

---

# USCILab3D: A Large-scale, Long-term, Semantically Annotated Outdoor Dataset

---

Kiran Lekkala\*, Henghui Bao\*, Peixu Cai, Wei Zer Lim, Chen Liu, Laurent Itti  
University of Southern California  
Los Angeles, CA 90089, USA  
klekkala, henghuib, itti@usc.edu

## Abstract

1 In this paper, we introduce the **USCILab3D dataset**, a large-scale, annotated out-  
2 door dataset designed for versatile applications across multiple domains, including  
3 computer vision, robotics, and machine learning. The dataset was acquired using a  
4 mobile robot equipped with 5 cameras and a 32-beam, 360° scanning LIDAR. The  
5 robot was teleoperated, over the course of a year and under a variety of weather  
6 and lighting conditions, through a rich variety of paths within the USC campus  
7 (229 acres =  $\sim$  92.7 hectares). The raw data was annotated using state-of-the-  
8 art large foundation models, and processed to provide multi-view imagery, 3D  
9 reconstructions, semantically-annotated images and point clouds (267 semantic  
10 categories), and text descriptions of images and objects within. The dataset also  
11 offers a diverse array of complex analyses using pose-stamping and trajectory  
12 data. In sum, the dataset offers 1.4M point clouds and 10M images ( $\sim$  6TB of  
13 data). Despite covering a narrower geographical scope compared to a whole-city  
14 dataset, our dataset prioritizes intricate intersections along with denser multi-view  
15 scene images and semantic point clouds, enabling more precise 3D labelling and  
16 facilitating a broader spectrum of 3D vision tasks. For data, code and more details,  
17 please visit our website.

## 18 1 Introduction

19 With the recent advancements in 3D vision techniques, the integration of three-dimensional perception  
20 has become integral to many interdisciplinary domains. Unlike the abundant resources available  
21 for 2D vision, the lack of comprehensive datasets for 3D vision poses a significant challenge to  
22 researchers. The progress in this field can be significantly propelled by leveraging large-scale datasets,  
23 which offer adaptability across a spectrum of downstream tasks.

24 In this paper, we present USCILab3D — a large-scale, long-term, semantically annotated outdoor  
25 dataset. USCILab3D comprises over 10 million images and 1.4 million semantic point clouds,  
26 rendering it suitable for a wide range of vision tasks.

27 Differing from smaller-scale semantic datasets or larger-scale undetailed ones, our dataset not only  
28 encompasses a wide array of outdoor multi-view scene images but also provides detailed semantic  
29 annotations, facilitating enhanced understanding and utilization of 3D perception techniques. Given  
30 the massive scale of our new dataset, as detailed below, we have thus far focused on leveraging

---

\*Equal Contribution.

31 the latest foundation models to compute detailed annotations. Our workflow using these models is  
32 detailed below.

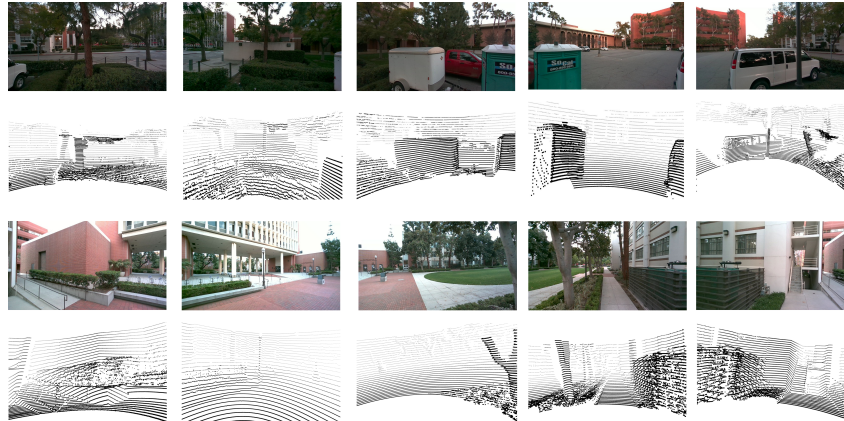


Figure 1: **Images with the respective 3D pointclouds** Our adjacent five cameras provide comprehensive coverage with overlap at the same timeframe, ensuring the captured information’s redundancy. We also show the corresponding point cloud view for every image.

## 33 2 Related datasets

34 Several large-scale scene datasets have been developed in recent years for indoor settings [17; 24; 19].  
35 Additionally, several datasets have focused on outdoor city navigation[16]. Furthermore, some  
36 datasets are generated using simulators [7; 22]. These attempt to solve the above problems, although  
37 presenting their challenges: While they offer controlled environments, there exists a noticeable gap in  
38 scene quality compared to real-world scenes.

### 39 2.1 Multi-view Scene dataset

40 Multi-view scene datasets are typically used for novel view synthesis tasks with generative models  
41 such as Neural Radiance Fields (NeRF) [15] and 3D Gaussian Splatting [12]. The LLFF dataset  
42 [14] is an early multi-view scene dataset that includes both indoor and outdoor scenes, with fewer  
43 than 1,000 low-resolution images. The DTU [11] and ScanNet [6] datasets contain between 30K and  
44 2,500K images, but they are limited to indoor scenes. The ETH3D dataset [21] provides high-quality  
45 outdoor scenes but has sparse scans and fewer than 1,000 images. Tanks and Temples [13] addresses  
46 these limitations by offering 147,000 high-quality outdoor images, which are commonly used in  
47 novel view synthesis benchmarks.

