

000 YawP³: Yaw-invariant Parametrization and 001 Panoramic Planar-reconstruction

004 Anonymous ECCV submission
005

006 Paper ID 5951
007

009
010 **Abstract.** This paper presents the first method for indoor planar recon-
011 struction from a panoramic image. We leverage three unique character-
012 istics in indoor scenes: 1) most planar surfaces are either horizontal or
013 vertical, 2) the 2D orientation of vertical planes co-vary with the change
014 of the yaw angle of the camera, and 3) most vertical planes share com-
015 mon 2D orientations. To inherently ingrain these priors, our model first
016 segments horizontal/vertical planes (HV-planes) for separate treatments.
017 Most importantly, we propose a novel yaw-invariant parameterization
018 for vertical planes to solve the yaw ambiguity problem of 360° and ef-
019 fectively cluster vertical planar segments with a shared 2D orientation.
020 Finally, we fuse the predicted geometry information and plane instance
021 segmentation into a piece-wise planar reconstruction. We evaluate our
022 method on our newly extracted panoramic piece-wise HV-planar dataset
023 derived from three large-scale RGB-D panorama datasets. For benchmark
024 baselines, we train two state-of-the-art planar models on our dataset
025 with modifications to help them adapt to our 360° HV-planar dataset.
026 Our method significantly outperforms all baselines and achieves superior
027 visually pleasing reconstruction of indoor scenes.
028

029 **Keywords:** Planar reconstruction, 3D reconstruction, panoramic image,
030 360° dataset
031

032 1 Introduction 033

034 Reconstructing planar surfaces from single-view images have many applications
035 such as interior modeling, virtual and augmented reality, robot navigation, scene
036 understanding, and 3D reconstruction. Generally, an indoor scene can be approx-
037 imated with a small set of dominant planes, which makes planar reconstruction
038 suitable for 3D indoor modeling.

039 Recent works on planar reconstruction [12–14] built upon recent state-of-the-
040 art instance segmentation methods achieve promising results. However, these
041 works are mostly trained on the planar datasets derived from ScanNet [5] and
042 NYUv2 [16] with a small field-of-view (FoV). Data of this kind require multiple
043 images to reconstruct entire scenes, which would cost more computational time
044 and resources. As 360° devices get popularized, the amount of 360° data has
significantly increased. With 360° images as input, 3D reconstruction of an entire
scene can be done with only one snapshot on the input data. Considering the

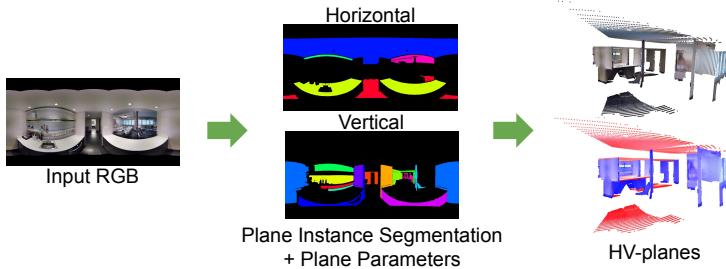


Fig. 1: Planar surfaces of man-made structures are mostly horizontal or vertical with respect to the gravity direction (HV-planes). Given an RGB panorama (left), we propose to model a 3D scene as HV-Planes (middle). The red and blue points represent pixels belonging to H-planes and V-planes respectively (right).

benefits of 360° data in planar reconstruction and the gap of existing literature in this research field, we believe the task of planar reconstruction from a panoramic image is worthy of investigation.

In this paper, we propose YawP³, the first deep neural network addressing 360° planar reconstruction. We manage to simplify the tasks of planar instance segmentation based on surface orientation grouping, which provides geometric awareness to the succeeding module of instance segmentation. Moreover, we propose a yaw-invariant parameter in V-plane prediction to solve the yaw ambiguity problem in omnidirectional images. We show that the yaw-invariant parameterization brings about significant improvements in 3D evaluation metrics that consider geometric accuracy. To evaluate our proposed method, we present a 360° planar dataset as well as baseline models adapted from state-of-the-art models of perspective views. We further simplify this task by presenting novel HV-plane annotations as data ground truth to reduce the degree-of-freedom in planar parameters. Fig. 1 illustrates the key concepts of the design. Our model splits the prediction into H/V branches based on the new representation and combines the yaw-invariant representation into V-branch since H-plane only consists of one degree of freedom, which is perpendicular to yaw-direction. The plane geometry information is fused with planar instance segmentation for final reconstruction. Our model achieves state-of-the-art performance on the newly presented benchmark.

The contributions of this work are manifold. In terms of **technical contribution**, the proposed YawP³ framework consists of *i*) a novel yaw-invariant vertical plane parameter representation that addresses the 360° yaw ambiguity and also boosts other existing methods; *ii*) a new geometry-aware mechanism of plane-based segmentation which exploits the geometry information to improve planar instance segmentation. In terms of **system contribution**, we construct a new 360° piece-wise planar benchmark, which focuses on horizontal and vertical planes, where the proposed YawP³ outperforms the two adapted state-of-the-art planar reconstruction methods that are originally designed for perspective images.

