

# Texture mapping 3D planar models of indoor environments with noisy camera poses

Peter Cheng, Michael Anderson, Stewart He, Avideh Zakhor

University of California, Berkeley

## ABSTRACT

Automated 3D modeling of building interiors is used in applications such as virtual reality and environment mapping. Texturing these models allows for photo-realistic visualizations of the data collected by such modeling systems. We perform these tasks using a backpack-mounted data acquisition system. Camera poses obtained by this system and used for texturing often suffer from inaccuracies however, resulting in visible discontinuities when successive images are projected onto a surface for texturing. Existing methods to stitch images together are often computationally expensive and work independently of pose estimates and geometry information. By intelligently selecting images for texturing, we develop an efficient method to adjust camera poses in 2D, followed by two different methods to composite images together, based on the geometry of surfaces being textured. The effectiveness of our method is demonstrated on a number of different indoor environments.

**Keywords:** Texture Mapping, Reconstruction, Image Stitching, Mosaicing

## 1. INTRODUCTION

Three-dimensional modeling of indoor environments has a variety of applications in fields such as architecture, conservation, and simulation. Applying accurate textures to these models is important to add context to digital visualizations, and allows for image-based analysis of environments. In this paper, we perform data acquisition with a backpack-mounted system carried by an ambulatory human.<sup>1</sup> This backpack system contains a series of cameras, laser scanners, and inertial measurement units, used to solve the simultaneous localization and mapping (SLAM) problem, which produces backpack (and therefore camera) poses over time as well as a point cloud representing the surrounding environment.<sup>1-3</sup> For environment geometry, we work with models obtained by fitting low-resolution planar surfaces to these point clouds.<sup>4</sup>

While a human-carried system provides advantages over more common wheeled systems in terms of agility and portability, it often leads to relatively high localization error even after applying sophisticated localization techniques.<sup>1</sup> Furthermore, our model geometry usually contains some error as well, and is often tuned to ignore low-resolution features, such as small walls, furniture, and other minor non-planar features commonly found in indoor environments. As a result, methods for texture mapping that rely on accurate camera poses and projection surfaces generally produce poor results, while more sophisticated methods to accurately align images in 3D lead to high runtime on our typically very large datasets. In contrast, we present a method, which given a high number of images, searches for a subset of them that can be successfully aligned with 2D adjustments alone.

The remainder of the paper is organized as follows. Section 2 explains how images are projected onto our geometry and explains a simple approach towards selecting images for texturing. Section 3 covers existing approaches to image stitching, and their performance on our datasets. Section 4 contains our approach towards efficient 2D image alignment, followed by Section 5, which describes our two methods of compositing images. Section 6 contains results and conclusions.

## 2. SIMPLE TEXTURE MAPPING

The geometry of the texture mapping process for a planar surface is shown in Figure 1. We are provided with a set of  $M$  images to texture a target plane, where each image has a camera matrix  $P_i$  for  $i = 1..M$ , which translates a 3D point in the world coordinate system to a 2D point or pixel in image  $i$ 's coordinates. A camera matrix  $P_i$  is composed of the camera's intrinsic parameters, containing focal length and image center, as well as extrinsic parameters which specify the rotation and translation of the camera's position in 3D world coordinates

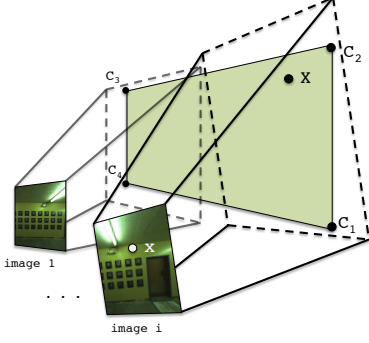


Figure 1: Surfaces to be textured are specified in 3D space by corners  $C_i$ . Images are related to each surface through the camera matrices  $P_{1..m}$ .

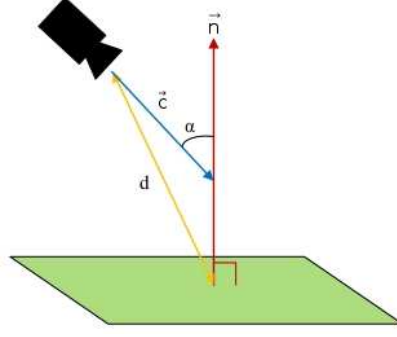


Figure 2: We minimize camera angle  $\alpha$  and distance  $d$  by maximizing the scoring function  $\frac{1}{d}(-1 \cdot \vec{c}) \cdot \vec{n}$

at the time that image  $i$  is taken. These extrinsic parameters are determined by the backpack hardware and localization algorithms<sup>1-3</sup> and are substantially noisy.

