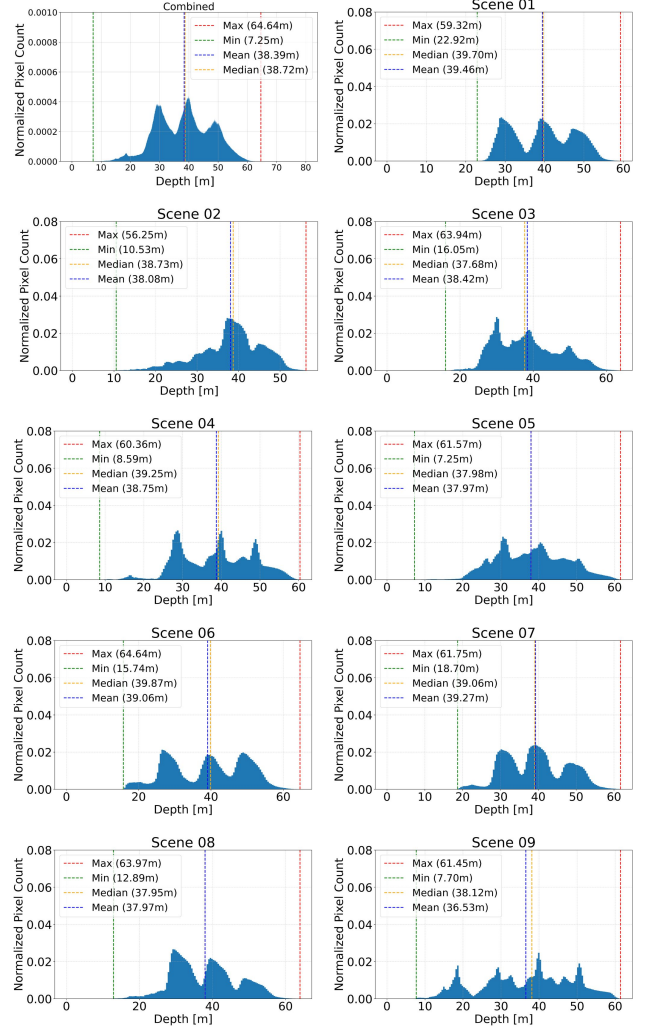


Figure 7. Depth map histograms for each of the 9 scenes in the OccuFly dataset. Zoom in for best view.

serves as an important contribution, providing a benchmark that can catalyze future research on robust SSC from aerial perspectives.

## 9.2. Metric Monocular Depth Estimation

In addition to the altitude-wise evaluation in Tab. 4, we design a reciprocal, cross-domain, zero-shot generalization experiment to quantify transfer between terrestrial (driving) and aerial domains. Results are reported in Tab. 10. Specifically, we evaluate DAv2-metric (see Sec. 5.1), fine-tuned on terrestrial Virtual KITTI 2 (V-KITTI-2) [7], on the OccuFly test set, and we evaluate our DAv2-OccuFly, fine-tuned on aerial OccuFly, on the V-KITTI-2 test set. Our model outperforms the baseline in metric-sensitive measures, indicating more accurate absolute depth predictions. In contrast, the baseline achieves better performance in scale-invariant metrics, which reflects higher fidelity in relative scene ge-

Table 11. Per-altitude and class-wise SSC performance of CGFormer [84] on the OccuFly test set, including semantic class frequencies. Best and second-best results are bold and underlined, respectively.

