# FloorNet: A Unified Framework for Floorplan Reconstruction from 3D Scans

Chen Liu[*]        Jiaye Wu[*]              Yasutaka Furukawa
Washington University in St. Louis        Simon Fraser University
{chenliu,jiaye.wu}@wustl.edu              furukawa@sfu.ca

**Abstract.** The ultimate goal of this indoor mapping research is to automatically reconstruct a floorplan simply by walking through a house with a smartphone in a pocket. This paper tackles this problem by proposing FloorNet, a novel deep neural architecture. The challenge lies in the processing of RGBD streams spanning a large 3D space. FloorNet effectively processes the data through three neural network branches: 1) PointNet with 3D points, exploiting the 3D information; 2) CNN with a 2D point density image in a top-down view, enhancing the local spatial reasoning; and 3) CNN with RGB images, utilizing the full image information. FloorNet exchanges intermediate features across the branches to exploit the best of all the architectures. We have created a benchmark for floorplan reconstruction by acquiring RGBD video streams for 155 residential houses or apartments with Google Tango phones and annotating complete floorplan information. Our qualitative and quantitative evaluations demonstrate that the fusion of three branches effectively improves the reconstruction quality. We hope that the paper together with the benchmark will be an important step towards solving a challenging vector-graphics reconstruction problem. Code and data are available at https://github.com/art-programmer/FloorNet.

**Keywords:** Floorplan Reconstruction; 3D Computer Vision; 3D CNN

## 1   Introduction

Architectural floorplans play a crucial role in designing, understanding, and remodeling indoor spaces. Their drawings are effective in conveying geometric and semantic information of a scene. For instance, we can quickly identify room extents, the locations of doors, or object arrangements (geometry). We can also recognize the types of rooms, doors, or objects easily through texts or icon styles (semantics). Unfortunately, more than 90% of houses in North America do not have floorplans. The ultimate goal of the indoor mapping research is to enable

---

[*]The first two authors contribute equally on this work.

automatic reconstruction of a floorplan simply by walking through a house with a smartphone in a pocket.

The Consumer-grade depth sensors have revolutionized indoor 3D scanning with successful products. Matterport [1] produces detailed texture mapped models of indoor spaces by acquiring a set of panorama RGBD images with a specialized hardware. Google Project Tango phones [2] convert RGBD image streams into 3D or 2D models. These systems produce detailed geometry, but fall short as floorplans or architectural blue-prints, whose geometry must be concise and respect underlying scene segmentation and semantics.

Reconstruction of the floorplan for an entire house or an apartment with multiple rooms poses fundamental challenges to existing techniques due to its large 3D extent. A standard approach projects 3D information onto a 2D lateral domain [3], losing the information of height. PointNet [4,5] consumes 3D information directly but suffers from the lack of local neighborhood structures. A multi-view representation [6,7] avoids explicit 3D space modeling, but has been mostly demonstrated for objects, rather than large scenes and complex camera motions. 3D Convolutional Neural Networks (CNNs) [8,9] also show promising results but have been so far limited to objects or small-scale scenes.

This paper proposes a novel deep neural network (DNN) architecture Floor-Net, which turns a RGBD video covering a large 3D space into pixel-wise predictions on floorplan geometry and semantics, followed by an existing Integer Programming formulation [10] to recover vector-graphics floorplans. FloorNet consists of three DNN branches. The first branch employs PointNet with 3D points, exploiting the 3D information. The second branch uses a CNN with a 2D point density image in a top-down floorplan view, enhancing the local spatial reasoning. The third branch uses a CNN with RGB images, utilizing the full image information. The PointNet branch and the point-density branch exchange features between the 3D points and their corresponding cells in the top-down view. The image branch contributes deep image features into the corresponding cells in the top-down view. This hybrid DNN design exploits the best of all the architectures and effectively processes the full RGBD video covering a large 3D scene with complex camera motions.

We have created a benchmark for floorplan reconstruction by acquiring RGBD video streams for 155 residential houses or apartments with Google Tango phones and annotated their complete floorplan information including architectural structures, icons, and room types. Extensive qualitative and quantitative evaluates demonstrate the effectiveness of our approach over competing methods.

In summary, the main contributions of this paper are two-fold: 1) Novel hybrid DNN architecture for RGBD videos, which processes the 3D coordinates directly, models local spatial structures in the 2D domain, and incorporates the full image information; and 2) A new floorplan reconstruction benchmark with RGBD videos, where many indoor scene databases exist [11,12,1] but none tackles a vector-graphics reconstruction problem, which has immediate impact on digital mapping, real estate, or civil engineering applications.

## 2 Related work

We discuss the related work in three domains: indoor scene reconstruction, 3D deep learning, and indoor scan datasets.

**Indoor scene reconstruction:** The advancements in consumer-grade depth sensors have brought revolutionary changes to indoor 3D scanning. KinectFusion [13] enables high-fidelity 3D scanning for objects and small-scale scenes. Whelan et al. [14] extends the work to building-scale scans. While being accurate with details, these dense reconstructions fall short as CAD models, which must have 1) concise geometry for efficient data transmission and 2) proper segmentations/semantics for architectural analysis or effective visualization.

Towards CAD-quality reconstructions, researchers have applied model-based approaches by representing a scene with geometric primitives. Utilizing the 2.5D property of indoor building structures, rooms can be separated by fitting lines to points in a top-down view [15,16,17]. Primitive types have been extended to planes [18,19,20,21,22] or cuboids [23]. While they produce promising results for selected scans, they critically rely on the low-level geometry analysis for primitive detection, which faces challenges with noisy and incomplete 3D data. Our approach conducts global analysis of the entire input by DNNs to detect primitive structures much more robustly.

Another line of research studies the top-down scene reconstruction with shape grammars from a single image [24] or a set of panorama RGBD images [3,25]. Crowdsensing data such as images and WiFi-fingerprints are also exploited in building scene graphs [26,27,28,29]. While semantic segmentation [11,4,5] and scene understanding [30] are popular for indoor scenes, there has been no robust learning-based method for vector-graphics floorplan reconstruction. This paper provides such a method and its benchmark with the ground-truth.