Because our backpack system takes photos at a rate of 5 Hz, hundreds of images are available for texturing each surface in our model. Our goal in designing a texture mapping process is to decide which of these images should be used, and where their contents should map onto the final texture, in order to minimize any visual discontinuities or seams that would suggest that the plane’s texture is not composed of a single continuous image.

## 2.1 Tile-based Texture Mapping

Ignoring the fact that the camera matrices  $P_{1..M}$  are inaccurate, a simple texturing approach can be performed by discretizing the target plane into small square tiles, and choosing an image to texture each tile directly.

We choose to work with rectangular units to ensure that borders between any two distinct images in the final texture are either horizontal or vertical. Since most strong environmental features inside buildings are horizontal or vertical, any visible seams in our texture will intersect them minimally and be less noticeable.

In order to select an image for texturing a tile  $t$ , we first gather a list of candidate images that contain all four of its corners, which we can quickly check by projecting  $t$  into each image using the  $P_i$  camera matrices. Furthermore, each candidate image must have been taken at a time when its camera had a clear line-of-sight to  $t$ , which can be calculated using standard ray-polygon intersection tests between the camera location,  $t$ , and every other plane.<sup>5</sup>

Once we have a list of candidate images for  $t$ , we define a scoring function in order to objectively select the best image for texturing  $t$ . Since resolution decreases and camera pose errors become more extreme over distance, we wish to minimize the distance between cameras and the surfaces they texture. Additionally, we desire images that are projected perpendicularly onto the plane, maximizing the resolution and amount of useful texture available in their projections, as well as minimizing any parallax effects due to real-world geometry not accurately represented by our digital model. In other words, we wish to minimize the angle between the tile’s normal vector and the camera axis for images selected for texturing that tile. These two criteria can be met by maximizing the function  $\frac{1}{d}(-1 \cdot \vec{c}) \cdot \vec{n}$  as shown in Figure 2. Specifically,  $d$  is the distance between the centers of a camera and a tile, and  $\vec{n}$  and  $\vec{c}$  are the directions of the plane’s normal and the camera axis respectively.

As Figure 3(a) demonstrates, there are many image boundaries with abrupt discontinuities between tiles, due to significant misalignment between images, though most of it appears to be reconcilable using 2D transforms. Given our high number of input images, a better image selection procedure will further improve results, as described in Section 5, but a more significant and reliable improvement can be found through first improving the alignment of our images.



(a)

(b)



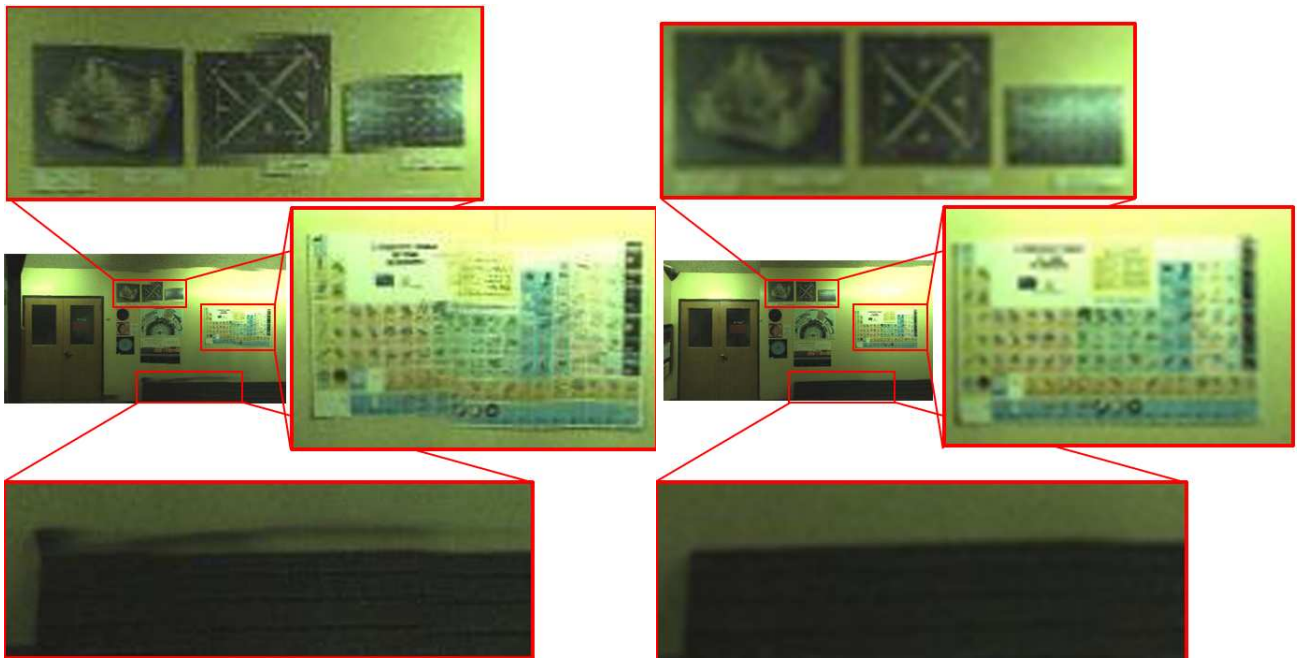
(c)