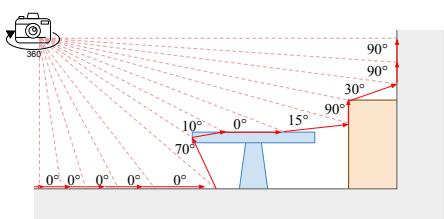
090 2 Related Work

091 Reconstructing piece-wise 3D planar surfaces from an image could be decomposed
092 into two problems: segmenting plane instances and estimating plane geometric
093 parameters. To solve these problems, PlaneNet [13] reconstructs a fixed number
094 of planes by estimating plane parameters and plane segmentation masks both in
095 an instance-wise manner by jointly training CNN and DCRF as [27]. Recent state-
096 of-the-art approaches relax the constraint on the number of planes by exploiting
097 popular frameworks in instance segmentation. PlaneRCNN [12] modifies the two-
098 stage architecture Mask R-CNN [9] with object category classification replaced
099 by plane geometry prediction, followed by a network to refine the segmentation
100 masks. PlanarReconstruction [24] predicts per-pixel plane parameters and adopts
101 the method similar to [3, 7, 10, 24], which trains a network to map each pixel to
102 embedding space and then clusters the embedded pixels to generate instances.
103 In the field of piece-wise planar reconstruction, geometric clues such as plane
104 orientation could be useful for segmenting plane instances. However, existing
105 works utilize limited plane geometry information in their instance segmentation
106 processes. For instance, PlanarReconstruction [24] integrates the segmented plane
107 instance masks into geometry prediction with the instance-level loss for plane
108 parameters, but the segmentation process is still independent to the estimated
109 geometry. PlaneRCNN [12] uses the initial estimation of segmentation maps and
110 depth as input to train a refinement network for segmentation. In contrast, our
111 method explicitly involves the estimated plane orientations in our plane instance
112 segmentation process.

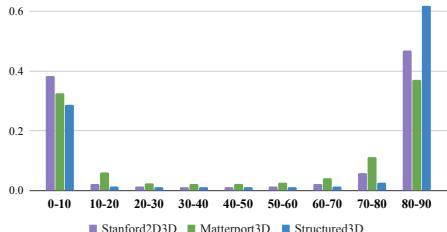
113 Previous works for planar reconstruction focus on perspective images, while
114 this paper targets 360° images. In 360° domain, estimating per-pixel surface
115 normals from an equirectangular image is challenging for the CNN layers as
116 the yaw-rotation of the 360° camera is ambiguous. In other words, 360° camera
117 yaw-rotation affects plane parameters, but the counterpart left-right circular
118 shifting on the equirectangular image is less aware by the CNN layers due to
119 its translation invariant property. Although surface normals are fundamental to
120 many 360° applications aside from planar reconstruction, none of the existing
121 works which estimate surface normals from an equirectangular image [22, 17, 20,
122 6] discuss the 360° camera yaw ambiguous problem. To this end, we propose
123 a method that targets 360° images, with a yaw-invariant parameterization of
124 vertical plane orientation for our applications.

127 3 PanoHV Dataset

128 In this section, we first introduce some large-scale panoramic public datasets
129 that are used to construct our dataset. We then show the statistical analysis
130 on these datasets to support the gravity aligned assumption and the validity
131 of scene approximation with HV-planes, followed by the description of the pre-
132 processing mechanism. Finally, we briefly describe the ground truth extraction
133 algorithm. Please refer to supplementary material for the full description of our



(a)



(b)

Fig. 2: The histograms of angles yielded by adjacent points on a vertical ‘slice’ of a point cloud. We divide the angles from 0° to 90° into nine bins. An angle is calculated between the horizon and the vector connecting two 3D points projected from the corresponding top-bottom adjacent pair of pixels in the image. This local analysis provides a strong cue to support the gravity aligned assumption and will be used in benchmark construction. Please refer to Sec. 5 for more details. (a) An illustration of a vertical ‘slice’ taken from a point cloud. (b) Statistics of angles quantized as histograms

data preparation algorithm. All panoramic images in this work are represented in the equirectangular format if not specified otherwise.

3.1 360° Dataset Sources

Inspired by [12, 13], we construct our dataset from three public 360° RGB-D datasets, including Structure3D [26] in synthetic environments, and Matterport3D [4] and Stanford2D3D [2] in real-world scenes. These three datasets consist of large-scale aggregation of panoramic RGB images and depth maps. Our annotations of planar masks and parameters are derived from the ground-truth depths.

3.2 Gravity Aligned and HV-plane Assumption

Similar to the *Manhattan* assumption made by most of the previous methods (*e.g.* [15, 18, 28]), we assume the gravity aligned property for simplifying the task of *indoor piece-wise planar reconstruction from a single panorama*. More specifically, we aim to approximate the indoor scene by horizontal planes (**H-planes** with the normal parallel to the gravity direction) and vertical planes (**V-planes** with the normal perpendicular to the gravity direction). Although such an assumption complies with our intuition, we further perform two quantitative analyses on the aforementioned three large-scale indoor panorama datasets. These datasets include various indoor scenes (*e.g.*, classroom, office, living room, kitchen) to justify the validity and applicability of the assumption.

- 1) For every pair of adjacent pixels in the y-direction of an image, we calculate the angle between the horizon and the vector connecting the two 3D points

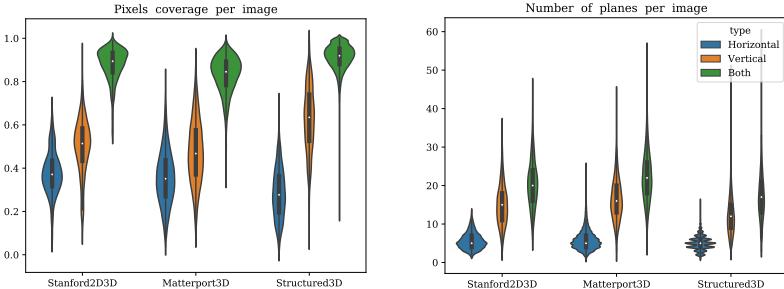


Fig. 3: The violin plot illustrates the statistics of the constructed HV-plane benchmark. We show all the combinations between the three panoramic datasets (X-axis) and the type of planes (blue bars for H-planes, orange bars for V-planes, and green bars for both). The left canvas shows the per-image pixel coverage of HV-planes. The Y-axis is the percentage of covered pixels with respect to valid pixels (having depth). The right canvas shows the number of HV-planes per image

that correspond to those two neighboring pixels. The histogram of angles computed from all images on the three indoor panoramic datasets is shown in Fig. 2. Imagine that we take a vertical ‘slice’ of the 3D point cloud and move vertically between the adjacent points along with the slice—Most of the moving directions will be either horizontal (0° - 10°) or vertical (80° - 90°). Such a cue suggests the gravity aligned nature of indoor scene structures.