One way to recover the mentioned vector-graphics floorplan models is from rasterized floorplan images [10]. We share the same reconstruction target, and we utilize their Integer Programming formulation in our last step to recover the final floorplan. Nevertheless, instead of a single image as input, our input is a RGBD video covering a large 3D space, which requires a fundamentally different approach to process the input data effectively.

**3D deep learning:** The success of CNN on 2D images has inspired research on 3D feature learning via DNNs. Volumetric CNNs [31,32,6] are straightforward extensions of CNN to a 3D domain, but there are two main challenges: 1) data sparsity and 2) computational cost of 3D convolutions. FPNN [33] and Vote3D [34] attempt to solve the first challenge, while OctNet [8] and O-CNN [35] address the computational costs via octree representations.

2D CNNs with multi-view renderings have been successful for object recognition [6,7] and part segmentation [36]. They effectively utilize all the image information but are so far limited to regular (or fixed) camera arrangements. The extension to larger scenes with complex camera motions is not trivial.

PointNet [4] directly uses 3D point coordinates to exploit the sparsity and avoid quantization errors, but it does not provide an explicit local spatial reason-

ing. PointNet++ [5] hierarchically groups points and adds spatial structures by farthest point sampling. Kd-Networks [37] similarly group points by a KD-tree. These techniques incur additional computational expenses due to the grouping and have been limited at object-scale. For scenes, they need to split the space into smaller regions (e.g., 1m×1m blocks) and process each region independently [4,5], potentially hurting global reasoning (e.g., identifying long walls for corridors or avoiding two kitchens for an apartment).

**Indoor scan dataset:** Affordable depth sensing hardware enables researchers to build many indoor scan datasets. The ETH3D dataset contains only 16 indoor scans [38], and its purpose is for multi-view stereo rather than 3D point-cloud processing. The ScanNet dataset [11] and the SceneNN dataset [39] capture a variety of indoor scenes. However, most of their scans contain only one or two rooms, not suitable for the floorplan reconstruction problem.

Matterport3D [40] builds high quality panorama RGBD image sets for 90 luxurious houses. 2D-3D-S dataset [41] provides 6 large-scale indoor scans of office spaces by using the same Matterport camera. However, they focus on 2D and 3D semantic annotations, and do not address a vector-graphics reconstruction problem. Meanwhile, they require an expensive specialized hardware (i.e., Matterport camera) for high-fidelity 3D scanning, while we aim to tackle the challenge by consumer-grade smartphones with low data quality.

Lastly, a large-scale synthetic dataset, SUNCG [42], offers a variety of indoor scenes with CAD-quality geometry and annotations. However, they are synthetic and cannot model the complexity of real scenes or replace the real photographs. We provide the benchmark with full floorplan annotations and the corresponding RGBD videos from smartphones for 155 residential units.

## 3 FloorNet

The proposed FloorNet converts a RGBD video with camera poses into pixel-wise floorplan geometry and semantics information, which is an intermediate floorplan representation introduced by Liu et al. [10]. We first explain the intermediate representation for being self-contained, then provide the details.

### 3.1 Preliminaries

The intermediate representation consists of the geometry and the semantics information. The geometry part contains room-corners, object icon-corners, and door/window end-points, where the locations of each corner/point type are estimated by a 256×256 heatmap in the 2D floorplan image domain, followed by a standard non-maximum suppression. For example, a room corner is either I-, L-, T-, or X-shaped depending on the number of incident walls, making the total number of feature maps to be 13 considering their rotational variants. The semantics part is modeled as 1) 12 feature maps as a probability distribution function (PDF) over 12 room types, and 2) 8 feature maps as a PDF over 8 icon
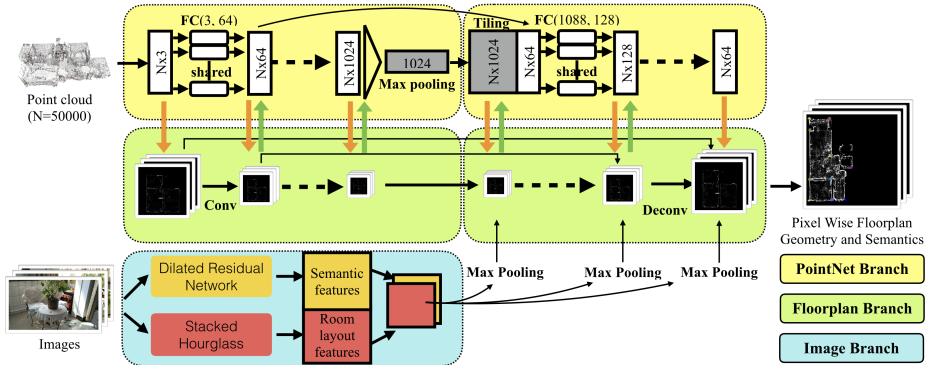
**Fig. 1.** FloorNet consists of three DNN branches. The first branch uses PointNet [4] to directly consume 3D information. The second branch takes a top-down point density image in a floorplan domain with a fully convolutional network [43], and produces pixel-wise geometry and semantics information. The third branch produces deep image features by a dilated residual network trained on the semantic segmentation task [44] as well as a stacked hourglass CNN trained on the room layout estimation [45]. The PointNet branch and the floorplan branch exchanges intermediate features at every layer, while the image branch contributes deep image features into the decoding part of the floorplan branch. This hybrid DNN architecture effectively processes an input RGBD video with camera poses, covering a large 3D space.

types. We follow their approach and use their Integer Programming formulation to reconstruct a floorplan from this representation at the end.

## 3.2 Triple-branch hybrid design

Floornet consists of three DNN branches. We employ existing DNN architectures in each branch without modifications. Our contribution lies in its hybrid design: how to combine them and share intermediate features (See Fig. 1).

**PointNet Branch:** The first branch is PointNet [4] that directly takes 3D points, where we use the original architecture without modifications. We randomly subsample 50,000 points for each data. We manually rectify the rotation before annotation, while aligning the gravity direction with the Z-axis. We add translation to move the center of mass to the origin.