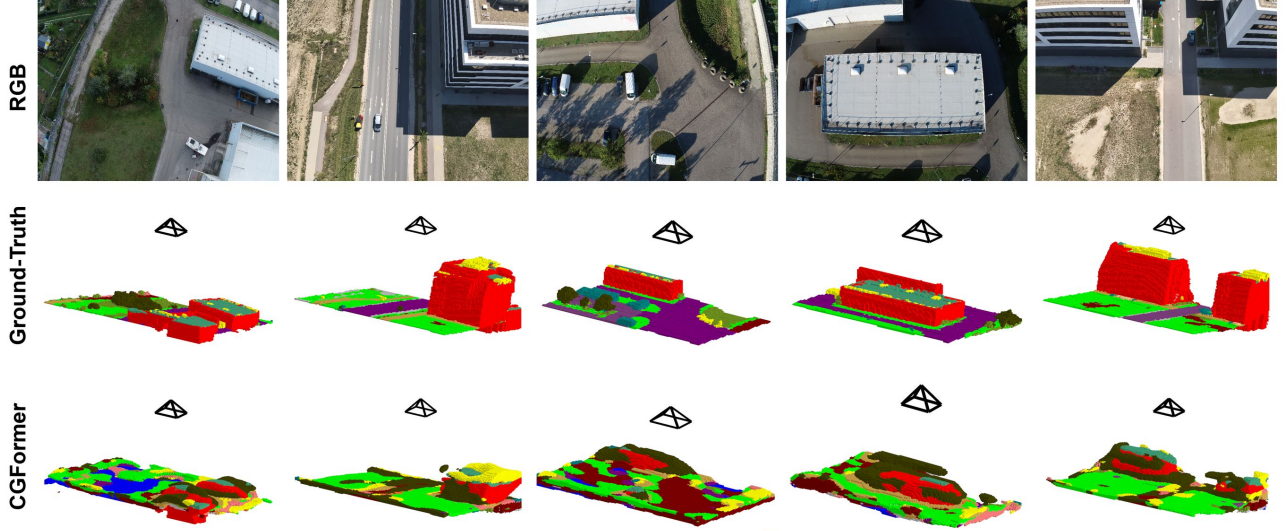| Altitude [m] | Road (0.9377%) | Walkway (1.0560%) | Dirt (1.2102%) | Gravel (0.7960%) | Rock (0.0210%) | Grass (4.1978%) | Vegetation (2.3234%) | Tree (6.8357%) | Ground-Obs. (1.7566%) | Person (0.0002%) | Bicycle (0.0046%) | Vehicle (0.5195%) | Water (1.1322%) | Building (73.7457%) | Roof (1.7417%) | Cables (0.0017%) | Cable-Tower (0.0040%) | Flying-Animals (0.0001%) | Parking-Lot (1.4349%) | Constructions (0.1689%) | Cranes (0.0052%) | Truck (0.1067%) | mIoU | IoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 30 | **0.0387** | **0.0122** | <u>0.0017</u> | 0.0016 | 0.0001 | <u>0.0346</u> | <u>0.0047</u> | **0.0054** | **0.0098** | 0.0000 | 0.0000 | **0.0025** | 0.0000 | **0.0959** | **0.0068** | 0.0000 | 0.0000 | 0.0000 | **0.0001** | 0.0005 | 0.0000 | 0.0000 | **0.0097** | **0.2948** |
| 40 | <u>0.0221</u> | <u>0.0055</u> | 0.0009 | <u>0.0023</u> | <u>0.0002</u> | 0.0103 | 0.0035 | 0.0024 | 0.0009 | 0.0000 | 0.0000 | <u>0.0018</u> | 0.0000 | 0.0767 | <u>0.0017</u> | 0.0000 | 0.0000 | 0.0000 | 0.0000 | <u>0.0006</u> | 0.0000 | 0.0000 | 0.0059 | 0.2082 |
| 50 | 0.0028 | <u>0.0055</u> | **0.0031** | **0.0044** | **0.0005** | **0.0851** | **0.0092** | <u>0.0028</u> | <u>0.0041</u> | 0.0000 | 0.0000 | 0.0002 | 0.0000 | <u>0.0790</u> | 0.0016 | 0.0000 | 0.0000 | 0.0000 | **0.0001** | **0.0015** | 0.0000 | 0.0000 | <u>0.0091</u> | <u>0.2600</u> |
| all | 0.0143 | 0.0061 | 0.0024 | 0.0035 | 0.0004 | 0.0535 | 0.0076 | 0.0028 | 0.0029 | 0.0000 | 0.0000 | 0.0011 | 0.0000 | 0.0788 | 0.0019 | 0.0000 | 0.0000 | 0.0000 | 0.0001 | 0.0013 | 0.0000 | 0.0000 | 0.0080 | 0.2406 |



Figure 8. Additional qualitative visualization results of CGFormer [84] on the OccuFly test set.

ometry. In terms of metric depth estimation, these results show that our proposed DAv2-OccuFly model generalizes better to the driving domain than DAv2-metric generalizes to the aerial domain. Since downstream tasks, such as Semantic Scene Completion, require accurate metric depth, our model is therefore better suited to these scenarios, providing a more reliable absolute scale. Taken together, these findings highlight OccuFly's value as a benchmark for cross-domain metric monocular depth estimation and downstream SSC.