(d)



(e)

(f)



(g)

(h)

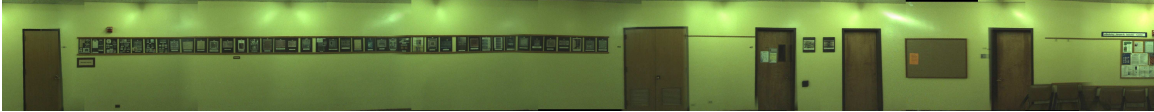
Figure 3: (a): Tile-based texturing. (b): Tile-based texturing after image alignment. (c): Tile-based texturing after image alignment with caching. (d): Shortest path texturing after image alignment). (e,f): Blending applied to (c) and (d). (g,h): Zoomed in views of discontinuities in (e) vs. in (f).



(a)



(b)



(c)

Figure 4: Texture alignment via (a) the graph-based localization refinement algorithm, (b) the AutoStitch software package, and (c) our method.

### 3. EXISTING APPROACHES TO IMAGE ALIGNMENT

Stitching together multiple images to produce a larger, seamless image is a commonly performed task, with many successful approaches over the past years. Generally, parts of images are matched to each other, usually through direct comparisons or feature detection and matching. Images are then transformed to maximize matches, often by calculating homographies between pairs of images, or by iteratively adjusting camera poses in 1 to 6 degrees of freedom.

Feature matching has a number of advantages over direct matching that make it more suitable for our often non-planar data, and rotational differences.<sup>6</sup> Feature matching however, works best when multiple unique visual references exist in the environment that can be detected in multiple images. In contrast, our indoor environments have a high prevalence of bare surfaces, as well as repeating textures, such as similar windows, doors, and wall-mounted decorations, that cause difficulty in disambiguating features. This lack of reliable reference points often results in errors when matching images together.

Additionally, our datasets often contain long chains of images, which leads to error accumulation when image correspondences are not exact and are compounded together. For example, when matching a long chain of images through homographies, a pixel in the  $n$ th image is translated into the first image’s coordinates by multiplying by the  $3 \times 3$  matrices  $H_1 H_2 H_3 \dots H_n$ . Errors in any of these homography matrices are propagated to all further images, resulting in drift.

Past work with our data acquisition backpack integrated image stitching with an iterative localization algorithm, and performed homography-based camera pose refinement.<sup>1</sup> When run on long chains of images, especially where features are sparse, it produces distorted textures, as seen in Figure 4(a), due to the problems described above. This approach is also not closed-form, and its iterative camera adjustment process over our large datasets leads to prohibitively long computation time.

The AutoStitch software package performs homography-based alignment as well, with additional provisions that attempt to reduce drift.<sup>7,8</sup> Though AutoStitch also performs well in areas with dense features, it can not handle areas without features, and has trouble aligning wall sections with even short segments of bare texture. The example in Figure 4(b) was generated after manual tuning, and areas with fewer visual features or repeating texture patterns simply failed outright with AutoStitch.

