

Exploiting 2D Floorplan for Building-scale Panorama RGBD Alignment

Erik Wijmans Yasutaka Furukawa
Washington University in St. Louis
jewijmans@cs.wustl.edu

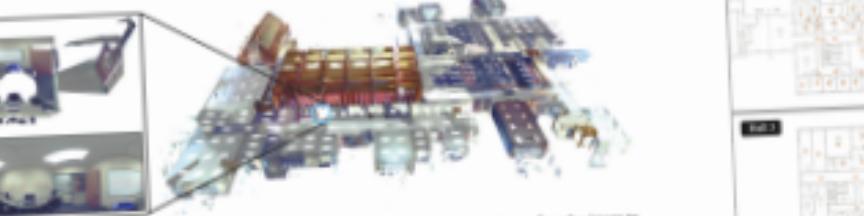


Figure 1: The paper tackles building-scale panorama RGBD image alignment. Our approach utilizes a floorplan image to significantly reduce the number of necessary scans and hence human operating costs.

Abstract

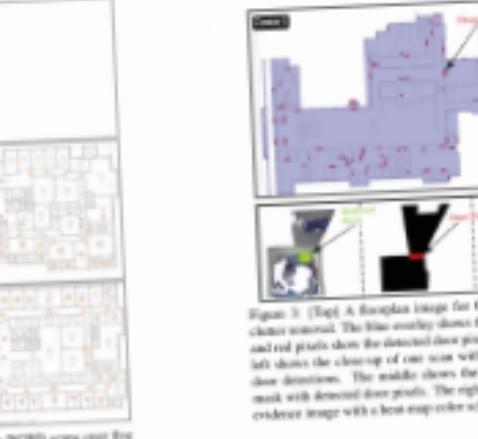
This paper presents a novel algorithm that utilizes a 2D floorplan to align panoramic RGBD scans. While effective panoramic RGBD alignment techniques exist, such a system requires extremely dense RGBD image sampling. Our approach can significantly reduce the number of necessary scans with the aid of a floorplan image. We formulate a novel Multi-Branch Random Field inference problem or a scan placement over the floorplans, as opposed to the conventional scan-to-scan alignment. The technical contributions lie in multi-modal image correspondence over thousands of scans and volumetric floorplans as well as a novel coverage potential avoiding an extensive masking bias. The proposed approach has been evaluated on five challenging large indoor spaces. In the best of our knowledge, we present the first effective system that utilizes a 2D floorplan image for building-scale 3D pointcloud alignment. The source code and the data will be shared with the community to further enhance indoor mapping research.

1. Introduction

3D scanning hardware has made remarkable progress in recent years, where successful products exist in industry for commercial applications. In particular, Panorama RGBD

scanners have found real-world application as the system produces both 3D geometry and immersive panoramic images. For instance, Face 3D [1] is a professional-grade panoramic RGBD scanner, which can reach more than 100 meters and produce 100 million points per scan within a millimeter accuracy. The device is perfect for 3D measurement, documentation, or surveillance in indoor mapping, real engineering or GIS applications. Maptek [3, 4] is an emerging low-end solution that can reach only 5 meters, but it reaches quicker (i.e., 1 to 2 minutes per scan), and has demonstrated compelling results for Real Estate markets.

Given the focus on the 3D scanning hardware, automated panoramic RGBD alignment has become a crucial technology. The Maptek system provides a robust solution but requires extremely dense scanning (e.g., one scan every 2 to 5 meters). Dense scanning becomes infeasible for high-end scanners (e.g., Face 3D [1]), whose single scan could take thirty minutes or an hour depending on the resolution. However, these scanners are the only option for large buildings such as department stores, airport terminals, or hotel lobbies, simply due to the required operating ranges (e.g., 20 to 30 meters). Furthermore, the precision of these high-end scanners is necessary for quantitative recovery of metric information for scientific and engineering data analysis. In practice, with high-end 3D scanning devices, people use calibration objects such as big bright balls and/or utilize



due to acquire building-scale panoramic RGBD scans over five locations. Only one scan has been acquired in each room in (a) and the use of floorplan-image essential for our problem.