**Semantic Point Cloud** | **Image Overlap** | **Semantic Point Cloud** | **Image Overlap**

**Overlap Color Table:**

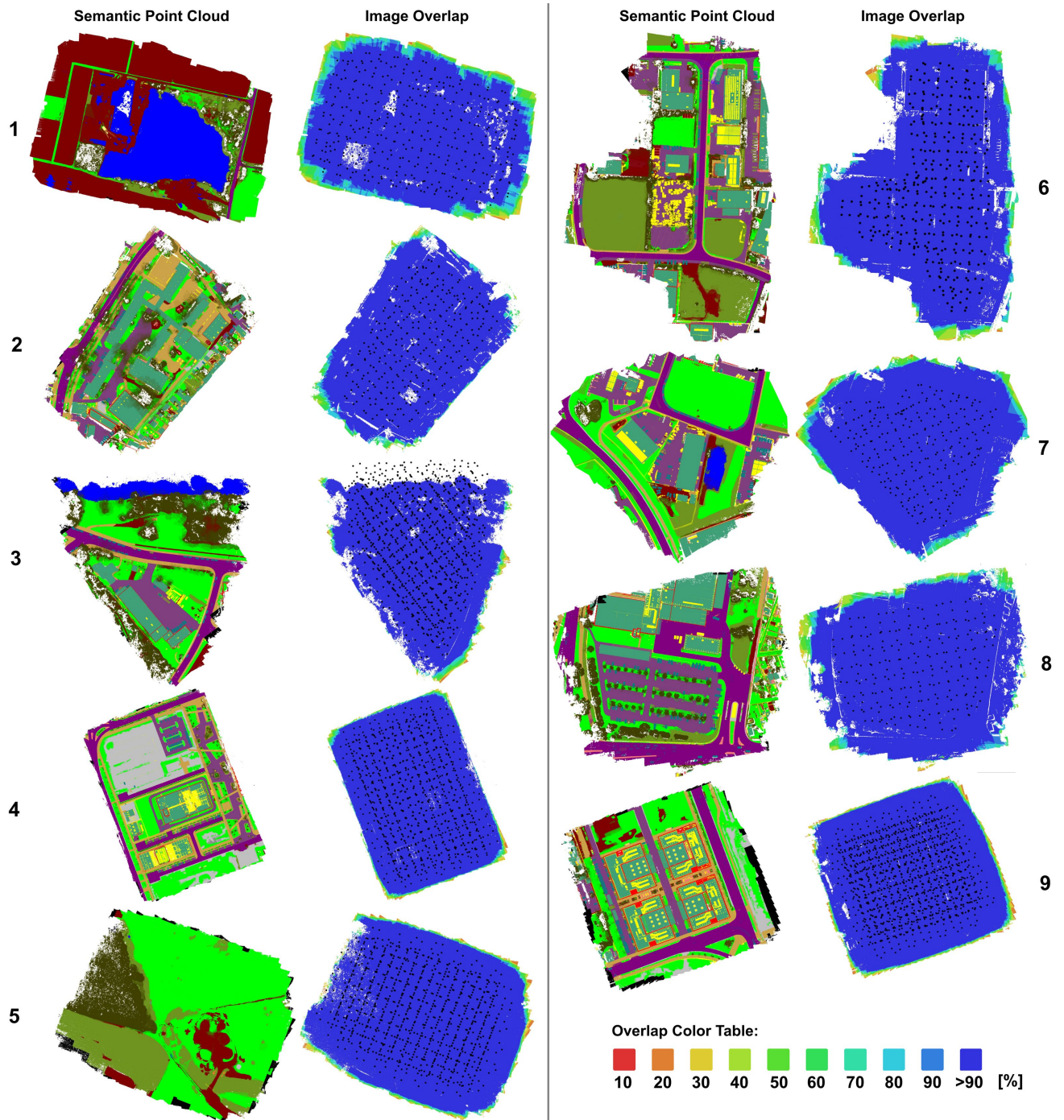| 10 | 20 | 30 | 40 | 50 | 60 | 70 | 80 | 90 | >90 | [%] |