- 2) We approximate the scenes with our H-planes and V-planes extraction algorithms, which will be described later. The results of statistics over indoor panoramic datasets are shown in Fig. 3. We can see that, in general, more than 80% of the pixels in an image can be covered by roughly 20 HV-planes. This result suggests that our HV-plane approximation derived from gravity aligned planes are suitable for modeling the gist of an indoor scene.

3.3 Panoramic Image Alignment

We assume all images are aligned with the gravity direction. In case that g-sensor and tripod are not equipped with the 360° camera, and the image is not aligned, we can use the voting-based algorithm mentioned in [8, 23, 25, 29] for panoramic image alignment and vanishing point (VP) detection. Panoramas generally provide enough context to extract the gravity direction, and we will not lose any pixel or introduce any padding after image alignment.

3.4 Ground-truth HV-planes Preparation

Following the gravity aligned and HV-plane assumptions, we present the first 360° planar dataset with training, validation, and test sets. The annotations include HV-plane masks and plane parameters, and all images and annotations are

in the same resolution of 512×1024 . We use the same local analysis as in Fig. 2 to classify each pixel into H-pixel, V-pixel, or other. RANSAC is then performed on H-pixels and V-pixels to extract instance mask and plane geometry for H-planes and V-planes, respectively. The statistical information of the constructed dataset is depicted in Fig. 3. Please refer to supplementary material for detailed ground truth extracting algorithm.

4 Approach

Our task is to reconstruct HV-planes from a single 360° RGB image. We first segment an image into H-planar, V-planar, and non-planar regions (Sec. 4.1). Based on the HV-planar segmentation and 360° image geometry, pixel-level depth map, and V-planar orientation with the novel yaw-invariant representation are transformed into planes geometry (Sec. 4.2). A geometry-aware plane instance segmentation is proposed to exploit the aforementioned estimation to first group pixels with similar planar orientations, and then identify individual plane instances within each group by pixel embedding (Sec. 4.3). An overview of our approach is depicted in Fig. 4. We use a shared backbone with multiple heads for different prediction targets.

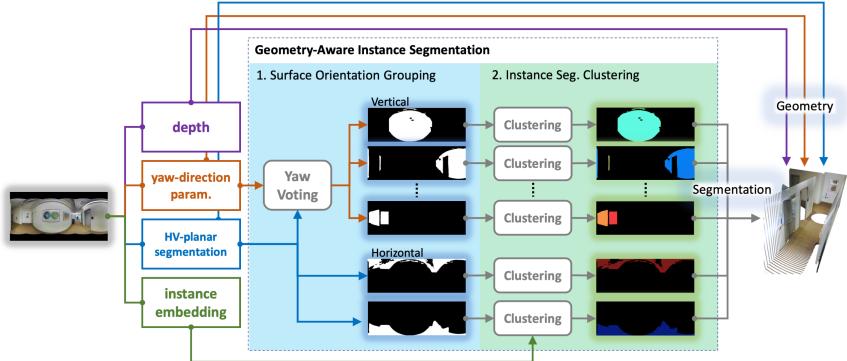


Fig. 4: The pipeline of the proposed YawP³ framework. The deep model predicts pixel-level geometry information, HV-planar segmentation, and instance embedding. In the proposed geometry-aware plane instance segmentation process, we first group pixels with a shared orientation and then identify individual planes from each group via instance embedding

4.1 HV-Planar Segmentation

To distinguish the H-planar and V-planar regions, we predict each pixel as ‘H-planar,’ ‘V-planar,’ or ‘Non-planar’ with two binary classifiers. Both classifiers

are trained to be activated on pixels of H-planes and V-planes respectively, with Binary Cross Entropy loss. For inference, a pixel is recognized as ‘Non-planar’ when the corresponding probabilities from both classifiers are below a certain threshold (we set to 0.5 in this work); otherwise, it will be classified as ‘H-planar’ or ‘V-planar’ according to which classifier gives a higher probability. Based on the segmented H-planar pixels and V-planar pixels, we exploit their prior to give different treatments in the following plane reconstruction process.

4.2 Planar Geometry Estimation

We train the model to estimate two types of pixel-level geometry information—depth d and V-planar surface normal in radian θ . Two heads are attached to the final layer of the backbone for per-pixel d, θ estimation. In below, a planar geometry is denoted as $\vec{n} = [x \ y \ z]$, which is the unit surface normal multiplied by the plane offset; image coordinate is denoted as $u \in [-\pi, \pi]$ and $v \in [-\frac{\pi}{2}, \frac{\pi}{2}]$. Note that the geometry re-projection described below is pixel-level where we use subscript to denote a pixel.

H-planar geometry. The unit surface normal of an H-plane is either $[0 \ 0 \ 1]$ or $[0 \ 0 \ -1]$ (corresponding to horizontal plane above or below the camera respectively) and can be determined accordingly as the pixel is located at the upper half or the bottom half of an equirectangular image. The H-plane offset of a pixel with index i can be derived from d_i by

$$z_i = d_i \cdot \sin(v_i).$$

The training loss for H-planes is $L_H = \frac{1}{|I_H|} \sum_{i \in I_H} |z_i - z_i^*|$ where I_H is the set of all H-planar indices and z_i^* is the ground truth H-plane offset of the pixel. In testing phase, the offset of a detected H-plane instance mask M is determined by the median of $\{z_i \mid i \in I_M\}$ where I_M is all the indices of the instance mask.