### 4. 2D IMAGE ALIGNMENT

In this section, we describe our method of efficient and robust image alignment. Though state-of-the-art techniques towards image stitching work in 3D, we elect to perform image alignments entirely in 2D. Besides the



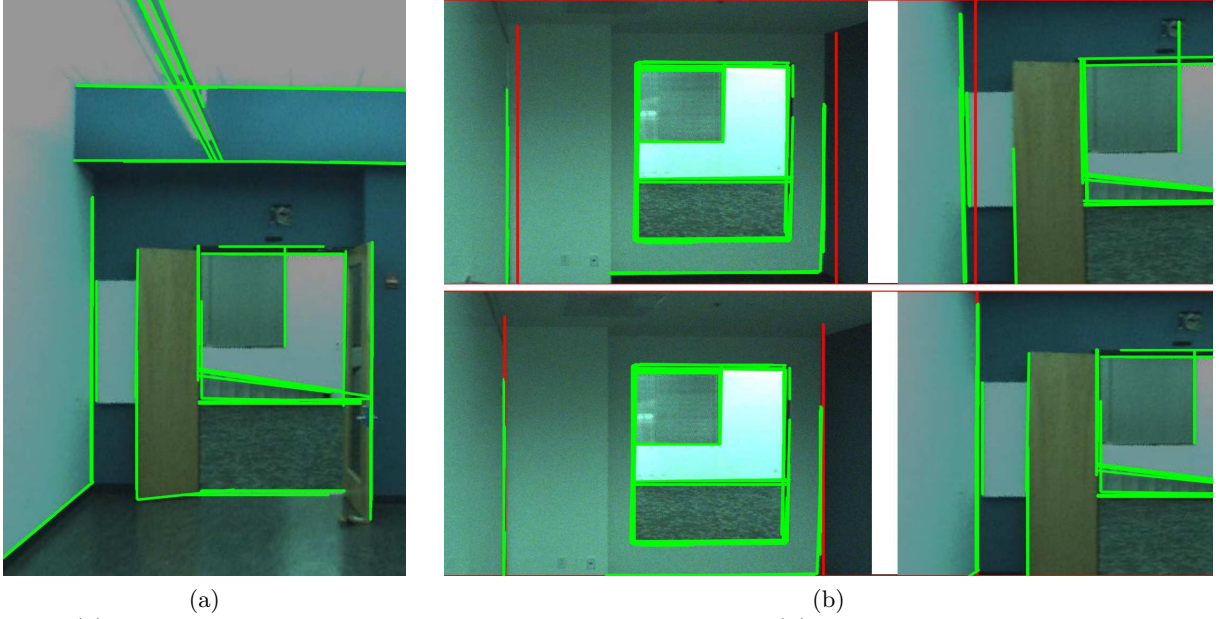


Figure 5: (a) Lines detected in image projections via Hough transform. (b) In this example, we project two images at once for the sake of demonstration, where the red lines are geometry-based lines where surfaces intersect, and the green lines are lines detected in the image. Above are the original projections using input camera poses, and below are the projections after rotation and translation for alignment.

benefits of reduced complexity, we have found 2D alignment to work well for two main reasons. First, because we have such a large number of images to choose from, the compositing approaches in Section 5 are able to select a subset of images that align well, given only 2D alignments. Second, the nature of our input data is such that localization error chiefly occurs in two dimensions, in the same plane as the surface being projected onto. Our data acquisition system contains cameras facing to the sides of the operator. The operator makes efforts to walk parallel to walls, and the localization and model-generation algorithms that provide our input conform the operator’s path to be straight and parallel to detected surfaces.<sup>3,4</sup> Since the operator also walks in a manner to keep the backpack balanced, errors in roll and yaw are minimal. Thus, our highest errors stem from uncertainty in the operator’s pitch, equivalent to rotation around the camera axes, as well as location along surfaces. These equate to 2D rotation and translation in a surface’s plane.

Our 2D alignment procedure consists of three parts. First, all images are projected onto the surface and lines within these projections are detected. These lines are then rotated and shifted to match geometry-based lines comprising the surface’s boundary and intersection with other surfaces. Second, occlusion checks are performed to remove invalid parts of each image for the target surface. These two steps are image-independent and performed in parallel. Third, we detect SIFT feature matches between pairs of images and solve a weighted linear least squares problem in 2D to maximize matches.