### 48 2.2 Semantic Scene dataset

49 **Indoor datasets** Datasets like [17; 24] represent large-scale 3D reconstruction datasets tailored for  
50 research in indoor robotic navigation and scene understanding. Matterport [4] is a large-scale RGB-D  
51 indoor dataset containing 10,800 panoramic views from 194,400 RGB-D images of 90 building-scale  
52 scenes. However, this dataset is limited to indoor environments and offers only 20 labels for scene  
53 annotation. In contrast, our dataset encompasses approximately 10 million images and over 4000  
54 labels, providing extensive coverage of outdoor scenes. Moreover, the inclusion of ground-truth point  
55 clouds in our dataset enhances the accuracy of alignment between 2D images and 3D annotations,  
56 surpassing the alignment capabilities of other datasets.

57 **Outdoor datasets** SemanticKITTI [3] is a widely used dataset for semantic segmentation and scene  
58 understanding in outdoor environments. It consists of dense point cloud sequences collected by a  
59 mobile LiDAR scanner which is similar to us. However, SemanticKITTI’s semantic annotations are  
60 confined to only 25 categories. In contrast, leveraging multimodal model outputs, our dataset enables

61 the labeling of almost every element within the scene, providing a comprehensive understanding of  
 62 outdoor environments.

63 Our dataset addresses the limitations of the above datasets by providing large-scale outdoor scenes  
 64 with diverse weather and lighting conditions, along with ground-truth semantic point clouds (Table  
 65 1). Leveraging multimodal foundational models, we accurately label 2D images and align them in 3D  
 66 space, resulting in precise 3D annotations.

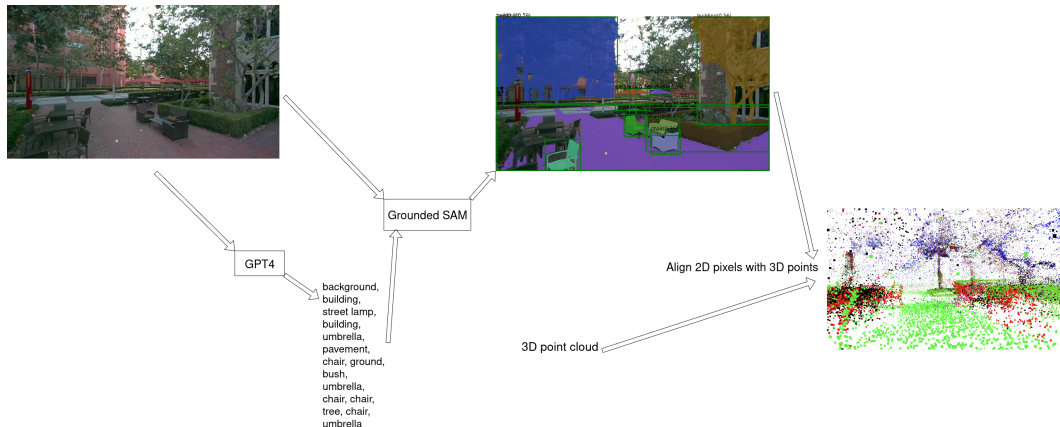


Figure 2: **The pipeline of our semantic annotations method** We use GPT4 and Grounded-SAM to create pixel-wise semantic labels and align the 2D and 3D points.

Dataset	Frames	Indoor	Outdoor	LiDAR Point Cloud	Semantic
LLFF[14]	< 1K images	✓	✓	✗	✗
DTU[11]	30K images	✓	✗	✗	✗
ScanNet[6]	2,500K images	✓	✗	✗	✗
Tanks and Temples[13]	147K images	✓	✓	✗	✗
ETH3D[21]	<1K images	✓	✗	✗	✗
Matterport3D[4]	195K images	✓	✗	✗	✓
Habitat[17]	-	✓	✗	✗	✓
iGibson[24]	-	✓	✗	✓	✓
SemanticKITTI[3]	23K scans	✗	✓	✓	✓
USCILab3d (ours)	10M images 1.4M scans	✗	✓	✓	✓

Table 1: Comparison of the existing datasets with our USCILab3D dataset.

### 67 3 Dataset collection

68 This section outlines our robot platform and data collection approach. Our robot, Beobot-v3, utilizes  
 69 multiple cameras and a LiDAR sensor for simultaneous data capture. We collect data across the USC  
 70 University Park campus and synchronize streams for analysis.

#### 71 3.1 Robot platform

72 We build our robot Beobot-v3 to collect the dataset, as shown in Figure 3. We use five Intel Realsense  
 73 D455 cameras and Velodyne HDL-32E LiDAR. The RGB images, featuring a field of view (FOV) of  
 74  $90 \times 65^\circ$  and a resolution of  $1280 \times 720$  pixels, are captured at a rate of 15 frames per second (FPS).  
 75 Utilizing a 1 MP RGB sensor, these images ensure high-quality visual data acquisition. Furthermore,  
 76 the LiDAR scans the environment at a rate of 10 Hz, capturing precise point clouds that complement  
 77 the visual data. These point clouds offer comprehensive 3D spatial information essential for scene  
 78 understanding and navigation tasks. Because of microcomputer’s limit, camera 1 and LiDAR are

79 controlled by one microcomputer, and other cameras are controlled by their own microcomputer. All  
80 microcomputers are all controlled by a central computer, our data collection system orchestrates the  
81 simultaneous scanning and recording process. As the LiDAR initiates scanning, capturing a 360°  
82 view of the environment, the data is saved directly into the system and five cameras capture images in  
83 tandem, storing them in separate ROS bag files.

### 84 3.2 Dataset collected over the entire USC campus

85 Our dataset is meticulously collected across the entirety of the USC University Park campus. Spanning  
86 an expansive area of 229 acres (0.93 km<sup>2</sup>), the campus makes our dataset diverse. From the varied  
87 architecture of its buildings to the network of roads, stairs, trails, paths, gardens, and sidewalks,  
88 each corner offers a unique scene. By dynamically selecting its route, the robot explores the full  
89 extent of the campus' diverse terrain, from thoroughfares to hidden nooks, creating a rich variety of  
90 surroundings.