V-planar geometry. We derive the V-planar unit surface normal from the estimated θ as $[\cos \theta_i \ \sin \theta_i \ 0]$. Following the manner of H-planar geometry, the V-planes offset is derived by

$$o_i = d_i \cdot \cos(v_i) \cdot [\cos \theta_i \ \sin \theta_i] \cdot [\cos(u_i) \ \sin(u_i)]^T.$$

Thus, $\vec{n}_i = o_i \cdot [\cos \theta_i \ \sin \theta_i \ 0]$ is the geometry representation of the V-plane. The V-planar loss is $L_V = \frac{1}{|I_V|} \sum_{i \in I_V} \|\vec{n}_i - \vec{n}_i^*\|$, where \vec{n}_i^* is the ground truth V-planar normal vector. We also train θ directly using cosine similarity loss $L_\theta = 1 - [\cos \theta_i \ \sin \theta_i] \cdot [\cos \theta_i^* \ \sin \theta_i^*]^T$ where θ_i^* is derived from \vec{n}_i^* . Similar to the testing phase of H-plane geometry, we take the dimensional wise median of \vec{n} in a V-plane instance region.

Yaw-invariant parameterization for V-planes. The accuracy of the estimated V-planar orientation is critical in our framework. It affects the geometry quality and also relates to the proposed geometry-aware plane instance segmentation process (Sec. 4.3). However, V-planar orientations estimation in a 360° image

315 is a challenging task for CNN layers. The reason is that the V-planar surface
 316 orientation co-varies with the 360° camera yaw-rotation (which corresponds to
 317 left-right circular shifting on the equirectangular image), but the CNN layers are
 318 less aware of it due to the translation invariant property. Adding u -coordinate to
 319 the input image as an additional channel could be a workaround. To achieve better
 320 performance, we propose re-parameterizing θ_i into residual form with respect
 321 to the pixel yaw viewing angle u_i such that it is invariant to the 360° camera
 322 yaw-rotation. More specifically, $\theta'_i = \theta_i^* - u_i = \arctan2(y_i^*, x_i^*) - u_i$ is the re-
 323 parameterized yaw-invariant V-planar orientation representation. The proposed
 324 representation enables the model to infer the V-planar orientation without the
 325 knowledge about the 360° camera yaw-rotation.

328 4.3 Geometry-Aware Plane Instance Segmentation

331 Motivated by the strong prior of indoor scenes where most plane instances share
 332 a small number of distinct orientations, we propose to integrate plane surface
 333 information into the main process of plane instance segmentation. Our geometry-
 334 aware plane instance segmentation consists of two stages: *i*) aggregate pixels
 335 by surface orientation grouping; *ii*) in each group, a pixel embedding clustering
 336 algorithm is conducted to identify individual plane instance.

337 **Stage 1: Surface orientation grouping.** The analysis on our dataset
 338 (see Fig. 5) reveals that most plane instances in an indoor scene share similar
 339 planar orientations with other planes. Hence, the per-pixel V-planar surface
 340 orientations are distributed primarily around a small number of angles. Utilizing
 341 such regularity, we divide pixels of an image into groups of similar plane normals,
 342 which preliminarily separates plane instances of dissimilar orientation. For H-
 343 planes, the two surface orientation groups (*i.e.* $\vec{n} = (0, 0, \pm 1)$) can be easily
 344 determined by the HV-planar segmentation mask (Sec. 4.1). Specifically, pixels
 345 classified as H-planar are assigned to two groups according as they are located
 346 above or below the equator (corresponding to H-planes above or below the
 347 camera center). For V-planes, we apply a voting process to detect the prominent
 348 V-planes surface orientation peaks. Namely, we quantize the estimated V-planar
 349 orientations of pixels classified as V-planar into circular bins following by a
 350 peak-finding algorithm (see supplementary material for detail). V-planar pixels
 351 are assigned to the nearest peak to form surface orientation groups.

353 **Stage 2: Pixel embedding clustering.** To further identify the plane in-
 354 stances in each surface orientation group, we exploit pixel embedding, which is
 355 widely used for instance segmentation [3, 7, 10] and also found to be effective for
 356 plane instances segmentation [24]. We follow the loss in [24] and train the model
 357 to pull the pixels of a plane instance toward their centroid in the embedding
 358 space and push the centroids of different planes away from each other. In testing
 359 phase, We apply the mean shift clustering to each orientation group separately.

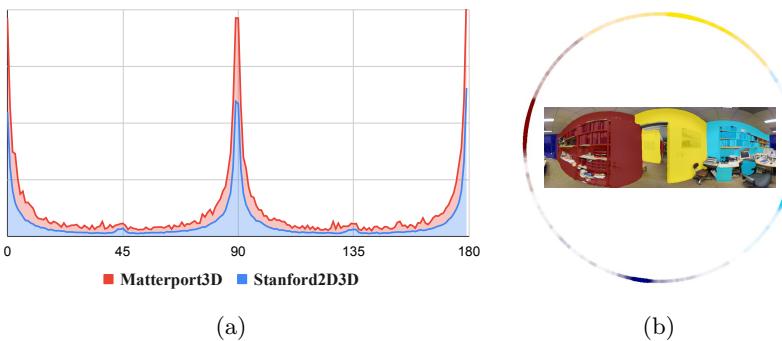


Fig. 5: (a) Statistic of angles between all pairs of planes in an image. The results are averaged across the entire Stanford2D3d and subsampled Matterport3D. (b) One result of surface orientation grouping. Different colors denote different groups. The outer circle represents the distribution of yaw-rotations, where each point on the circle is a vote by a pixel

5 Experiments

5.1 Baselines Construction

We manage to adapt two competitive planar reconstruction approaches – PlanarReconstruct [24]¹ and PlaneRCNN [12]² – to 360° H-FoV panoramic images based on their official implementation. The following is the description of the general changes for both baselines and the respective adaptations.