2. Related work

Two approaches exist for indoor 3D scanning: “RGBD scanning” or “Panorama RGBD scanning”. RGBD-scanning continuously moves a depth camera and scans a scene. This has been the major choice among Computer Vision researchers [26, 7, 22] after the success of Kinect Fusion [13]. The input is a RGBD video stream, while Simultaneous Localization and Mapping (SLAM) is the core technology. Panorama RGBD scanning has been rather successful in industry, because (i) data acquisition is easy (i.e., picking a 2D position as opposed to 6 DoF navigation in RGBD scanning); (ii) alignment is easier thanks to the panoramic field of views; and (iii) the system produces panoramic images, essential for many visualization applications. Structure from Motion (SfM) is the core technology to this approach. This paper provides an automated solution for Panorama RGBD alignment, and the remainder of the section focuses on the description of the SfM techniques, where we refer the reader to a survey article [23] for the SLAM literature.

Structure from Motion (SfM) addresses the problem of metric information for scientific and engineering data analysis. In practice, with high-end 3D scanning devices, people use calibration objects such as big bright balls and/or utilize

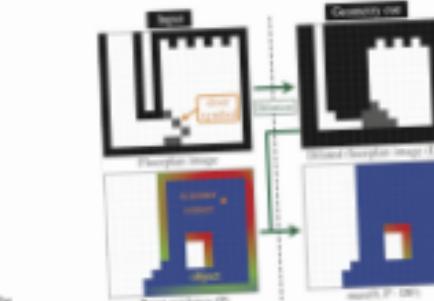


Figure 3: (Top) A floorplan image for Center 1 after the scan is placed. The blue overlay shows the building mask, and red pixels show the detected door pixels. (Bottom) The left shows the close-up of one scan with the result of 3D door detection. The middle shows the free-space image mask with detected door pixels. The right shows the point-evidence image with a heat-map color scheme.

2.1. The average penalty over all the door-pixels in the evidence image is the semantic penalty.

Geometric cue: Measuring the consistency between the floorplan image and the point evidence image is a real challenge: (i) A floorplan image contains extra symbols that are not in evidence images; (ii) An evidence image contains obstacles that are not in a floorplan image; (iii) The size of a floorplan (e.g., line thickness) may vary; and (iv) Both are essentially line-drawings, making the comparison sensitive to small errors. In practice, we have found that the following consistency potential provides a robust metric:

• We detect doors both in a floorplan image and 3D scans. For a floorplan, we randomly specify a bounding box containing a door symbol, and use a standard template matching. For 3D scans, we use a heuristic to identify short-specific 3D patterns directly in 3D points.

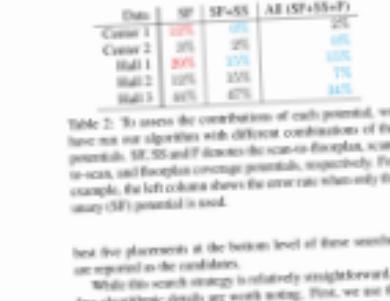
4. MRF formulation

The multi-modal nature of the problem makes our formulation fundamentally different from existing ones [21, 22]. The first critical difference lies in the definition of the variables. In existing approaches, a variable encodes a 3D relative placement between a pair of scans [21, 22]. In our formulation, a variable encodes a 2D absolute placement of a single scan over a floorplan image.

Let $\theta = \{\theta_1, \theta_2, \dots\}$ be our variables, where θ_i encodes the 2D placement of a single scan. θ_i consists of two components: (1) rotation, which takes one of the four angle values (0, 90, 180, or 270 degrees); and (2) translation, which is a pixel coordinate in the floorplan image (P). Then, θ_i defines the location of the scan in the floorplan image ($P - \theta_i$) divided by the size of elements in the original evidence image (P) indicates the extent of the discrepancy. We sweep the role of the floorplan and a point-evidence image, compute the other discrepancy measure, and take the average.