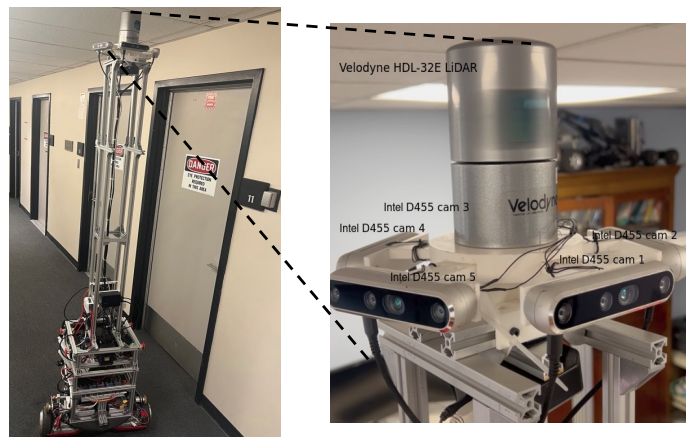


Figure 3: **Overview of the data collection robot and its hardware.** Beobot-v3 is a differential-drive, non-holonomic mobile robot, equipped with five Intel Realsense D455 cameras and one Velodyne HDL-32E LiDAR sensor used to collect the dataset.

91 The data collection occurred in many daytime sessions, with a preference for sunrise or sunset periods  
92 to avoid crowds and mitigate harsh sunlight that could degrade image quality. However, a small  
93 portion of the captured images may still exhibit the effects of powerful sunshine. The sample images  
94 are shown in Figure 4.

95 Our data collection efforts span from March 11, 2023, to March 16, 2024, encompassing 12 months.  
96 Over this time frame, the environment undergoes dynamic changes, including variations in weather,  
97 seasons, and alterations to the campus landscape, such as ongoing construction projects. This  
98 deliberate scheduling ensures that our dataset encapsulates a diverse range of environmental scenarios,  
99 enriching the dataset with a wide array of conditions for robust training and evaluation of algorithms.

### 100 3.3 Synchronization of cameras and LiDAR

101 To address the synchronization issue between the LiDAR and cameras due to the control of different  
102 microprocessors, we implement a synchronization process. Given that the LiDAR operates on the  
103 same system clock as camera 1, we only need to synchronize the remaining cameras with camera  
104 1. To achieve this, we employ a method based on feature detection and optical flow tracking. At  
105 the onset of each session, the scene remains static. Leveraging ShiTomasi corner detection [25],  
106 we identify key features in the camera images. Subsequently, using the Lucas-Kanade optical flow  
107 algorithm, we track the movement of these features over consecutive frames. If the displacement  
108 of these features exceeds a predefined threshold, indicative of the robot initiating movement, we  
109 designate this time as the session's start time.

110 Once the start time is determined for camera 1, we synchronize the start times of the remaining  
111 cameras by aligning them with the start time of camera 1. This ensures temporal coherence across  
112 all camera feeds, enabling accurate alignment of the visual and LiDAR data streams. Through  
113 this synchronization process, we establish temporal consistency across all data sources, facilitating  
114 coherent analysis and interpretation of the collected data.



Figure 4: **Sample snapshots from our dataset of various days.** These are images obtained from randomly sampling across the entire dataset.

### 115 3.4 Sensor calibration

116 By aligning the coordinate systems of the Velodyne LiDAR and the camera, we ensure that the  
117 geometric transformation from 3D to 2D space is accurate. With this calibrated setup, we can  
118 assign semantic labels to the 3D points based on the information extracted from the images. The  
119 accurate alignment between the Velodyne-frame and camera-frame ensures that the projected points  
120 correspond to the correct regions in the images, enabling us to leverage the semantic information  
121 obtained from the images to label the 3D points accurately.

## 122 4 Dataset annotation

123 In this section, we describe methods used as part of the pipeline for our semantic annotations of 3D  
124 point clouds. A high-level overview is shown in Figure 2.

### 125 4.1 GPT4-based candidate labels and clustering

126 We use GPT-4 [1] to detect the semantic labels in an image. Since images are obtained at 15Hz and  
127 the robot moves at a velocity close to 1 m/s, it is redundant and expensive to query the semantic  
128 labels for all images through GPT-4 model. Given that the image frequency is 15Hz, for about every  
129 225 images from one camera, we extract the the images of five cameras at that time. Given that the  
130 camera records at 15Hz, a 15-second interval of movement (typically less than 12 meters) ensures a  
131 small scene variation.

132 We then pass 5 images, each from every camera to GPT-4, and prompt it to estimate the semantic  
133 labels of the images using the following prompt *"List every possible semantic class that exists in the  
134 scene. List only the names and nothing else."* After standardizing and filtering the output, we obtain a  
135 total of 4162 labels. But most labels are meaningless or have similar meaning. We then again use  
136 GPT-4 to perform clustering and categorization on the estimated semantic labels.

137 After removing the meaningless labels and merging semantically equivalent labels, we obtained 257  
138 unique labels. Then, for all images we asked GPT-4 to extract objects from the image again, now  
139 with prompt is "I will give you a list of semantic class, list every possible semantic class that exists in  
140 the scene. List only the names and nothing else, split by comma." This yields the final label list for  
141 each image.