**Floorplan Branch:** The second branch is a fully convolutional network (FCN) [43] with skip connections between the encoder and the decoder, which takes a point-density image in the top-down view. We compute a 2D axis-aligned bounding box of the Manhattan-rectified 3D points to define a rectangular floorplan domain, while ignoring the 2.5% outlier points and expanding the rectangle by 5% in each of the four directions. The rectangle is placed in the middle of the $256 \times 256$ square image in which the geometry and semantics feature maps are produced. The input to the branch is a point-density image in the same domain.
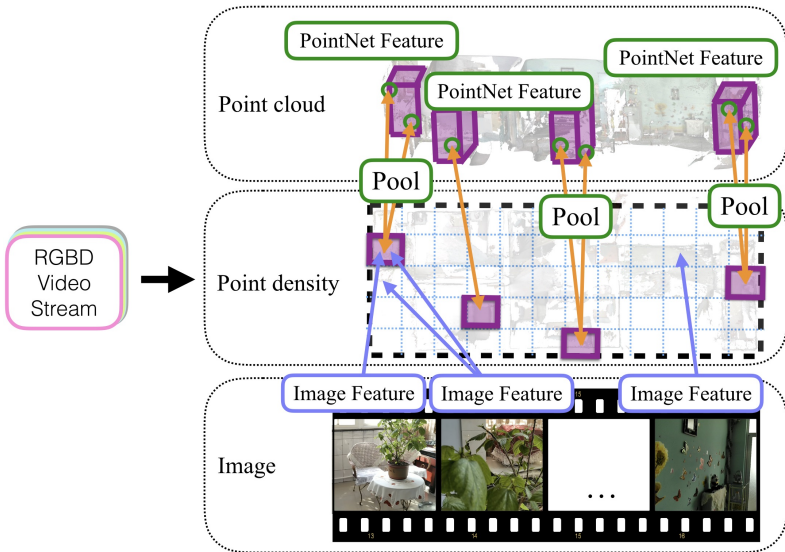
**Fig. 2.** FloorNet shares features across branches to exploit the best of all the architectures. PointNet features at 3D points are pooled into corresponding 2D cells in the floorplan branch. Floorplan features at 2D cells are unpooled to the corresponding 3D points in the PointNet branch. Deep image features are pooled into corresponding 2D cells in the floorplan branch, based on the depthmap and the camera pose information.

**Image Branch:** The third branch computes deep image features through two CNN architectures: 1) Dilated residual network (DRN) [44] trained on semantic segmentation; and 2) stacked hourglass CNN (HG) [45] trained on room layout estimation. We have used ScanNet [11] benchmark to train DRN and LSUN [46] benchmark to train HG.

### 3.3 Intra-branch feature sharing

Different branches learn features in different domains (3D points, the floorplan, and images). FloorNet offers three intra-branch feature sharing by pooling and unpooling operations, based on the camera poses and 3D information (See Fig. 2).

**PointNet to floorplan pooling:** This pooling module takes features of unordered points from each layer of the PointNet branch and produces a 2D top-down feature map in the corresponding layer of the floorplan branch. The module simply spreads point features into cells defined by the output top-down feature map, then computes the sum of the features in each cell. Though the projection could be performed along each of the three axes, we focus on the vertical projection since our goal is to reconstruct a floorplan. A constructed feature map has the same dimension as the layer in the floorplan branch, and is simply con-

catenated to the current feature stack. The time complexity of the projection pooling module is linear in the number of 3D points.

**Floorplan to pointNet unpooling:** This module reverses the above pooling operation. It simply copies and adds a feature of the floorplan cell into each of the corresponding 3D points that project inside the cell. The time complexity is again linear in the number of points.

**Image to floorplan pooling:** The image branch produces two deep image features of dimensions 512x32x32 and 256x64x64 from DRN and HG for each video frame. We first unproject image features to 3D space by their depthmaps and camera poses, and then apply the same 3D to floorplan pooling operation to aggregate 3D features to the floorplan branch. We conduct the image branch pooling for every 10 frames in the video sequence.

### 3.4 Loss functions

Our network outputs pixel-wise predictions on the floorplan geometry and semantics information in the same resolution $256 \times 256$. For geometry heatmaps (i.e., room corners, object icon-corners, and door/window end-points), a sigmoid cross entropy loss is used. The ground-truth heatmap is prepared by putting a value of 1.0 inside a disk of radius 11 pixels around each ground-truth pixel. For semantic classification feature maps (i.e., room types and object icon types), a pixel-wise softmax cross entropy loss is used.

## 4 Floorplan reconstruction benchmark

This paper creates a benchmark for the vector-graphics floorplan reconstruction problem from RGBD videos with camera poses. We have acquired roughly two-hundreds 3D scans of residential units in the United States and China using Google Tango phones (Lenovo Phab 2 Pro and Asus ZenFone AR) (See Fig. 3). After manually removing poor quality scans, we have annotated the complete floorplan information for the remaining 155 scans: 1) room-corners as points, 2) walls as pairs of room-corners, 3) object icons and types as axis-aligned rectangles and classification labels, 4) doors and windows (i.e., openings) as line-segments on walls, and 5) room types as classification labels for polygonal areas enclosed by walls. The list of object types is {*counter, bathtub, toilet, sink, sofa, cabinet, bed, table, refrigerator*}. The list of room types is {*living room, kitchen, bedroom, bathroom, closet, balcony, corridor, dining room*}. Table 1 provides statistics of our data collections.

Reconstructed floorplans are evaluated on three different levels of geometric and semantic consistency with the ground-truth. We follow the work by Liu et al. [10] and define the low- and mid-level metrics as follows.