#### 4.1 Geometry-based Alignment

After computing each image’s projection onto the target surface, as described in Section 2, we detect line segments in the image projections using Hough transforms, as shown in Figure 5(a). Experience and intuition show that walls in indoor environments often contain linear features that are either horizontal or vertical, corresponding to doors, windows, posters, etc. Thus, when texturing walls, we rotate images in 2D such that dominant lines are made to be horizontal or vertical. This can be further generalized by instead orienting lines to be parallel with a surface’s boundaries, which makes it applicable to non-rectangular surfaces, such as ceilings, floors, or slanted walls.

Furthermore, at this point in time, image occlusions have not been accounted for. As a result, some image projections contain texture that should project to an adjacent surface, sometimes with a visible linear boundary where the two surfaces meet. If this linear boundary can be detected by Hough transform, the image is rotated and shifted in 2D such that the visual boundary between two surfaces in an image projection matches the physical boundary in our digital model. An example of such an adjustment is in Figure 5(b) and 5(c).

To perform rotation, we use a RANSAC<sup>9</sup> framework to determine an optimal rotation angle between lines in our image and lines in our geometry. RANSAC allows us to ignore erroneous lines in our image, such as the slanted lines corresponding to the ladder in Figure 5(a). Once this rotation is discovered and applied, we gather all pairs of image line segments and geometry line segments with less than  $0.1^\circ$  difference in orientation and less than 250 mm distance at their furthest points. If 2 or more pairs exist, we select the 2 that contain the longest noncollinear image-based line segments, and perform the corresponding fixed translation to match them. If we only have 1 pair, we perform the minimal shift to match the lines of that pair, and make note of their orientation, for usage in Section 4.3.

## 4.2 Image Occlusion

Now that images have been aligned to geometry where possible, we perform a simple recursive occlusion procedure on each image. This is done by performing the intersection tests in Section 2 in a regularly spaced grid. Where four corners of a rectangular region are occluded, texture is removed. Where no corners are occluded, the recursion stops. Where there is a mixture of both, the rectangular region is subdivided into four, and the same process is performed on each. By performing Section 4.1's alignment procedure before occlusion, texture belonging to other surfaces is accurately removed, which is necessary for the next section.

## 4.3 2D Feature Alignment

Our next step is to align images by searching for corresponding points between all pairs of overlapping images. We use feature alignment rather than pixel or intensity-based alignment due to the differences in lighting as well as possible occlusion among our images, both of which feature alignment is less sensitive to.<sup>6,10,11</sup> We use SiftGPU<sup>12</sup> for its high performance on both feature detection as well as pairwise matching. These matches determine  $dx$  and  $dy$  distances between each pair of features for two image projections, though these distances may not always be the same for different features. Since indoor environments often contain repetitive features such as floor tiles or doors, we need to ensure that SIFT-based distances are reliable. First, we only perform alignment on parts of images that overlap given the original noisy poses. Second, we discard feature matches that correspond to an image distance greater than 200 mm from what the noisy poses estimate. In order to utilize the remaining feature matches robustly, RANSAC<sup>9</sup> is again used to estimate the optimal  $dx_{i,j}$  and  $dy_{i,j}$  distances between two images  $i$  and  $j$ . We use a 5 mm threshold for RANSAC, so that SIFT matches are labeled as outliers if their distance is not within 5 mm of the sampled average distance.

We now use the  $dx_{i,j}$  and  $dy_{i,j}$  distances between each pair of images to refine their positions using weighted linear least squares. The variables we wish to solve for are the  $x_i$  and  $y_i$  positions of our images, while our equations are the SIFT-based distances between pairs of images, images fixed to geometry with 0 or 1 degrees of freedom, and the original camera poses. Ignoring weighting and the original camera poses for now, an example setup for solving  $\min_{\vec{\beta}} \|A\vec{\beta} - \vec{\gamma}\|_2^2$  with 3 images is shown below.

$$A = \begin{pmatrix} -1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 & 0 \\ 0 & -1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & -1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & -m & 0 & 0 & 1 & 0 \end{pmatrix} \quad \vec{\beta} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ y_1 \\ y_2 \\ y_3 \end{pmatrix} \quad \vec{\gamma} = \begin{pmatrix} dx_{1,2}, \\ dy_{1,2}, \\ dx_{2,3}, \\ dy_{2,3}, \\ t1_x, \\ t1_y, \\ -mx_2 + ty_2 \end{pmatrix}$$

In this scenario, a SIFT-based distance of  $dx_{1,2}$ ,  $dy_{1,2}$  was calculated between images 1 and 2. This corresponds to the first and second row of  $A$ , while the third and fourth row of  $A$  represent the same for images 2 and 3.