### 142 4.2 Grounded-SAM masks on pixel space

143 After we obtain the candidate labels, for equally spaced subset of images, we use those labels as  
144 an input to the Grounded-SAM model [18] to detect and segment the image by pixel. Since we are

Category	Elements
Vehicle	vehicle, bicycle, van, truck, motorcycle, golf cart, bus, car, skateboard
Nature	sky, grass, tree, shrub, shrubbery, hedge, trunk, tree trunk, green area, birds, bush, yard, plant sun, palm, rock, soil, leaf, leaves, water, flower, branch, bushes, vegetation, bird, ivy
Human	person, hand
Ground	pavement, curb, gravel, rail, sidewalk, street, walkway, floor, road, pedestrian walkway, crosswalk ramp, garden, ground, pathway, paving stone, golf course, parking lot, drainage grate, mulch
Structure	monument, structure, courtyard, fountain, public space, construction, emergency station ceiling, fence, gate, wall, balcony, container, stadium, lattice, shed, house, construction pipe, roof, building, sports field, campus, toilet, baseball field, architecture site, parking structure, garage, scaffolding, archway, call station
Street Furniture	bench, pole, feeding station, patio, handicap, barrier, hydrant, construction cone, construction barrier lamp post, lamp, trash can, receipt, sign, parking meter, public art, statue, sculpture bollard, bus stop, park bench
Architectural Elements	drain cover, manhole cover, vent, air vent, arch, sill, doorway, baluster, security camera, electric box corridor, stair, ventilation grill, door handle, entrance, post, air unit, pillar, balustrade, handrail window, door, elevator, gutter, bleachers, tank, generator, utility meter
General Objects	umbrella, table, chair, stroller, furniture, board, bottle, canopy, outdoor gear, advertisement, station pot, rack, flag, locker, ladder, garbage, bulletin board, pallet, planter, equipment, tent, base, hat curtain, blinds, cardboard, box, tire, wheels, bag, bed, frame, bucket, painting, poster, machine shadow, reflection, traffic cone, parking space line, space line, road marking
Signs and Symbols	parking symbol, stop sign, street sign, road sign, symbol, plaque, banner, graffiti, waste container signboard, security camera, camera, warning sign, fire safety sign, transportation sign handicap sign, closed sign, exit sign, parking sign, reservation sign, rec sign
Materials	concrete, brick, construction materials, stone, wood, plastic, metal, glass, iron, materials
Lighting	outdoor lighting, light, street light, indoor light, lantern, sunlight, shade
Miscellaneous	cover, trash, outdoor, chain, unit, security, exterior, fire, electric, meter, lettering, phone, debris, railway text, potted, space, portable, cone, slight, cross, marker, grate, blea, stoller, units, picnic, electrical cable, basin, pavilion, ster, bal, field, curve, bod, bay, pal, firent, box, exit, baseball, image, rec, sports public, piping, grill, guttering, utility, call, case, recacle, gut, hydra, air line, tile, cardboard, patch, reservoir, valve

Table 2: **Clustering of the semantic labels.** We use GPT-4 to cluster 267 labels into 12 categories using the prompt "Could you help me classify by following category: Vehicle, Nature, Human, Ground, Structure, Street Furniture, Architectural Elements."

145 using a differential-drive robot and it could potentially rotate left or right, images may look very  
 146 different quite rapidly, so we merge the five image labels from GPT-4 and pass to next step. After  
 147 conducting our experiments, we found that the presence of unrelated labels (not visually represented  
 148 in the images) does not significantly influence the results of Grounded-SAM. This observation is  
 149 reflected in Figure 5 through the percentage of incorrect pixel labels in the masks of 2 images. We  
 150 show the top 50 frequent objects and their pixel percentage in images of our dataset in Figure 6.



Figure 5: **(a) Robustness of Grounded SAM to prompts (left).** Comparison of the semantic masks obtained using different prompts for the same image by Grounded-SAM model, showing the robustness of the model. **(b) Percentage of incorrect pixel labels (right).** Quantitative measures to show robustness through the change in the percentage of incorrect pixel labels with additional prompts.

### 151 4.3 Post-processing after Grounded-SAM

152 Grounded-SAM’s output is not always using the same vocabulary as our input labels, e.g., one may  
 153 prompt it for ‘vehicle’ but obtain a segmented ‘car’. It may also generate meaningless words or  
 154 words having similar meaning. To address this, we perform clustering and categorization as in section  
 155 4.1 again to merge all similar labels. Additionally, we manually merge and remove some words.  
 156 Ultimately, we obtain 267 labels and 12 categories (Table 2).

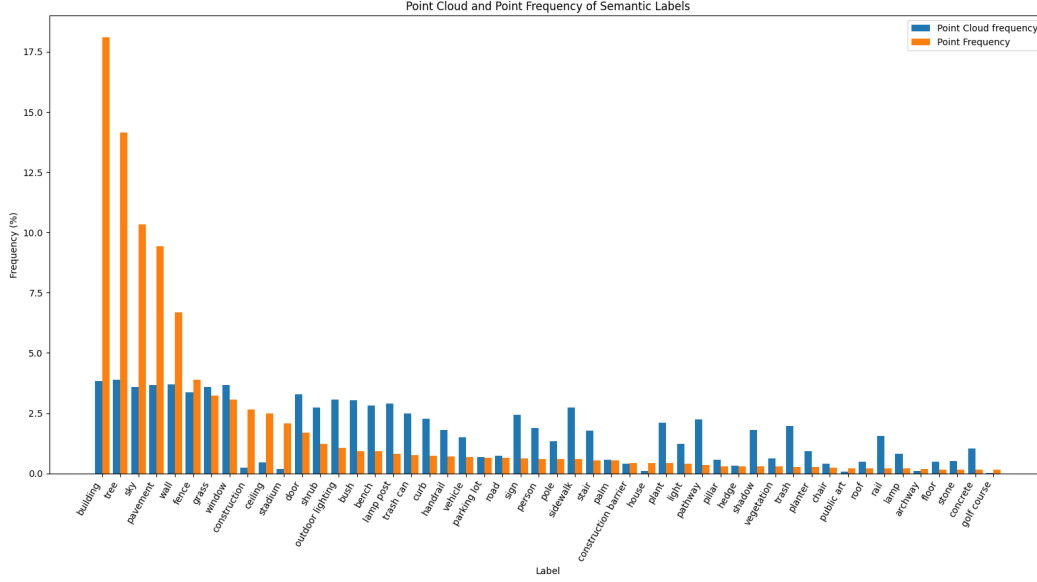


Figure 6: **Histogram of the semantic labels frequency in point cloud scans and points.** Top 50 frequently estimated semantic classes in points(orange), and corresponding point cloud scan frequency

#### 157 4.4 Projecting 2D semantic masks to 3D pointcloud

158 From the LIDAR data, we reconstruct 3D trajectories of the robot throughout the dataset. Essentially,  
 159 we compute a pose transformation for each LiDAR scan in the dataset. We then interpolate the LiDAR  
 160 poses to the camera images using the extrinsic parameters corresponding to the transformation of  
 161 each camera with respect to the LiDAR sensor. This results in a pose estimate for every camera image  
 162 in the dataset.