Common adaptation. The x -axis of a perspective image is linear, but in our 360° H-FoV panorama case, the x -axis is circular. We exploit the left-right circular padding [19] for all CNN layers with 360° H-FoV features as input. The u -coordinate is concatenated as one of the input channels to alleviate the ambiguity of yaw-rotation on 360° images. All methods, including ours, use Resnet101-FPN [11] as the backbone for a fair comparison.

PlaneRCNN [12]. We carefully update RoIAlign and NMS in Mask R-CNN to deal with the left-right circular coordinate system. To refine a b-box from an anchor box, the left-right movement can be ambiguous, so we train the model by choosing the direction with the minimum shifting. We have normal clusters for both vertical and horizontal planes. The resolution of the segmentation refinement module is scaled according to the resolution of the original input.

PlanarReconstruct [24]. We duplicate all heads for H-planes and V-planes. We enforce the plane parameters $\vec{n}_x = 0, \vec{n}_y = 0$ for H-planes, and $\vec{n}_z = 0$ for V-planes by ignoring the prediction of the corresponding output channels.

¹ <https://github.com/svip-lab/PlanarReconstruct>

² <https://github.com/NVLabs/planercnn>

405 Hyperparameters (*i.e.*, number of sampled points and number of iterations) of
 406 the Efficient Mean Shift are scaled according to the change of image resolution.
 407

408 5.2 Data Split

410 We follow the official setting to split the scenes into training, validation, and
 411 testing sets, but remove data with too many missing pixels in depth value
 412 while extracting the ground-truth HV-planes (refer to Sec. 3 for ground truth
 413 extraction detail). Finally, Stanford2D3D [2] contains 1,040 images for training
 414 and 372 images for validation; Matterport3D [4] contains 7,275, 1,189, and 1,005
 415 for training/validation/test respectively; Structured3D [26] contains 18,332 for
 416 training, 1,771 for validation, and 1,691 for testing.

417 418 5.3 Evaluation Metrics

419 Following previous works [12, 21, 24], we evaluate the performance of plane
 420 instances segmentation with some common clustering metrics [1]: Adjusted Rand
 421 Index (ARI↑), Variation of Information (VI↓), and Segmentation Covering (SC↑)
 422 which only consider the segmentation result on the 2D image. To evaluate 3D
 423 reconstruction quality, plane and pixel recall are used under different geometric
 424 thresholds and a segmentation criteria. We report the results by averaging the
 425 plane recall and pixel recall under depth threshold of 5cm, 10cm, 20cm, 30cm
 426 and 60cm.

427 428 5.4 Results

429 **Quantitative evaluation.** In Table 1, we report comparisons between our
 430 approach and two competitive baselines on three presented 360° datasets. The
 431 performance of Stanford2D3D [2] is shown on the validation set as it lacks test
 432 set. Our method achieves state-of-the-art performance consistently on all metrics.
 433 For the metrics that only relate to segmentation quality, our method shows more
 434 improvement on Matterport3D, which contains more complex scenes captured
 435 in luxury houses comparing to the performance gain on Stanford2D3D with
 436 mainly official buildings and the synthetic Structured3D dataset. For metrics
 437 considering 3D reconstruction quality, our approach outperforms both baselines
 438 by a large margin. The results in Table 1 clearly demonstrate the effectiveness of
 439 our proposed method.

440 **Qualitative results.** In Fig. 7, we show the 3D models reconstructed by
 441 PlaneRCNN [12], PlanarReconstruct [24], and the proposed YawP³. Owing to
 442 our yaw-invariant representation, the geometric quality of YawP³ is generally
 443 better than the others. In addition, the planes reconstructed by ours are aligned
 444 better with each other. In contrast, using a detection-based method to predict
 445 segmentation mask and plane normal separately for each instance, PlaneRCNN
 446 has larger gaps between plane instances, which is especially obvious for the large
 447 planes in the first, second, and sixth examples. In general, we achieve higher
 448 quality in both geometry and plane instance segmentation.

450
451 Table 1: The performance of our method and the two baselines reported on the
452 three datasets

Method	Segmentation Quality			Per-pixel recall ↑	Per-plane recall ↑
	ARI↑	VI↓	SC↑		
Matterport3D [4] test set					
P.RCNN [12]	0.574	2.022	0.632	0.473	0.336
P.Recon. [24]	0.673	1.944	0.640	0.498	0.381
Ours	0.686	1.894	0.660	0.544	0.410
Stanford2D3D [2] validation set					
P.RCNN [12]	0.682	1.677	0.703	0.452	0.297
P.Recon. [24]	0.765	1.536	0.733	0.520	0.341
Ours	0.768	1.514	0.742	0.627	0.430
Structured3D [26] test set					
P.RCNN [12]	0.726	1.393	0.743	0.654	0.522
P.Recon. [24]	0.821	1.175	0.785	0.728	0.591
Ours	0.824	1.150	0.794	0.794	0.657

469 470 5.5 Ablation Study

471 In this section, several ablation studies are shown to further prove the effectiveness
472 of the proposed method and identify the source of improvement.