Rows 5 and 6 of  $A$  correspond to a fixed location of  $tx_1, ty_1$  for image 1, while Row 7 corresponds to a constraint to a line of slope  $m$ , with current location  $tx_2, ty_2$ , both possible results of the geometry alignment procedure in Section 4.1. If we do not have enough SIFT matches, or lack images matched to geometry, our matrix becomes rank-deficient and our problem cannot be solved. As a result, we add rows for each image corresponding to the original noisy image locations, and downweight them heavily, e.g. with a weight of 0.01 versus 1.

Because our problem is linear, it is easily and quickly solved, and after applying the resulting shifts, our images overlap and match each other with far greater accuracy. Applying the simple mapping scheme in Section 2.1 to the same wall used in that section results in Figure 3(b), which has far fewer discontinuities, though errors due to lighting differences and parallax effects are still visible.

## 5. IMAGE COMPOSITING

In this section, we revisit the tile-based texturing approach from Section 2, with an added caching mechanism to reduce image boundaries. This method works well given all manner of camera poses and surfaces, but for optimal cases where we have large sections of usable texture from images, we propose a superior method that further reduces image boundaries.

### 5.1 Tile-Mapping with Caching

From Section 2.1, we saw that discontinuities occur where adjacent tiles are textured by different images. Though Section 4’s image alignment removes many such discontinuities, there are still cases where seams are visible due to imprecise matching or other factors such as model-based errors. To reduce the cases where this happens, it makes sense to take into account image selections made by neighboring tiles while texture mapping a given tile. By using the same image across tile boundaries, we can eliminate a discontinuity altogether. If this is not possible because a tile is not visible in images used by neighboring tiles, using similar images across tile boundaries also leads to less noticeable discontinuities.

Essentially a caching mechanism, we select the best image for a tile  $t$  by searching through two subsets of images for a viable candidate, before searching through the entire set of valid images. The first subset of images is the images selected by adjacent tiles that have already been textured. We must first check which of these images can map to  $t$ , and then of those, we make a choice according to the scoring function in Figure 2. Before reusing this image, we ensure it meets the criteria  $\alpha < 45^\circ$ , in order to ensure a high resolution projection, with  $\alpha$  as the camera angle as shown in Figure 2.

If no satisfactory image is found in the first subset, we check a second subset of images, consisting of those taken near the ones in the first subset, both spatially and temporally. These images are not the same as the ones used for neighboring tiles, but are taken at a similar location and time, suggesting that their localization and projection are quite similar, and thus likely matched more cleanly. If no viable image is found according to the same criteria as before, we search the entire set of candidate images, selecting based on the same scoring function from Figure 2.

The result of this caching approach is shown in Figure 3(c), where seams are now reduced as compared to Figure 3(b). However, some discontinuities are still present, as visible in the posters on the wall with breaks in their borders.

As mentioned earlier, our data comes from a mobile backpack system. Human operators can not carry the backpack in a perfectly upright position and are bent forwards at 15 to 20 degrees with respect to the vertical direction. As a result, cameras facing sideways are head on with respect to vertical walls, while cameras oriented towards the top or bottom of the backpack are at an angle with respect to horizontal floors and ceilings. This is depicted in Figures 6(a) and 6(b). These oblique camera angles for horizontal surfaces translate into textures that span large areas once projected, as shown in Figure 6(b). Using the tile-based texture mapping criteria from Figure 2, such projections have highly varying scores depending on the location of a tile on the plane. Thus, the tiling approach in this section is a good choice for texturing floors and ceilings, as it uses the parts of image projections that maximize resolution and accuracy for their respective plane locations, e.g. areas near point A and not near point B, in Figure 6(b).