Figure 9. Scene-wise image overlap during data collection for all scenes 1-9 of the OccuFly dataset. **Left:** Top-down view of the semantic point cloud. **Right:** Image overlap with camera centers depicted as black dots. Note that we remove scene borders with <90 % overlap to ensure geometric and semantic fidelity, as discussed in Sec. 7. Zoom in for best view.
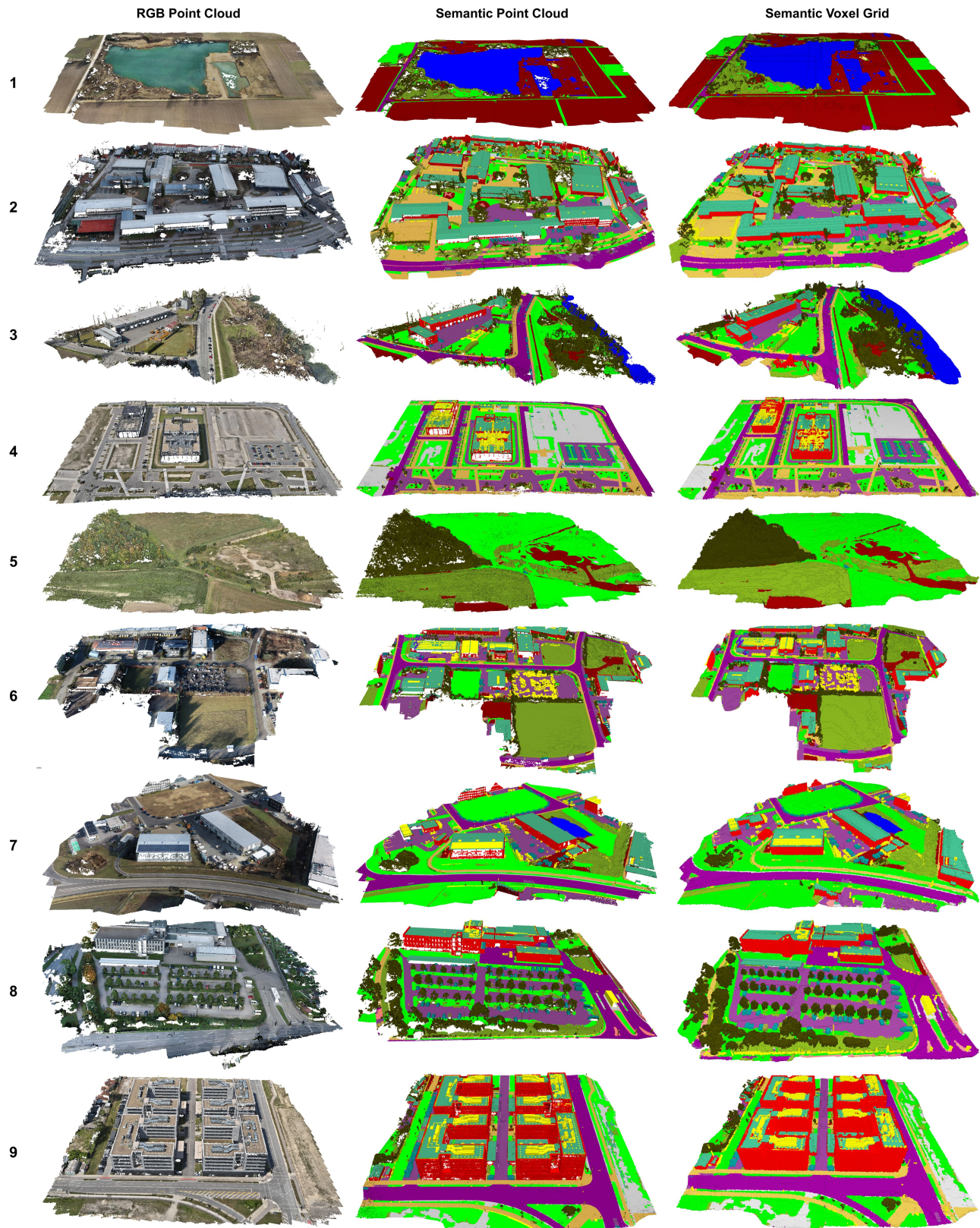
Figure 10. Scene-level outputs of our proposed data generation framework for all scenes 1-9 of the OccuFly dataset. **Left:** RGB pointcloud from 3D reconstruction (Sec. 3.3.1). **Center:** Semantic point cloud from semantic annotation (Sec. 3.3.2). **Right:** Semantic voxel grid from densification and voxelization (Sec. 3.3.3). Zoom in for best view.