• The low-level metric is the precision and recall of room corner detections. A corner detection is declared a success if its distance to the ground-truth is below 10 pixels and the closest among all the other room corners.
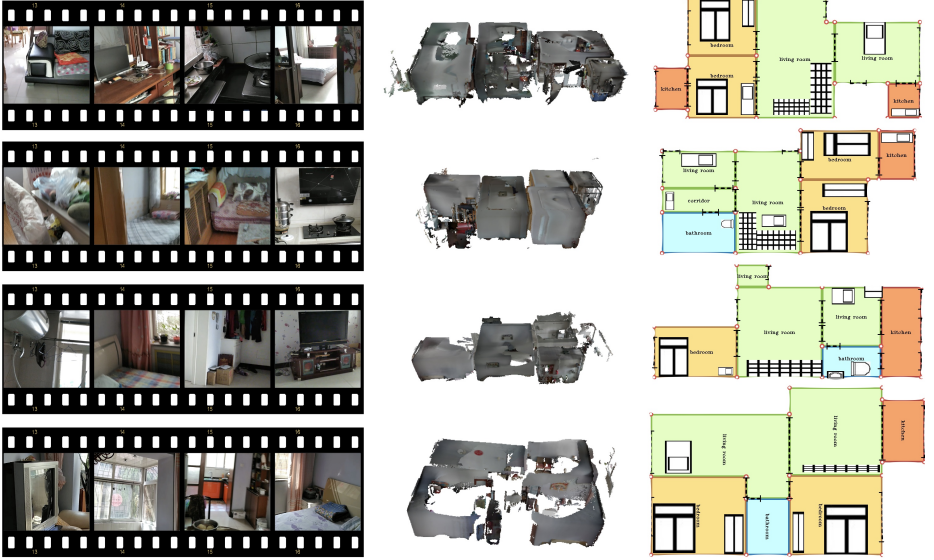
**Fig. 3.** Floorplan reconstruction benchmark. From left to right: Subsampled video frames, colored 3D point clouds, and ground-truth floorplan data. The floorplan data is stored in a vector-graphics representation, which is visualized with a simple rendering engine (e.g., rooms are assigned different colors based on their types, and objects are shown as canonical icons).

• The mid-level metric is the precision and recall of detected openings (i.e., doors and windows), object-icons, and rooms. The detection of an opening is declared a success if the largest distance of the corresponding end-points is less than 10 pixels. The detection of an object (resp. a room) is declared a success if the intersection-over-union (IOU) with the ground-truth is above 0.5 (resp. 0.7).

• Relationships of architectural components play crucial roles in evaluating indoor spaces. For example, one may look for apartments where bedrooms are not connected to a kitchen. A building code may enforce every bedroom to have a quick evacuation route to outside through windows or doors in case of fire. We introduce the high-level metric as the ratio of rooms that have the correct relationships with the neighboring rooms. More precisely, we declare that a room has a correct relationship if 1) a room is connected to the correct set of rooms through doors, where two rooms are connected if their common walls contain at least one door, 2) each room (i.e., the room and its neighboring rooms) has an IOU score larger than 0.5 with the corresponding ground-truth, and 3) each room has the correct room type.

**Table 1.** Dataset statistics. From left to right: the number of rooms, the number of icons, the number of openings (i.e., doors or windows), the number of room-corners, and the total area. The average and the standard deviation are reported for each entry.

|  | #room | #icon | #opening | #corner | area |
|---|---|---|---|---|---|
| Average | 5.2 | 9.1 | 9.9 | 18.1 | $63.8[m^2]$ |
| Std | 1.8 | 4.5 | 2.9 | 4.2 | $13.0[m^2]$ |

# 5 Implementation details

## 5.1 DNN Training

Among the 155 scans we collected, we randomly sample 135 for training and leave 20 for testing. We perform data augmentation by random scaling and rotation every time we feed a training sample. First, we apply rescaling to the point-cloud and the corresponding annotation with a random factor uniformly sampled from a range $[0.5, 1.5]$. Second, we randomly apply the rotation around the z axis by either $0^o$, $90^o$, $180^o$, or $270^o$.

We use the official code for each DNN module, that is, PointNet [4], FCN [43] for the Floorplan branch, and DRN [44] and SH [45] for the Image branch. We pre-train DRN on the semantic segmentation task with ScanNet database [11] and SH on the room layout estimation task with LSUN [46]. DRN and SH are fixed during the FloorNet training, and we optimize only the PointNet and the floorplan branches.

FloorNet has three types of loss functions. We have observed that enabling all the loss functions in a single training would lead to worse testing performance, because the model overfits with the icon loss. Instead, we train the model with each loss function one by one. When using the icon loss, we limit the training to be at most 600 iterations and use early-stopping based on the testing loss, where 1 iteration contains 20 batches. [1]

Training of FloorNet takes around 2 hours with a TitanX GPU. On the average, the training continues for about 600 iterations, consuming $1,620,000 = 135(\text{samples}) \times 600(\text{iterations}) \times 20(\text{batches})$ augmented training samples. It is initially to our surprise that FloorNet generalizes even from a small number of 3D scans. However, FloorNet makes pixel-wise predictions, which are mostly low-level vision tasks. Each 3D scan contains about 10 object-icons, 10 openings, and a few dozen room corners, which probably lead to the good generalization performance together with data augmentation, where similar phenomena were observed by Chen et al. [10]

---

[1]  We incorporated synthetic dataset SUNCG [42] and/or real dataset Matterport3D [40] for training with the icon loss, while using their semantic segmentation information to produce icon annotations. However, the joint-training still experiences overfitting, while this simply early-stopping heuristic works well in our experiments.
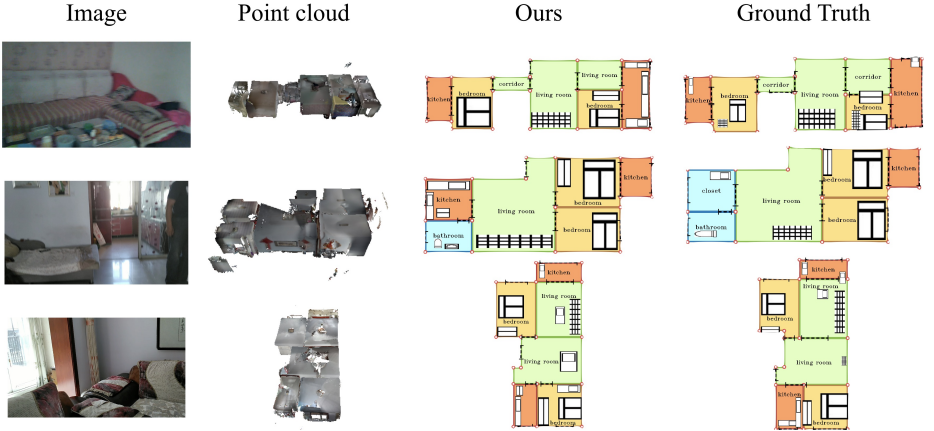
| Image | Point cloud | Ours | Ground Truth |
|-------|-------------|------|--------------|