473 **Can the proposed 360° yaw-invariant plane parameters benefit other
474 baselines?** As mentioned in Sec. 4.2, one challenge of 360° planar orientation
475 estimation is that the non-horizontal normals are related to camera yaw rotation,
476 but the CNN layers are less aware of the counterpart 2D left-right circular shifting
477 on the equirectangular image. In addition to the input u -coordinate workaround,
478 we propose a 360° camera yaw-invariant representation for V-plane normals,
479 which ensures the independence between the ground truth normal and camera
480 yaw rotation (see Sec. 4.2 for detail). The proposed representation is applied to
481 other baselines, and the results are reported in Table 2. We observe consistent
482 improvements on all metrics by applying the proposal on PlaneRCNN [12]. PlanarReconstruct
483 [24] also achieves better 3D reconstruction quality with a slightly
484 degraded 2D segmentation. Please note that even by applying the proposed
485 yaw-invariant to other baselines, our method still achieves the best performance,
486 which suggests that the proposed plane representation is not the only source for
487 our superior results.

488
489 **Can the proposed planar orientation parameterization along resist the
490 yaw ambiguity of 360° camera?** To show the yaw ambiguity problem of
491 360° images and how the proposed yaw-invariant representation resists to it, we
492 conduct an ablation study on our method with all the combinations of input
493 u -coordinate and the proposed representation. In Table 3, we show the per-pixel

495
496 Table 2: Ablation study of the proposed yaw-invariant parameterization on the
497 two baselines
498

Method	yaw-invariant	Segmentation Quality			Per-pixel recall ↑	Per-plane recall ↑
		ARI↑	VI↓	SC↑		
Matterport3D [4] validation set						
P.RCNN	✓	0.531	2.118	0.613	0.471	0.313
		0.571	2.053	0.632	0.495	0.326
P.Recon	✓	0.634	2.102	0.615	0.484	0.362
ours	✓	0.625	2.171	0.602	0.495	0.369
		0.654	2.044	0.640	0.545	0.398
Stanford2D3D [2] validation set						
P.RCNN	✓	0.682	1.677	0.703	0.452	0.297
		0.683	1.665	0.705	0.481	0.313
P.Recon	✓	0.766	1.536	0.733	0.520	0.341
ours	✓	0.763	1.537	0.728	0.555	0.368
		0.768	1.514	0.742	0.627	0.430

513
514
515 normal error in degree over the V-planar region of Stanford2D3D [2] validation set.
516 As mentioned in Sec. 4.2, V-plane normals depend on 360° camera yaw-rotation,
517 but CNN layers are less aware of it. Consequently, the model with neither input
518 u -coordinate nor the proposed yaw-invariant representation leads to inferior V-
519 planar geometry quality. On the other hand, our yaw-invariant reparameterization,
520 which is designed to be invariant to the 360° camera yaw-rotation, can achieve
521 superior orientation quality.
522

523
524
525 Table 3: Ablation study of the V-plane normal error under different settings. A
526 row with the first column checked indicates that the model has an extra input
527 of image u -coordinate. A checkmark on the second column means the proposed
528 yaw-invariant representation is used. Other than the two factors, all models share
529 the same setting and are trained on Stanford2D3D [2] training set. The results
530 are evaluated on Stanford2D3D validation set
531

input u -coord.	yaw-invariant	V-plane orientation error (deg°)↓
		7.29
✓		5.97
	✓	5.18
✓	✓	5.21

The effectiveness of the surface orientation grouping stage. To examine the impact of the surface orientation grouping stage, we modify PlanarReconstruction [24], the baseline which only adopts pixel embedding, by integrating our proposed yaw-invariant representation into their method. In Table 2, the comparison between the results of ours and the modified PlanarReconstruction [24] suggests that that the geometry-aware plane instance segmentation apart from yaw-invariant representation also improves the performance. We show some representative qualitative comparisons in Fig 6.

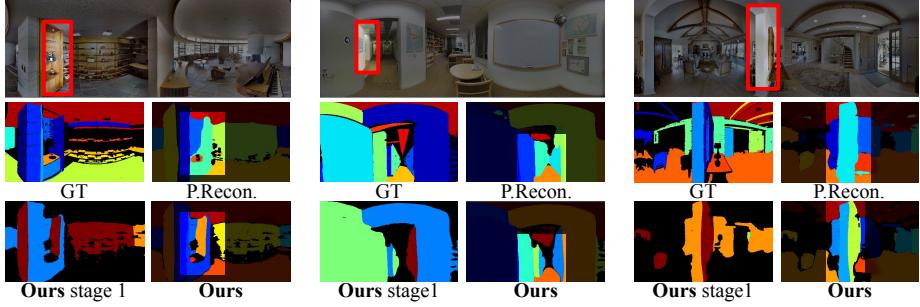


Fig. 6: Three representative examples to demonstrate the effectiveness of the proposed geometry-aware plane instance segmentation. In comparison with the baseline considering only pixel embedding, which fails in the highlighted regions, our geometry-aware plane instance segmentation can separate the undetected plane from its neighboring planes by the first stage orientation voting. Please see Sec. 4.3 for the intuition and the technical detail

6 Conclusion and Future Work

This paper presents YawP³, a novel deep neural network targeting indoor 3D planar reconstruction from a single panoramic image with HV-planes. YawP³ applies surface orientation grouping to panoramic images, which improves the instance segmentation module through geometric awareness. The proposed yaw-invariant V-plane parameter of YawP³ is capable of both solving surface normal ambiguity of panoramic images and boosting the performance of reconstruction. We created the first 360° planar reconstruction benchmark with datasets, and baseline models for experimentation and comparison. The presented planar datasets are derived from large-scale RGB-D datasets and the baseline models are modified from representative perspective view planar models to adapt to 360° data. Our model outperforms the two baselines by a large margin and achieves state-of-the-art on this newly presented benchmark. Two future directions that are worth investigating are *i*) finding better approaches for reconstructing HV-planes and *ii*) completing occluded and non-estimated regions in HV-planes.