4.2. Scan-to-scan consistency potential

Different from standard MRF formulation, we do not know which pairs of variables (i.e., scans) should have interaction, because our variables encode the placements of the scans. Therefore, we set up a potential for every pair of scans. The potential measures the photometric and geometric



Data	SP	SP+S	AB (SP+SP+F)
Center 1	22%	22%	22%
Center 2	20%	23%	23%
Hall 2	12%	15%	15%
Hall 3	10%	17%	14%

Table 2: To assess the contribution of each potential, we have run our algorithm with different combinations of the potentials. SP, SS and F denotes the scan-on-floorplan, scan-to-scan, and scan-to-scan coverage potentials, respectively. For example, the left column shows the error rate when only the unary (SP) potential is used.



and 3D point cloud, and 2D colored free-space masks are shown.

5. Experimental results and discussions

We have used C++ for implementation and Intel Core i7 CPU with 16GB RAM PC. These computational expensive steps are pre-processing, unary-potential evaluation, and TRW optimization, where the running time is roughly proportional to the number of the input scans. For large datasets with 20 to 30 scans, these steps roughly take 5 hours, 2.2 hours, and 50 minutes, respectively. The pre-processing is the bottleneck due to 3D processing of the floorplan images, which can be parallelized if necessary.

Our method has successfully aligned most of the scans. We have not scanned the right wing of the building in Center 1 and Center 2 (See Fig. 2), which makes a large space for scans to be misplaced. Nonetheless, our method has only one misplacement in that area (Center 1). Note that Hall 3 is an exception in which we make many errors due to the glitch in the unary potential.

Figure 5 illustrates typical failure modes of Damer transformations, which tends to concentrate scans in large rooms. Our analysis is that a large room tends to have non-architectural lines or symbols in 3D-in an empty canvas, which makes the distance transform image contain small values, and allow loose-scans placements.

Figure 6 illustrates the effects of the semantic cue in the unary potential (i.e., the door detection). Failure cases are full of inconsistencies and repetitions, which makes the computation of pure geometry (i.e., geometric cues) susceptible to local minima. The figure demonstrates a representative case, where the door detections break such an ambiguity.

When the placement is ambiguous even with the geometric and the semantic cues, we rely on the MRF optimization with the full three potentials. Figure 7 compares the final scan placements, with or without the floorplan coverage potential. The floorplan coverage potential tends to avoid “stacking” and evenly distribute the placements.

Our method is not perfect and has exposed several failure modes. First, our approach tends to make mistakes for small storage-style rooms, where a small room with a lot of clutter makes the semantic cue very noisy. Second, there are genuinely ambiguous cases where the scene geometry, appearance, and local locations are exactly the same. Lastly, our method has made major errors in Hall 3, simply because the floorplan has not reflected annotations in the past. Unfortunately, it was difficult to identify erroneous scans based on the potentials. As the presence of problematic scans, the MRF optimization seems to shuffle around scan placements including correct ones to achieve a low-energy state. Nonetheless, the total potential, in particular, the magnitude of the total potential divided by the number of scans is a good indicator of success. The quantity for Hall 3 is a few times larger than the others and indicates that “something is wrong”. Our main future work is to develop a robust algorithm to detect potentially erroneous scan placements, which will allow a quick user feedback to correct mistakes. We will share our source code and high-end building-scale datasets to further advance indoor mapping research.

[26] M. Stoll, M. Zollhöfer, S. Izquierdo, and M. Niessner. 3D reconstruction at scale using voxel hashing. *ACM Transactions on Graphics* (TOG), 32(6):180:1–180:17.

[27] A. Pogani and D. Stricker. Structure from motion using full spherical panoramic cameras. In *Computer Vision Workshops (ICCV Workshops) 2013 IEEE International Conference on*, pages 319–323. IEEE, 2013.

[28] G. Shabotov, R. Viola, and S. Savery. Structure from small baseline motion with-convex cameras. In *Computer Vision and Pattern Recognition Workshop, 2003. CVPRW 2003. Conference on*, volume 7, pages 35–42. IEEE, 2003.

[29] P. Thaler, J. Wagner, and K. Schindler. Fast registration of laser scans with a prior constraint onto what works and what does not. *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, 3(3):149, 2014.

[30] P. Van L. Nam and B. Wandt. Block assembly for global registration of building scans. *ACM Transactions on Graphics* (TOG), 33(6):237, 2014.

[31] Q.-N. Zhou, S. Miller, and V. Kalof. Elastic fragments for dense scan reconstruction. In *2010 IEEE International Conference on Computer Vision*, pages 475–480. IEEE, 2010.