**Fig. 4.** Floorplan reconstruction results.

## 5.2 Enhancement heuristics

We augment the Integer Programming Formulation [10] with the following two enhancement heuristics to deal with more challenging input data (i.e., large-scale raw sensor data) and hence more noise in the network predictions.

**Primitive candidate generation:** Standard non-maximum suppression often detects multiple room corners around a single ground-truth. After thresholding the room-corner heatmap by a value 0.5, we simply extract the highest peak from each connected component, whose area is more than 5 pixels. To handle localization errors, we connect two room-corners and generate a wall candidate when their corresponding connected components overlap along X or Y direction. We do not augment junctions to keep the number of candidates tractable.

**Objective function:** Wall and opening candidates are originally assigned uniform weights in the objective function [10]. We calculate the confidence of a wall (resp. opening) candidate by taking the average of the semantic heatmap scores of type "wall" along the line with width 7 pixels (resp. 5 pixels). We set the weight of each primitive by the confidence score minus 0.5, so that a primitive is encouraged to be chosen only when the confidence is at least 0.5.

## 6 Experiments

Figure 4 shows our reconstruction results on some representative examples. Our approach successfully recovers complex vector-graphics floorplan data including room geometries and their connectivities through doors. One of the major failure modes is in the icon detection. As the model training experiences overfitting in the icon loss, object detection generally requires more training data than low-level geometry (i.e., corners) detection [10]. We believe that more training data
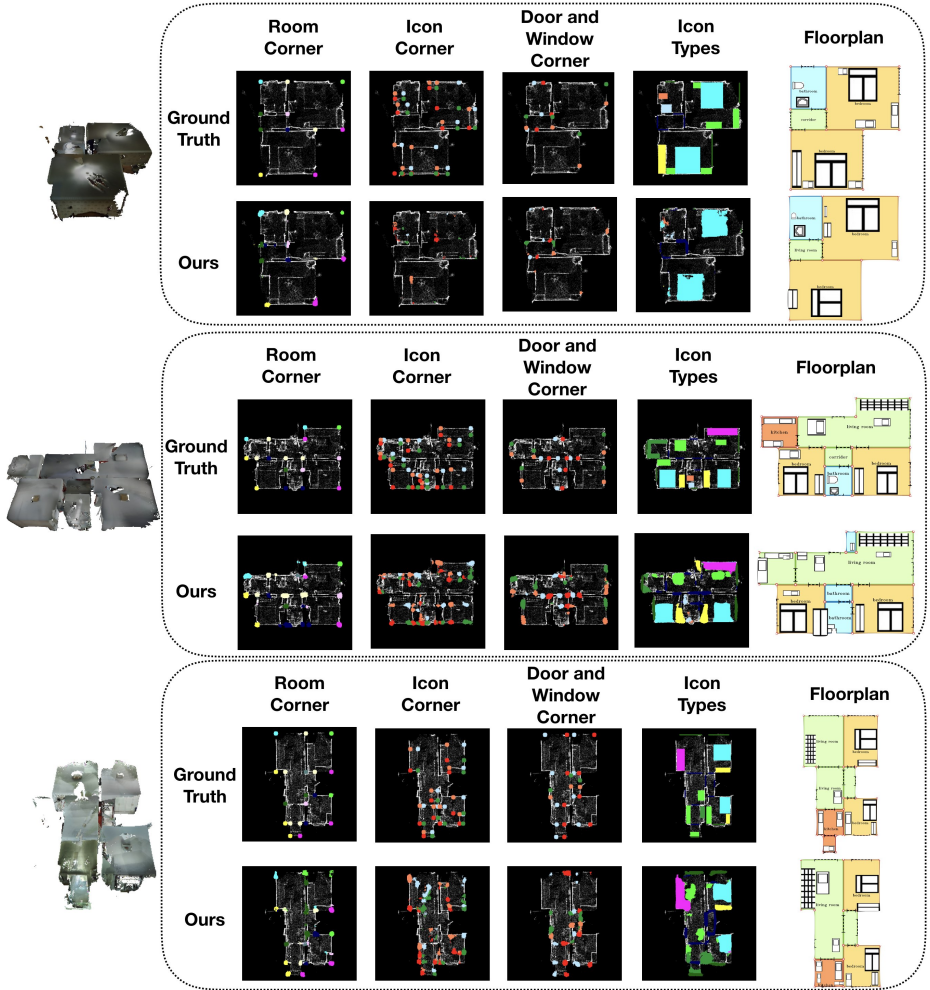
**Fig. 5.** Intermediate results. For each example, we show raw outputs of the networks (room corners, icon corners, opening corners, and icon types) compared against the ground-truth. In the second example, we produce a fake room (blue color) at the top due to poor quality 3D points. In the third example, reconstructed rooms have inaccurate shapes near the bottom left again due to noisy 3D points, illustrating the challenge of our problem.
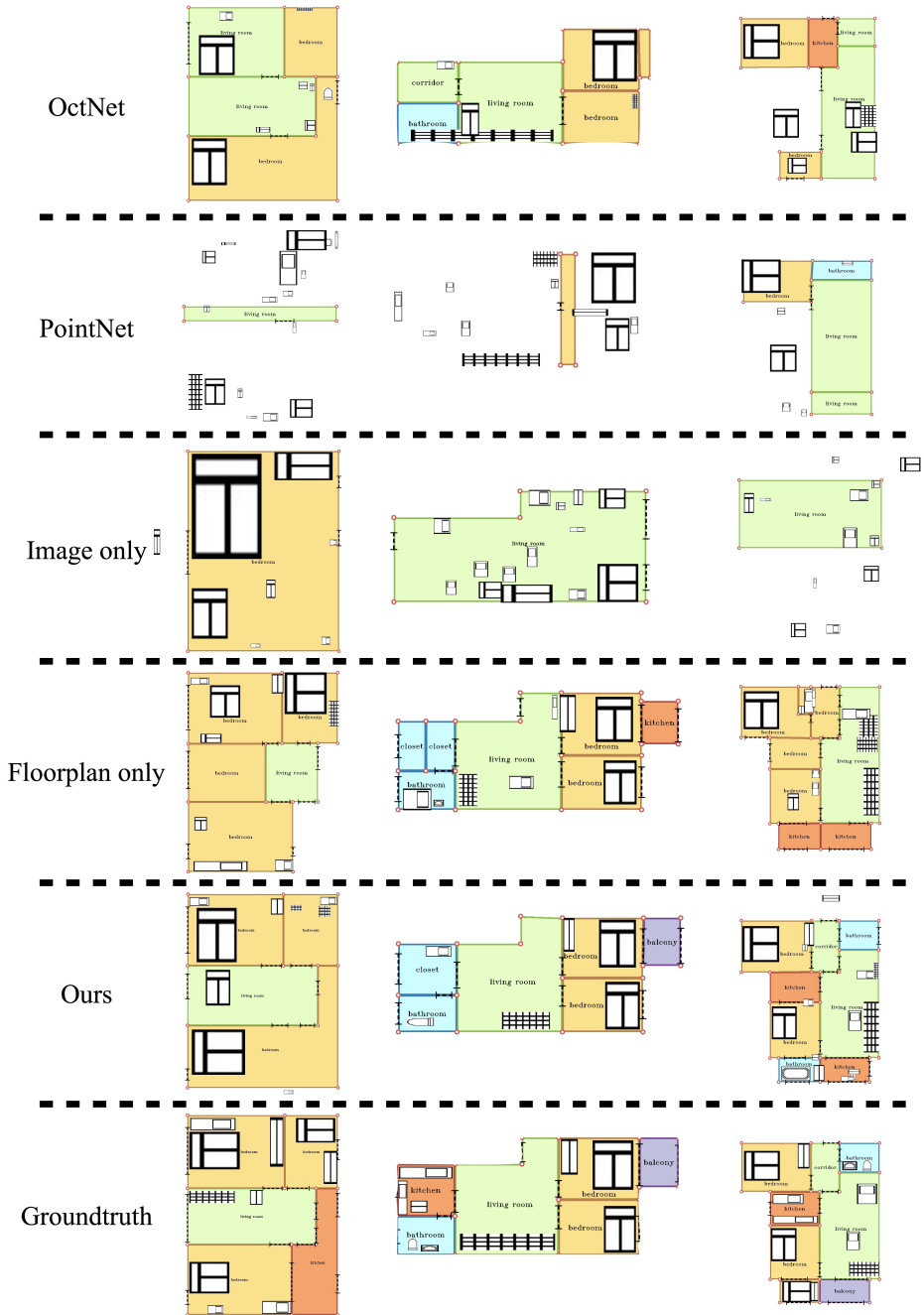
12