585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629

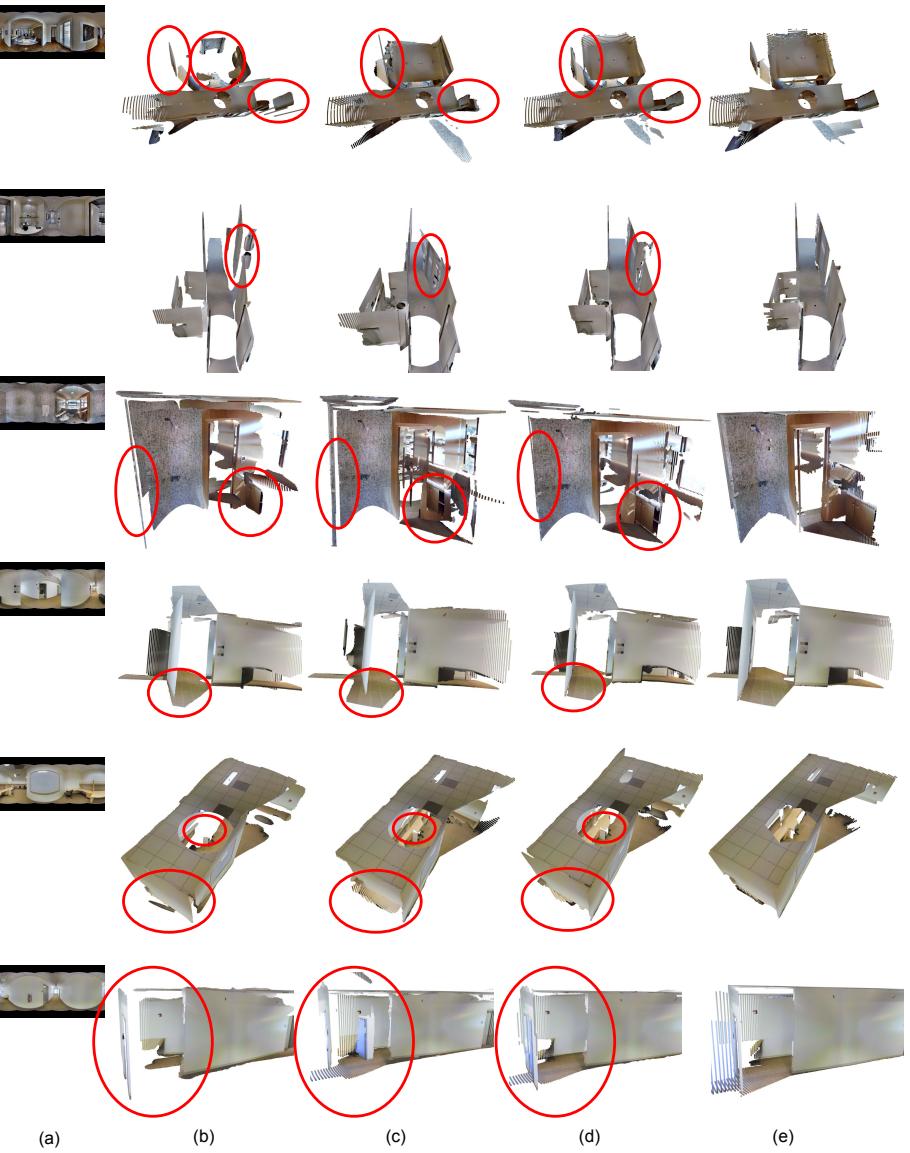


Fig. 7: We show some qualitative comparisons of the 3D reconstruction results among the three methods. We highlight the differences in red circles for easier to distinguish the reconstruction quality. (a) Input RGB, (b) PlaneRCNN [12], (c) PlanarReconstruction [24], (d) ours, and (e) ground truth

630 References

- 632 1. Arbelaez, P., Maire, M., Fowlkes, C.C., Malik, J.: Contour detection and hierarchical
image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.* **33**(5), 898–916 (2011)
- 633 2. Armeni, I., Sax, S., Zamir, A.R., Savarese, S.: Joint 2d-3d-semantic data for indoor
scene understanding. *CoRR* **abs/1702.01105** (2017)
- 634 3. Brabandere, B.D., Neven, D., Gool, L.V.: Semantic instance segmentation with a
discriminative loss function. *CoRR* **abs/1708.02551** (2017)
- 635 4. Chang, A.X., Dai, A., Funkhouser, T.A., Halber, M., Nießner, M., Savva, M.,
Song, S., Zeng, A., Zhang, Y.: Matterport3D: learning from RGB-D data in indoor
environments. In: 2017 International Conference on 3D Vision, 3DV 2017, Qingdao,
China, October 10-12, 2017. pp. 667–676 (2017)
- 636 5. Dai, A., Chang, A.X., Savva, M., Halber, M., Funkhouser, T.A., Nießner, M.:
Scannet: Richly-annotated 3d reconstructions of indoor scenes. In: 2017 IEEE
Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu,
HI, USA, July 21-26, 2017. pp. 2432–2443 (2017)
- 637 6. Eder, M., Moulon, P., Guan, L.: Pano popups: Indoor 3d reconstruction with a
plane-aware network. In: 2019 International Conference on 3D Vision, 3DV 2019,
Québec City, QC, Canada, September 16-19, 2019. pp. 76–84 (2019)
- 638 7. Fathi, A., Wojna, Z., Rathod, V., Wang, P., Song, H.O., Guadarrama, S., Mur-
phy, K.P.: Semantic instance segmentation via deep metric learning. *CoRR*
abs/1703.10277 (2017)
- 639 8. Fernandez-Labrador, C., Pérez-Yus, A., López-Nicolás, G., Guerrero, J.J.: Layouts
from panoramic images with geometry and deep learning. *IEEE Robotics and
Automation Letters* **3**(4), 3153–3160 (2018)
- 640 9. He, K., Gkioxari, G., Dollár, P., Girshick, R.B.: Mask R-CNN. In: IEEE International
Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017.
pp. 2980–2988 (2017)
- 641 10. Kong, S., Fowlkes, C.C.: Recurrent pixel embedding for instance grouping. In: 2018
IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt
Lake City, UT, USA, June 18-22, 2018. pp. 9018–9028 (2018)
- 642 11. Lin, T., Dollár, P., Girshick, R.B., He, K., Hariharan, B., Belongie, S.J.: Feature
pyramid networks for object detection. In: 2017 IEEE Conference on Computer
Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017.
pp. 936–944 (2017)
- 643 12. Liu, C., Kim, K., Gu, J., Furukawa, Y., Kautz, J.: PlaneRCNN: 3d plane detection
and reconstruction from a single image. In: IEEE Conference on Computer Vision
and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019.
pp. 4450–4459 (2019)
- 644 13. Liu, C., Yang, J., Ceylan, D., Yumer, E., Furukawa, Y.: PlaneNet: piece-wise planar
reconstruction from a single RGB image. In: 2018 IEEE Conference on Computer
Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22,
2018. pp. 2579–2588 (2018)
- 645 14. Newell, A., Huang, Z., Deng, J.: Associative embedding: End-to-end learning
for joint detection and grouping. In: Advances in Neural Information Processing
Systems 30: Annual Conference on Neural Information Processing Systems 2017,
4-9 December 2017, Long Beach, CA, USA. pp. 2277–2287 (2017)
- 646 15. Purkait, P., Zach, C., Leonardis, A.: Rolling shutter correction in manhattan world.
In: IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy,
October 22-29, 2017. pp. 882–890 (2017)