163 By utilizing the semantic map of every image obtained from Grounded SAM, we use ground truth  
 164 camera intrinsics and extrinsics to accurately project 3D point clouds onto 2D images, following  
 165 equation. Here,  $(X, Y, Z)$  represents the world coordinates of a point, while  $(x, y)$  denotes the  
 166 coordinates of the point projected onto the image plane, measured in pixels.  $r$  and  $t$  are rotation  
 167 and translation.  $c_x, c_y$  represents the principal point, and  $f_x, f_y$  are the focal lengths in pixels.  
 168 Subsequently, we align the 2D and 3D points to assign labels to the 3D points.

$$\begin{pmatrix} x \\ y \\ 1 \end{pmatrix} \sim \begin{pmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \end{pmatrix} \begin{pmatrix} X \\ Y \\ Z \\ 1 \end{pmatrix} \quad (1)$$

169 Considering the presence of moving objects and calibration errors, there may be some offset for each  
 170 projection. To reduce erroneous labels, we run DBSCAN clustering [8] on each label projection to  
 171 check whether the 3D points projected belong to a single cluster. If they do not, we only label the  
 172 cluster with the most points.

#### 173 4.5 Released data

174 We release the raw ROS Bagfiles, and extracted images, point cloud files, COLMAP [20] poses and  
 175 sparse reconstructions. The raw data consists of a set of sequences, each of which is collected during  
 176 a specific data recording session. To make the data more manageable, we divide each session into  
 177 different subsequences or "sectors", with each sector consisting of 1250 images and roughly 167  
 178 point cloud scans. In addition, we conducted face detection and applied blurring techniques to ensure  
 179 privacy protection on campus.

180 **Multi-view images** Each image is named according to the convention `cam[id]-[timestamp].jpg`.  
181 We estimate synchronized timestamps for all images within a sector, using the method mentioned  
182 in section 3.3. The wide field of view (FoV) of 90 degrees for each of the five cameras results in  
183 significant overlap between their respective images, as depicted in Figure 1. This substantial overlap  
184 ensures more robust Structure from Motion (SfM) reconstruction. By having multiple views of the  
185 same scene, the SfM algorithm can triangulate feature points more accurately, leading to a more  
186 precise reconstruction of the 3D environment. This overlap also aids in improving the accuracy of  
187 semantic labelling. By leveraging overlapping information from multiple viewpoints, inconsistencies  
188 or errors in semantic annotations of 3D points from 2D-pixel maps can be identified and rectified  
189 through cross-validation. This double-checking mechanism helps to enhance the reliability of  
190 semantic labels assigned to objects in the scene.

191 **Semantic instances and masks of images** In addition to the raw image data, we also provide  
192 semantic labels and label masks generated by Grounded-SAM for each image in the dataset. These  
193 labels offer valuable insights into the semantic understanding of the scene, allowing researchers to  
194 perform tasks such as semantic segmentation and object detection.

195 **Semantic annotated point cloud scans** As mentioned before, the pointcloud streams are captured at  
196 10Hz. Similar to KITTI Semantic [3], we extract each of the pointclouds scans and annotate the 3D  
197 points by assigning semantic labels to individual points based on the closest image’s label, using the  
198 method outlined in section 4.3. The color and corresponding label for each point are saved in a JSON  
199 file named `labels.json`, ensuring easy access and interpretation of the semantic annotations.

200 **Semantic annotated session point clouds** In addition to the individual semantic annotated point  
201 cloud scans, we have processed each session’s point cloud data using LeGO-LOAM [23] to generate  
202 merged point cloud of a sector. We mention the statistics of the distribution of points in each of  
203 the point cloud scans and the merged point clouds in the supplemental material. Unlike the point  
204 cloud scans, sector-based point clouds have more points and offer a comprehensive overview of the  
205 semantic annotated scene. Through these semantic point clouds, researchers can gain deeper insights  
206 into the semantic structure and composition of the environment.

207 **Pose annotations for images.** We release interpolated poses from LeGO-LOAM, and COLMAP  
208 Structure from Motion (SfM) [20]. The COLMAP SfM results can serve as inputs for some generative  
209 model like NeRF or 3D Gaussian Splatting. Further, by utilizing the poses computed by COLMAP,  
210 we aim to improve the precision of our annotations given the different sampling rates of the LiDAR  
211 (10Hz) and cameras (15Hz). This alignment is crucial for accurately projecting semantic labels onto  
212 the 3D points based on the information extracted from the images. We are currently investigating  
213 how to best merge the LiDAR and COLMAP poses, likely resulting in a unified set of poses indexed  
214 non-uniformly in time, for each image and for each point cloud. We expect that these unified poses  
215 will be released with the next version our dataset.

216 **Robotic dataset for visual navigation.** Our dataset comprises diverse sequences captured within  
217 a university environment, reflecting a range of real-world scenarios. Leveraging the compact form  
218 factor of our robot, we collected data across a variety of settings including roads, outdoor lobbies,  
219 ramps, and other typical campus landscapes. This dataset is particularly valuable for applications in  
220 visual navigation and is integrated into the comprehensive Open X-Embodiment dataset [5].