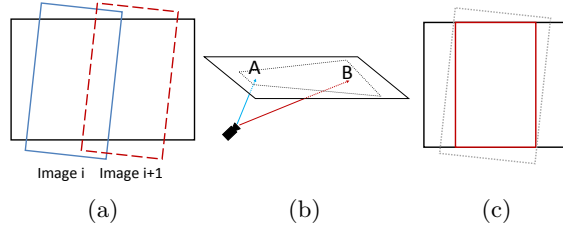


Figure 6: (a) Images for vertical planes are tilted, but their camera axes are more or less normal to their respective planes. (b) Camera axes for ceiling images are at large angles with respect to plane normals. (c) Wall images are cropped to be rectangular.

## 5.2 Shortest Path Texturing

For vertical walls, most images are taken from close distances and head-on angles, resulting in high resolution fronto-parallel projections. As a result, for each tile on a wall plane, the scoring function of Figure 2 is relatively flat with respect to candidate images, as they are all more or less head on. Thus, the scoring function is less significant for walls, and it is conceivable to use a different texturing strategy to directly minimize visible seams when texturing them. This is done by choosing the smallest possible set of images that (a) covers the entire plane and (b) minimizes the visibility of borders between them. A straightforward cost function that accomplishes the latter is the sum of squared differences (SSD) of pixels in overlapping regions between all pairs of images. Minimizing this cost function encourages image boundaries to occur either in featureless areas, such as bare walls, or in areas where images match extremely well.

The first step in this procedure is to obtain a list of images useful for texturing the wall. A simple way to do this is to use the images selected by the tiling process in Section 2.1. Such a list of images is guaranteed to cover the entire wall, and consists of desired camera poses overall.

To cover the entirety of a plane, our problem can be defined as minimally covering a polygon i.e. the planar surface, using other polygons of arbitrary geometry i.e. image projections, with the added constraint of minimizing the cost function between chosen images. This is a complex problem, though we can take a number of steps to simplify it.

Given that wall-texture candidate images are taken from more or less head-on angles, and knowing that only minor rotations are made in Section 4, we can reason that their projections onto the plane are approximately rectangular. By cropping them all to be rectangular, as shown in Figure 6(c), our problem becomes the conceptually simpler one of filling a polygon with rectangles, such that the sum of all costs between each pair of rectangles is minimal. We thus also retain the advantages of working with rectangular units, as explained in Section 2.1.

The operator’s path, and the location and orientation of the cameras on the acquisition backpack is such that images nearly always contain the entirety of the floor to ceiling range of wall planes. Images therefore rarely need to be projected with one above another when texturing wall planes. In essence, we need only to ensure lateral coverage of wall planes, e.g. from left to right, as our images provide full vertical coverage themselves. We can thus construct a Directed Acyclic Graph (DAG) from the images, with edge costs defined by the SSD cost function, and solve a simple shortest path problem to find an optimal subset of images with regard to the SSD cost function.<sup>13</sup>

Figure 7 demonstrates the construction of a DAG from overlapping images of a hallway wall. Images are sorted by horizontal location left to right, and become nodes in a graph. Directed edges are placed in the graph from left to right between overlapping images. The weights of these edges are determined by the SSD cost function. Next, we add two artificial nodes, one start node representing the left border of the plane, and one end node representing the right border of the plane. The left(right) artificial node has directed edges with equal and arbitrary cost  $C_0$  to(from) all images that meet the left(right) border of the plane. We now solve the shortest path problem from the start node to the end node. This results in a set of images completely covering the plane horizontally, while minimizing the cost of seams between images.



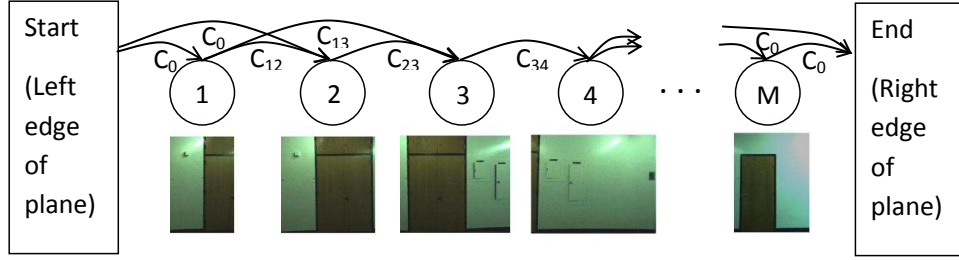


Figure 7: DAG construction for the image selection process.