- 675 16. Silberman, N., Hoiem, D., Kohli, P., Fergus, R.: Indoor segmentation and support
676 inference from RGBD images. In: Computer Vision - ECCV 2012 - 12th European
677 Conference on Computer Vision, Florence, Italy, October 7-13, 2012, Proceedings,
678 Part V. pp. 746–760 (2012) 675
679 17. Song, S., Zeng, A., Chang, A.X., Savva, M., Savarese, S., Funkhouser, T.A.:
680 Im2pano3d: Extrapolating 360° structure and semantics beyond the field of view.
681 In: 2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR
682 2018, Salt Lake City, UT, USA, June 18-22, 2018. pp. 3847–3856 (2018) 679
683 18. Sun, C., Hsiao, C., Sun, M., Chen, H.: HorizonNet: learning room layout with
684 1d representation and pano stretch data augmentation. In: IEEE Conference on
685 Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA,
686 June 16-20, 2019. pp. 1047–1056 (2019) 683
687 19. Wang, T., Huang, H., Lin, J., Hu, C., Zeng, K., Sun, M.: Omnidirectional CNN for
688 visual place recognition and navigation. In: 2018 IEEE International Conference on
689 Robotics and Automation, ICRA 2018, Brisbane, Australia, May 21-25, 2018. pp.
690 2341–2348 (2018) 688
691 20. Xu, J., Stenger, B., Kerola, T., Tung, T.: Pano2cad: Room layout from a single
692 panorama image. In: 2017 IEEE Winter Conference on Applications of Computer
693 Vision, WACV 2017, Santa Rosa, CA, USA, March 24-31, 2017. pp. 354–362 (2017) 691
694 21. Yang, F., Zhou, Z.: Recovering 3d planes from a single image via convolutional
695 neural networks. In: Computer Vision - ECCV 2018 - 15th European Conference,
696 Munich, Germany, September 8-14, 2018, Proceedings, Part X. pp. 87–103 (2018) 694
697 22. Yang, H., Zhang, H.: Efficient 3d room shape recovery from a single panorama. In:
698 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016,
699 Las Vegas, NV, USA, June 27-30, 2016. pp. 5422–5430 (2016) 697
700 23. Yang, Y., Jin, S., Liu, R., Kang, S.B., Yu, J.: Automatic 3d indoor scene modeling
701 from single panorama. In: 2018 IEEE Conference on Computer Vision and Pattern
702 Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018. pp. 3926–3934
703 (2018) 700
704 24. Yu, Z., Zheng, J., Lian, D., Zhou, Z., Gao, S.: Single-image piece-wise planar 3d
705 reconstruction via associative embedding. In: IEEE Conference on Computer Vision
706 and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019.
707 pp. 1029–1037 (2019) 704
708 25. Zhang, Y., Song, S., Tan, P., Xiao, J.: PanoContext: A whole-room 3d context
709 model for panoramic scene understanding. In: Computer Vision - ECCV 2014 -
710 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings,
711 Part VI. pp. 668–686 (2014) 709
712 26. Zheng, J., Zhang, J., Li, J., Tang, R., Gao, S., Zhou, Z.: Structured3D: A large
713 photo-realistic dataset for structured 3d modeling. CoRR [abs/1908.00222](https://arxiv.org/abs/1908.00222) (2019)
714 27. Zheng, S., Jayasumana, S., Romera-Paredes, B., Vineet, V., Su, Z., Du, D., Huang,
715 C., Torr, P.H.S.: Conditional random fields as recurrent neural networks. In: 2015
716 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile,
717 December 7-13, 2015. pp. 1529–1537. IEEE Computer Society (2015) 714
718 28. Zhou, Y., Qi, H., Zhai, Y., Sun, Q., Chen, Z., Wei, L., Ma, Y.: Learning to
719 reconstruct 3d manhattan wireframes from a single image (2019) 718
720 29. Zou, C., Colburn, A., Shan, Q., Hoiem, D.: LayoutNet: reconstructing the 3d room
721 layout from a single RGB image. In: 2018 IEEE Conference on Computer Vision
722 and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018.
723 pp. 2051–2059 (2018) 721