## 221 5 Benchmarks

### 222 5.1 Evaluation on Novel View Synthesis

223 We examine the current state-of-the-art (SOTA) Novel View Synthesis methods on several datasets:  
224 USCILab3D, ETH3D[21], Mip-NeRF360[2], Tanks&Temples[13], Deep Blending[9], and Deep  
225 Blending[9]. For each dataset, we run 3D Gaussian Splatting and evaluate the generated image  
226 quality using PSNR, SSIM, and L-PIPS metrics. For each scene, we use 7/8 of the data as the training  
227 set and 1/8 as the test set, then calculate the average result for each scene. Considering the large size  
228 of our dataset, we randomly extract one sector from each session to compute the average result.



229 Our dataset achieves superior PSNR, SSIM, and the best L-PIPS performance compared to other  
 230 datasets (Table 3). Among these datasets, ours is the only one that provides large-scale scenes,  
 231 making it suitable for a wider range of applications, such as simulators [10].

	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Resolution $\downarrow$	iteration
USCILab3D (ours)	26.02	0.86	0.20	1280 $\times$ 720	7000
ETH3D[21]	21.25	0.83	0.27	6048 $\times$ 4032	7000
Tanks&Temples [13]	21.20	0.77	0.28	980 $\times$ 540	7000
Mip-NeRF360[2]	25.19	0.75	0.25	1256 $\times$ 828	7000
Deep Blending[9]	27.01	0.87	0.32	1332 $\times$ 876	7000

Table 3: **Performance comparison of 3D Gaussian splatting on different datasets.** Our dataset achieves superior performance compared to other datasets. Although Deep Blending demonstrates a higher PSNR, it only contains 2.6K images.

## 232 5.2 Evaluation on Semantic Segmentation and Completion

233 We also evaluate our dataset using key tasks: semantic segmentation, panoptic segmentation, and  
 234 semantic scene completion. Semantic segmentation is crucial for understanding and labeling every  
 235 point in a 3D point cloud with a specific class, providing detailed insights into the composition of the  
 236 scene. Panoptic segmentation extends this by not only classifying each point but also distinguishing  
 237 between different instances of the same class. This is particularly valuable for environments with  
 238 multiple similar objects, enhancing the dataset’s utility in more complex and dynamic scenarios.  
 239 Lastly, semantic scene completion involves predicting the complete geometry and semantics of a  
 240 scene, including occluded and unobserved regions. This task is vital for creating comprehensive  
 241 and accurate representations of environments, which is indispensable for advanced applications in  
 242 augmented reality and spatial analysis. Due to page limitations, we have included the results in the  
 243 supplemental material.

## 244 6 Caveats

245 Thus far, our annotations have been machine-generated using the latest foundation models. Although  
 246 this may pose a few risks, nevertheless, to the best of our knowledge, our method is the first of its  
 247 kind to annotate 3D point clouds using image and text based foundational models without any manual  
 248 intervention. Casual inspection by authors suggests that the annotations are indeed of high quality.  
 249 However, we plan to validate them by hiring a group of human annotators to inspect and possibly  
 250 correct a fraction of the machine-generated annotations. We expect that this will be completed by  
 251 the time of publication.

## 252 7 Discussion and Conclusion

253 In this paper, we introduced the USCILab3D dataset, a comprehensive outdoor 3D dataset designed  
 254 to address the limitations of existing datasets in the domain of 3D scene understanding and navigation.  
 255 Our dataset offers a diverse array of complex intersections and outdoor scenes meticulously collected  
 256 across the USC University Park campus. With approximately 10 million images and 1.5 million  
 257 dense point cloud scans, our dataset prioritizes intricate areas, enabling more precise 3D labelling  
 258 and facilitating a broader spectrum of 3D vision tasks.

259 Moving forward, we believe that the USCILab3D dataset will serve as a valuable resource for re-  
 260 searchers and practitioners across various domains, including computer vision, robotics, and machine  
 261 learning. We anticipate that the dataset will stimulate further advancements in 3D vision-based  
 262 models and foster the development of robust algorithms capable of tackling real-world challenges in  
 263 outdoor environments.