**Fig. 6.** Qualitative comparisons against competing methods. The top is OctNet [8], a state-of-the-art 3D CNN architecture. The next three rows show variants of our FloorNet, where only one branch is enabled. FloorNet with all the branches overall produce more complete and accurate floorplans.

**Table 2.** Quantitative evaluations on low-, mid-, and high-level metrics against competing methods and our variants. The orange and cyan color indicates the best and the second best result for each entry.

| | wall | door | icon | room | relationship |
|---|---|---|---|---|---|
| PointNet [4] | 25.8/42.5 | 11.5/38.7 | 22.5/27.9 | 27.0/40.2 | 5.0 |
| Floorplan-branch | 90.2/88.7 | 70.5/78.0 | 43.4/42.8 | 76.3/75.3 | 50.0 |
| Image-branch | 40.0/83.3 | 15.4/47.1 | 21.4/17.4 | 25.0/57.1 | 0.0 |
| OctNet [8] | 75.4/89.2 | 36.6/82.3 | 32.8/48.8 | 62.1/72.0 | 13.5 |
| Ours w/o PointNet-Unpooling | 92.6/92.1 | 75.8/76.8 | 55.1/51.9 | 80.9/77.4 | 52.3 |
| Ours w/o PointNet-Pooling | 88.4/93.0 | 73.0/87.2 | 50.0/42.2 | 75.0/80.6 | 52.8 |
| Ours w/o Image-Pooling | 92.6/89.7 | 77.1/74.4 | 50.5/57.8 | 84.2/83.1 | 56.8 |
| Ours | 92.1/92.8 | 76.7/80.2 | 56.1/57.8 | 83.6/85.2 | 56.8 |

will overcome this issue. Another typical failures come from missing room corners due to clutter or incomplete scanning. The successful reconstruction of a room requires successful detection of every single room corner. This is a challenging problem and the introduction of higher level constraints may reveal a solution (e.g., if one sees a door, a room must be reconstructed on the other side even with missing corners).

Figure 6 and Table 2 qualitatively and quantitatively compare our method against competing techniques, namely, OctNet [8], PointNet [4], and a few variants of our FloorNet. OctNet and PointNet represent state-of-the-art 3D DNN techniques. More precisely, we implement the voxel semantic segmentation network based on the official OctNet library, [2] which takes 256x256x256 voxels as input and outputs 3D voxels of the same resolution. We then add three separate $5 \times 3 \times 3$ convolution layers with strides $4 \times 1 \times 1$ to predict the same pixel-wise geometry and semantics feature-maps with the same set of loss functions. PointNet is simply our FloorNet without the point density or the image input. Similarly, we construct a FloorNet variant by enabling only the 3D points (for the PointNet branch) or the point density image (for the floorplan branch) as the input.