264 **References**

- 265 [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni  
266 Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4  
267 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 268 [2] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman.  
269 Mip-nerf 360: Unbounded anti-aliased neural radiance fields. *CoRR*, abs/2111.12077, 2021.  
270 URL <https://arxiv.org/abs/2111.12077>.
- 271 [3] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall. Se-  
272 manticKITTI: A Dataset for Semantic Scene Understanding of LiDAR Sequences. In *Proc. of*  
273 *the IEEE/CVF International Conf. on Computer Vision (ICCV)*, 2019.
- 274 [4] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis  
275 Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in  
276 indoor environments. *International Conference on 3D Vision (3DV)*, 2017.
- 277 [5] Open X.-Embodiment Collaboration, Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex  
278 Bewley, Alexander Herzog, Alex Irpan, Alexander Khazatsky, Anant Raj, Anikait Singh,  
279 Anthony Brohan, Antonin Raffin, Ayzaan Wahid, Ben Burgess-Limerick, Beomjoon Kim,  
280 Bernhard Schölkopf, Brian Ichter, Cewu Lu, Charles Xu, Chelsea Finn, Chenfeng Xu, Cheng  
281 Chi, Chenguang Huang, Christine Chan, Chuer Pan, Chuyuan Fu, Coline Devin, Danny Driess,  
282 Deepak Pathak, Dhruv Shah, Dieter Buechler, Dmitry Kalashnikov, Dorsa Sadigh, Edward  
283 Johns, Federico Ceola, Fei Xia, Freck Stulp, Gaoyue Zhou, Gaurav S. Sukhatme, Gautam  
284 Salhotra, Ge Yan, Giulio Schiavi, Gregory Kahn, Hao Su, Haoshu Fang, Haochen Shi, Heni Ben  
285 Amor, Henrik I. Christensen, Hiroki Furuta, Homer Walke, Hongjie Fang, Igor Mordatch, Ilija  
286 Radosavovic, and et al. Open x-embodiment: Robotic learning datasets and RT-X models.  
287 *CoRR*, abs/2310.08864, 2023. doi: 10.48550/ARXIV.2310.08864. URL [https://doi.org/](https://doi.org/10.48550/arXiv.2310.08864)  
288 [10.48550/arXiv.2310.08864](https://doi.org/10.48550/arXiv.2310.08864).
- 289 [6] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas A. Funkhouser, and  
290 Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *2017 IEEE*  
291 *Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July*  
292 *21-26, 2017*, pages 2432–2443. IEEE Computer Society, 2017. doi: 10.1109/CVPR.2017.261.  
293 URL <https://doi.org/10.1109/CVPR.2017.261>.
- 294 [7] Alexey Dosovitskiy, Germán Ros, Felipe Codevilla, Antonio M. López, and Vladlen Koltun.  
295 CARLA: an open urban driving simulator. In *1st Annual Conference on Robot Learning,*  
296 *CoRL 2017, Mountain View, California, USA, November 13-15, 2017, Proceedings*, volume 78  
297 of *Proceedings of Machine Learning Research*, pages 1–16. PMLR, 2017. URL [http://](http://proceedings.mlr.press/v78/dosovitskiy17a.html)  
298 [proceedings.mlr.press/v78/dosovitskiy17a.html](http://proceedings.mlr.press/v78/dosovitskiy17a.html).
- 299 [8] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm  
300 for discovering clusters in large spatial databases with noise. In Evangelos Simoudis, Jiawei  
301 Han, and Usama M. Fayyad, editors, *Proceedings of the Second International Conference on*  
302 *Knowledge Discovery and Data Mining (KDD-96), Portland, Oregon, USA*, pages 226–231.  
303 AAAI Press, 1996. URL <http://www.aaai.org/Library/KDD/1996/kdd96-037.php>.
- 304 [9] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel J.  
305 Brostow. Deep blending for free-viewpoint image-based rendering. *ACM Trans. Graph.*, 37  
306 (6):257, 2018. doi: 10.1145/3272127.3275084. URL [https://doi.org/10.1145/3272127.](https://doi.org/10.1145/3272127.3275084)  
307 [3275084](https://doi.org/10.1145/3272127.3275084).
- 308 [10] Laurent Itti Henghui Bao\*, Kiran Lekkala\*. Real world navigation in a simulator: A benchmark,  
309 2024. URL <https://sites.google.com/usc.edu/real-world-navigation/home>.

- 310 [11] Rasmus Ramsbøl Jensen, Anders Lindbjerg Dahl, George Vogiatzis, Engin Tola, and Henrik  
311 Aanæs. Large scale multi-view stereopsis evaluation. In *2014 IEEE Conference on Computer  
312 Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014*, pages  
313 406–413. IEEE Computer Society, 2014. doi: 10.1109/CVPR.2014.59. URL <https://doi.org/10.1109/CVPR.2014.59>.
- 315 [12] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian  
316 splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139:1–139:14, 2023.  
317 doi: 10.1145/3592433. URL <https://doi.org/10.1145/3592433>.
- 318 [13] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: bench-  
319 marking large-scale scene reconstruction. *ACM Trans. Graph.*, 36(4):78:1–78:13, 2017. doi:  
320 10.1145/3072959.3073599. URL <https://doi.org/10.1145/3072959.3073599>.
- 321 [14] Ben Mildenhall, Pratul P. Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi  
322 Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis  
323 with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 2019.
- 324 [15] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi,  
325 and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *CoRR*,  
326 abs/2003.08934, 2020. URL <https://arxiv.org/abs/2003.08934>.
- 327 [16] Piotr Mirowski, Andras Banki-Horvath, Keith Anderson, Denis Teplyashin, Karl Moritz Her-  
328 mann, Mateusz Malinowski, Matthew Koichi Grimes, Karen Simonyan, Koray Kavukcuoglu,  
329 Andrew Zisserman, and Raia Hadsell. The streetlearn environment and dataset. *CoRR*,  
330 abs/1903.01292, 2019. URL <http://arxiv.org/abs/1903.01292>.
- 331 [17] Santhosh Kumar Ramakrishnan, Aaron Gokaslan, Erik Wijmans, Oleksandr Maksymets, Alexan-  
332 der Clegg, John M. Turner, Eric Undersander, Wojciech Galuba, Andrew Westbury, Angel X.  
333 Chang, Manolis Savva, Yili Zhao, and Dhruv Batra. Habitat-matterport 3d dataset (HM3D):  
334 1000 large-scale 3d environments for embodied AI. In Joaquin Vanschoren and Sai-Kit Ye-  
335 ung, editors, *Proceedings of the Neural Information Processing Systems Track on Datasets  
336 and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*, 2021.  
337 URL [https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/  
338 34173cb38f07f89ddebcb2ac9128303f-Abstract-round2.html](https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/34173cb38f07f89ddebcb2ac9128303f-Abstract-round2.html).
- 339 [18] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu  
340 Huang, Yukang Chen, Feng Yan, Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang, Hongyang  
341 Li, Qing Jiang, and Lei Zhang. Grounded SAM: assembling open-world models for diverse  
342 visual tasks. *CoRR*, abs/2401.14159, 2024. doi: 10.48550/ARXIV.2401.14159. URL <https://doi.org/10.48550/arXiv.2401.14159>.
- 344 [19] Denys Rozumnyi, Stefan Popov, Kevis-Kokitsi Maninis, Matthias Nießner, and Vittorio Fer-  
345 rari. Estimating generic 3d room structures from 2d annotations. In Alice Oh, Tristan  
346 Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Ad-  
347 vances in Neural Information Processing Systems 36: Annual Conference on Neural In-  
348 formation Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December  
349 10 - 16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/  
350 76bf913ad349686b2aa552a1c6ee0a2e-Abstract-Datasets\\_and\\_Benchmarks.html](http://papers.nips.cc/paper_files/paper/2023/hash/76bf913ad349686b2aa552a1c6ee0a2e-Abstract-Datasets_and_Benchmarks.html).
- 351 [20] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-Motion Revisited. In  
352 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- 353 [21] Thomas Schöps, Johannes L. Schönberger, Silvano Galliani, Torsten Sattler, Konrad Schindler,  
354 Marc Pollefeys, and Andreas Geiger. A multi-view stereo benchmark with high-resolution  
355 images and multi-camera videos. In *2017 IEEE Conference on Computer Vision and Pattern  
356 Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 2538–2547. IEEE  
357 Computer Society, 2017. doi: 10.1109/CVPR.2017.272. URL [https://doi.org/10.1109/  
358 CVPR.2017.272](https://doi.org/10.1109/CVPR.2017.272).

- 359 [22] Shital Shah, Debadepta Dey, Chris Lovett, and Ashish Kapoor. Airsim: High-fidelity visual  
360 and physical simulation for autonomous vehicles. In Marco Hutter and Roland Siegwart,  
361 editors, *Field and Service Robotics, Results of the 11th International Conference, FSR 2017,*  
362 *Zurich, Switzerland, 12-15 September 2017*, volume 5 of *Springer Proceedings in Advanced*  
363 *Robotics*, pages 621–635. Springer, 2017. doi: 10.1007/978-3-319-67361-5\_40. URL [https://doi.org/10.1007/978-3-319-67361-5\\_40](https://doi.org/10.1007/978-3-319-67361-5_40).
- 365 [23] Tixiao Shan and Brendan J. Englot. Lego-loam: Lightweight and ground-optimized lidar  
366 odometry and mapping on variable terrain. In *2018 IEEE/RSJ International Conference on*  
367 *Intelligent Robots and Systems, IROS 2018, Madrid, Spain, October 1-5, 2018*, pages 4758–  
368 4765. IEEE, 2018. doi: 10.1109/IROS.2018.8594299. URL <https://doi.org/10.1109/IROS.2018.8594299>.
- 370 [24] Bokui Shen, Fei Xia, Chengshu Li, Roberto Martín-Martín, Linxi Fan, Guanzhi Wang, Claudia  
371 Pérez-D’Arpino, Shyamal Buch, Sanjana Srivastava, Lyne Tchapmi, Micael Tchapmi, Kent  
372 Vainio, Josiah Wong, Li Fei-Fei, and Silvio Savarese. igibson 1.0: A simulation environment for  
373 interactive tasks in large realistic scenes. In *IEEE/RSJ International Conference on Intelligent*  
374 *Robots and Systems, IROS 2021, Prague, Czech Republic, September 27 - Oct. 1, 2021*, pages  
375 7520–7527. IEEE, 2021. doi: 10.1109/IROS51168.2021.9636667. URL <https://doi.org/10.1109/IROS51168.2021.9636667>.
- 377 [25] Jianbo Shi and Carlo Tomasi. Good features to track. In *Conference on Computer Vision*  
378 *and Pattern Recognition, CVPR 1994, 21-23 June, 1994, Seattle, WA, USA*, pages 593–600.  
379 IEEE, 1994. doi: 10.1109/CVPR.1994.323794. URL [https://doi.org/10.1109/CVPR.](https://doi.org/10.1109/CVPR.1994.323794)  
380 [1994.323794](https://doi.org/10.1109/CVPR.1994.323794).

## 381 Checklist

382 The checklist follows the references. Please read the checklist guidelines carefully for information on  
383 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or  
384 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing  
385 the appropriate section of your paper or providing a brief inline description. For example:

- 386 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 387 • Did you include the license to the code and datasets? **[No]** The code and the data are  
388 proprietary.
- 389 • Did you include the license to the code and datasets? **[N/A]**

390 Please do not modify the questions and only use the provided macros for your answers. Note that the  
391 Checklist section does not count towards the page limit. In your paper, please delete this instructions  
392 block and only keep the Checklist section heading above along with the questions/answers below.

393 1. For all authors...

- 394 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
395 contributions and scope? **[Yes]**
- 396 (b) Did you describe the limitations of your work? **[Yes]** See Section6 : Caveats
- 397 (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
- 398 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
399 them? **[Yes]**

400 2. If you are including theoretical results...

- 401 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
- 402 (b) Did you include complete proofs of all theoretical results? **[N/A]**

403 3. If you ran experiments (e.g. for benchmarks)...

- 404 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
405 mental results (either in the supplemental material or as a URL)? **[Yes]** In our website  
406 we provide data and code.
- 407 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
408 were chosen)? **[Yes]** In the supplemental material.
- 409 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
410 ments multiple times)? **[Yes]**
- 411 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
412 of GPUs, internal cluster, or cloud provider)? **[Yes]** In the supplemental material.

413 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 414 (a) If your work uses existing assets, did you cite the creators? **[Yes]**
- 415 (b) Did you mention the license of the assets? **[Yes]** In URL
- 416 (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**
- 417 (d) Did you discuss whether and how consent was obtained from people whose data you're  
418 using/curating? **[N/A]**
- 419 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
420 information or offensive content? **[Yes]** In section 4.

421 5. If you used crowdsourcing or conducted research with human subjects...

- 422 (a) Did you include the full text of instructions given to participants and screenshots, if  
423 applicable? **[N/A]**
- 424 (b) Did you describe any potential participant risks, with links to Institutional Review  
425 Board (IRB) approvals, if applicable? **[N/A]**
- 426 (c) Did you include the estimated hourly wage paid to participants and the total amount  
427 spent on participant compensation? **[N/A]**