The table shows that the floorplan branch is the most informative as it is the most natural representation for floorplan reconstruction task, while PointNet branch or Image branch alone does not work well. We also split the entire point clouds into $1m \times 1m$ blocks, train the PointNet-only model that makes predictions per block separately, followed by a simple merging. However, this performs much worse. OctNet performs reasonably well across low- to mid-level metrics, but does poorly on the high-level metric, where all the rooms and relevant doors must be reconstructed at high precision to report good numbers.

To further evaluate the effectiveness of the proposed FloorNet architecture, we conduct ablation studies by disabling each of the intra-branch pooling/unpooling

---

[2] OctNet library: https://github.com/griegler/octnet

operations. The bottom of Table 2 shows that the feature sharing overall leads to better results, especially for mid- to high-level metrics.

Finally, Figure 7 compares against a build-in Tango Navigator App [47], which generates a floorplan image real-time on the phone. Note that their system does not 1) produce room segmentations, 2) recognize room types, 3) detect objects, 4) recognize object types, or 5) produce CAD-quality geometry. Therefore, we quantitatively evaluate only the geometry information by measuring the line distances between the ground-truth walls and predicted walls. More precisely, we 1) sample 100 points from each wall line segment, 2) for each sampled point, find the closest one in the other line segment, and 3) compute the mean distance over all the sampled points and line segments. The average line distances are 2.72 [pixels] and 1.66 [pixels] for Tango Navigator App and our FloorNet, respectively. This is a surprisingly result to some extent, because our algorithm drops many confident line segments during Integer Programming, when their corresponding rooms miss one corner and are not reconstructed. On the other hand, it is an expected result as our approach makes full use of geometry and image information to recover a floorplan.



**Fig. 7.** Comparison against a commercial floorplan generator, Tango Navigator App. Top: Floorplan image from Tango. Bottom: Our results.

## 7 Conclusion

This paper proposes a novel DNN architecture FloorNet that reconstructs vector-graphics floorplans from RGBD videos with camera poses. FloorNet takes a hybrid approach and exploits the best of three DNN architectures to effectively process a RGBD video covering a large 3D space with complex camera motions. The paper also provides a new benchmark for a new vector-graphics reconstruction problem, which is missing in the recent indoor scene databases of Computer

Vision. Two main future works are ahead of us. The first one is to learn to enforce higher level constraints inside DNNs as opposed to inside a separate post-processing (e.g., Integer Programming). Learning high-level constraints likely require more training data and the second future work is to acquire more scans.

More than 90% of houses in North America do not have floorplans. We hope that this paper together with the benchmark will be an important step towards solving this challenging vector-graphics reconstruction problem, and enabling the reconstruction of a floorplan just by walking through a house with a smartphone. We will publicly share our code and data to promote further research.

## Acknowledgement

# References

1. : Matterport. https://matterport.com/
2. Lee, J., Dugan, R., et al.: Google project tango
3. Ikehata, S., Yang, H., Furukawa, Y.: Structured indoor modeling. In: Proceedings of the IEEE International Conference on Computer Vision. (2015) 1323–1331
4. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. arXiv preprint arXiv:1612.00593 (2016)
5. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In: Advances in Neural Information Processing Systems. (2017) 5105–5114
6. Qi, C.R., Su, H., Nießner, M., Dai, A., Yan, M., Guibas, L.J.: Volumetric and multi-view cnns for object classification on 3d data. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2016) 5648–5656
7. Su, H., Maji, S., Kalogerakis, E., Learned-Miller, E.: Multi-view convolutional neural networks for 3d shape recognition. In: Proceedings of the IEEE international conference on computer vision. (2015) 945–953
8. Riegler, G., Ulusoys, A.O., Geiger, A.: Octnet: Learning deep 3d representations at high resolutions. arXiv preprint arXiv:1611.05009 (2016)
9. Tatarchenko, M., Dosovitskiy, A., Brox, T.: Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. arXiv preprint arXiv:1703.09438 (2017)
10. Liu, C., Wu, J., Kohli, P., Furukawa, Y.: Raster-to-vector: Revisiting floorplan transformation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2017) 2195–2203
11. Dai, A., Chang, A.X., Savva, M., Halber, M., Funkhouser, T., Nießner, M.: Scannet: Richly-annotated 3d reconstructions of indoor scenes. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR). Volume 1. (2017)
12. Song, S., Yu, F., Zeng, A., Chang, A.X., Savva, M., Funkhouser, T.: Semantic scene completion from a single depth image. IEEE Conference on Computer Vision and Pattern Recognition (2017)
13. Newcombe, R.A., Izadi, S., Hilliges, O., Molyneaux, D., Kim, D., Davison, A.J., Kohi, P., Shotton, J., Hodges, S., Fitzgibbon, A.: Kinectfusion: Real-time dense surface mapping and tracking. In: Mixed and augmented reality (ISMAR), 2011 10th IEEE international symposium on, IEEE (2011) 127–136
14. Whelan, T., Kaess, M., Fallon, M., Johannsson, H., Leonard, J., McDonald, J.: Kintinuous: Spatially extended kinectfusion. (2012)
15. Okorn, B., Xiong, X., Akinci, B., Huber, D.: Toward automated modeling of floor plans. In: Proceedings of the symposium on 3D data processing, visualization and transmission. Volume 2. (2010)
16. Turner, E., Cheng, P., Zakhor, A.: Fast, automated, scalable generation of textured 3d models of indoor environments. IEEE Journal of Selected Topics in Signal Processing **9**(3) (2015) 409–421
17. Sui, W., Wang, L., Fan, B., Xiao, H., Wu, H., Pan, C.: Layer-wise floorplan extraction for automatic urban building reconstruction. IEEE transactions on visualization and computer graphics **22**(3) (2016) 1261–1277
18. Furukawa, Y., Curless, B., Seitz, S.M., Szeliski, R.: Manhattan-world stereo. In: Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on, IEEE (2009) 1422–1429

19. Furukawa, Y., Curless, B., Seitz, S.M., Szeliski, R.: Reconstructing building interiors from images. In: Computer Vision, 2009 IEEE 12th International Conference on, IEEE (2009) 80–87
20. Sinha, S., Steedly, D., Szeliski, R.: Piecewise planar stereo for image-based rendering. (2009)
21. Xiong, X., Adan, A., Akinci, B., Huber, D.: Automatic creation of semantically rich 3d building models from laser scanner data. Automation in Construction **31** (2013) 325–337
22. Mura, C., Mattausch, O., Villanueva, A.J., Gobbetti, E., Pajarola, R.: Automatic room detection and reconstruction in cluttered indoor environments with complex room layouts. Computers & Graphics **44** (2014) 20–32
23. Xiao, J., Furukawa, Y.: Reconstructing the worlds museums. International journal of computer vision **110**(3) (2014) 243–258
24. Zhao, Y., Zhu, S.C.: Image parsing with stochastic scene grammar. In: Advances in Neural Information Processing Systems. (2011) 73–81
25. Mura, C., Mattausch, O., Pajarola, R.: Piecewise-planar reconstruction of multi-room interiors with arbitrary wall arrangements. In: Computer Graphics Forum. Volume 35., Wiley Online Library (2016) 179–188
26. Gao, R., Zhao, M., Ye, T., Ye, F., Wang, Y., Bian, K., Wang, T., Li, X.: Jigsaw: Indoor floor plan reconstruction via mobile crowdsensing. In: Proceedings of the 20th annual international conference on Mobile computing and networking, ACM (2014) 249–260
27. Gao, R., Zhao, M., Ye, T., Ye, F., Luo, G., Wang, Y., Bian, K., Wang, T., Li, X.: Multi-story indoor floor plan reconstruction via mobile crowdsensing. IEEE Transactions on Mobile Computing **15**(6) (2016) 1427–1442
28. Luo, H., Zhao, F., Jiang, M., Ma, H., Zhang, Y.: Constructing an indoor floor plan using crowdsourcing based on magnetic fingerprinting. Sensors **17**(11) (2017) 2678
29. Jiang, Y., Xiang, Y., Pan, X., Li, K., Lv, Q., Dick, R.P., Shang, L., Hannigan, M.: Hallway based automatic indoor floorplan construction using room fingerprints. In: Proceedings of the 2013 ACM international joint conference on Pervasive and ubiquitous computing, ACM (2013) 315–324
30. Zhang, Y., Bai, M., Kohli, P., Izadi, S., Xiao, J.: Deepcontext: Context-encoding neural pathways for 3d holistic scene understanding. arXiv preprint arXiv:1603.04922 (2016)
31. Wu, Z., Song, S., Khosla, A., Yu, F., Zhang, L., Tang, X., Xiao, J.: 3d shapenets: A deep representation for volumetric shapes. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2015) 1912–1920
32. Maturana, D., Scherer, S.: Voxnet: A 3d convolutional neural network for real-time object recognition. In: Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on, IEEE (2015) 922–928
33. Li, Y., Pirk, S., Su, H., Qi, C.R., Guibas, L.J.: Fpnn: Field probing neural networks for 3d data. In: Advances in Neural Information Processing Systems. (2016) 307–315
34. Wang, D.Z., Posner, I.: Voting for voting in online point cloud object detection. In: Robotics: Science and Systems. (2015)
35. Wang, P.S., Liu, Y., Guo, Y.X., Sun, C.Y., Tong, X.: O-cnn: Octree-based convolutional neural networks for 3d shape analysis. ACM Transactions on Graphics (TOG) **36**(4) (2017) 72
36. Limberger, F.A., Wilson, R.C., Aono, M., Audebert, N., Boulch, A., Bustos, B., Giachetti, A., Godil, A., Le Saux, B., Li, B., et al.: Shrec'17 track: point-cloud

shape retrieval of non-rigid toys. In: 10th Eurographics workshop on 3D Object retrieval. (2017) 1–11

37. Klokov, R., Lempitsky, V.: Escape from cells: Deep kd-networks for the recognition of 3d point cloud models. In: 2017 IEEE International Conference on Computer Vision (ICCV), IEEE (2017) 863–872

38. Schöps, T., Schönberger, J.L., Galliani, S., Sattler, T., Schindler, K., Pollefeys, M., Geiger, A.: A multi-view stereo benchmark with high-resolution images and multi-camera videos. In: Proc. CVPR. Volume 3. (2017)

39. Hua, B.S., Pham, Q.H., Nguyen, D.T., Tran, M.K., Yu, L.F., Yeung, S.K.: Scenenn: A scene meshes dataset with annotations. In: 3D Vision (3DV), 2016 Fourth International Conference on, IEEE (2016) 92–101

40. Chang, A., Dai, A., Funkhouser, T., Halber, M., Nießner, M., Savva, M., Song, S., Zeng, A., Zhang, Y.: Matterport3d: Learning from rgb-d data in indoor environments. arXiv preprint arXiv:1709.06158 (2017)

41. Armeni, I., Sener, O., Zamir, A.R., Jiang, H., Brilakis, I., Fischer, M., Savarese, S.: 3d semantic parsing of large-scale indoor spaces. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2016) 1534–1543

42. Song, S., Yu, F., Zeng, A., Chang, A.X., Savva, M., Funkhouser, T.: Semantic scene completion from a single depth image. arXiv preprint arXiv:1611.08974 (2016)

43. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. (2015) 3431–3440

44. Yu, F., Koltun, V., Funkhouser, T.: Dilated residual networks. In: Computer Vision and Pattern Recognition. Volume 1. (2017)

45. Newell, A., Yang, K., Deng, J.: Stacked hourglass networks for human pose estimation. In: European Conference on Computer Vision, Springer (2016) 483–499

46. Zhang, Y., Yu, F., Song, S., Xu, P., Seff, A., Xiao, J.: Large-scale scene understanding challenge: Room layout estimation. accessed on Sep **15** (2015)

47. Inc., G.: Project tango. https://developers.google.com/tango/