In rare cases where the vertical dimension of the plane is not entirely covered by one or more chosen images, we are left with holes where no images are selected to texture. Since these holes are rare, and generally fairly small, we use a greedy approach, repeatedly filling the hole with images that result in the lowest SSD costs in a blending region around the hole, as discussed in Section 5.3. This method is not as optimal as a true 2D-coverage solution would be, but it is a fast approximation, and adequately handles the few holes we encounter.

With this completed, we have now mapped every location on the plane to at least one image, and have minimized the number of images, as well as the discontinuities at their borders. As seen in Figure 3(d), this shortest path method has fewer visible discontinuities than Figure 3(c) corresponding to the tile caching approach\*. This is especially evident when comparing the posters in the images. This shortest path approach directly reduces the cost of each image boundary, while the tile caching method uses a scoring function that only approximates this effect. Furthermore, this approach guarantees the best selection of images to minimize seams, while the sequential tile caching method may select images early on that turn out to be poor choices once subsequent tiles have been processed. This shortest path approach is also far less intensive in terms of memory usage and runtime, both during texture generation and rendering, as it does not require discretizing planes or images.

When texturing an entire 3D planar model, we apply the shortest path method on walls, due to its superior output when provided with head-on images. Floors and ceilings however, given their many images taken at oblique angles, are textured using the tile caching method.

### 5.3 Blending

We now apply a blending procedure to both texturing methods. Although the image alignment steps and image selection in both methods attempt to minimize all mismatches between images, there are occasional unavoidable discontinuities in the final texture due to different lighting conditions or inaccuracies in model geometry. These can however be treated and smoothed over by applying alpha blending over image seams. Whether the units we are blending are rectangularly-cropped images or rectangular tiles, we can apply the same blending procedure, as long as we have a guaranteed overlap between units to blend over.

For the tile caching method, we can ensure overlap by texturing a larger tile than needed for display. For example, for a rendered tile  $l_1 \times l_1$ , we can associate it with a texture  $(l_1 + l_2) \times (l_1 + l_2)$  in size. We have found  $l_2 = \frac{l_1}{2}$  to provide a balance between blending and keeping features unblurred. For the shortest path method, we have already ensured overlap between images. To enforce consistent blending however, we add a minimum required overlap distance of 200 mm while solving the shortest path problem in Section 5.2. Additionally, if images overlap in a region greater than the overlap distance, we only apply blending over an area equal to the overlap distance.

After performing linear alpha blending across overlapping regions, our texture mapping process is complete. Figures 3(e) and 3(f) show the blended versions of Figures 3(c) and 3(d) respectively. The remaining images in Figure 3 highlight differences between the two methods, showing that Figure 3(f) clearly has the best visual quality and the best texturing approach among the textures in Figure 3.

\*In Figure 3(d), we arbitrarily chose one image for texturing where images overlap, as blending will be discussed in section 5.3.

## 6. RESULTS AND CONCLUSIONS

Examples of ceilings and floors textured with the tile caching approach, and walls textured with the shortest path approach, are displayed in Figure 8. High resolution colored texture comparisons, as well as further examples and interactive walkthroughs are available at <sup>†</sup>.

As mentioned earlier, our approach is quite efficient. The top wall in Figure 8(a) was generated with  $7543 \times 776$  pixels, and spans a 40-meter long wall. Given 41000 input images, a 2.8GHz dual-core consumer-grade laptop takes approximately a minute to pick 36 candidate images, followed by another minute to perform both image alignment and the shortest path texturing method, though over 75% of that time is spent calculating SIFT matches within the SiftGPU framework. While not real-time, the process is capable of generating quick updates after changes in various parameters or modifications to input data, and if integrated directly into a modeling system, could provide live visualization as data is collected. Our full models consist of our input model file, textures that we generate, and a mapping of image points to model vertices. The models shown in Figure 8 are roughly 20 MB in size, and are visualized using the OpenSceneGraph toolkit,<sup>14</sup> which allows for export to many common model formats, as well as efficient visualization, even in web browsers or mobile devices.

In this paper, we have developed an approach to texture mapping models with noisy camera localization data. We are able to refine image locations based on geometry references and feature matching, and robustly handle outliers. Using the tile-based mapping approach, we can texture both simple rectangular walls as well as complex floor and ceiling geometry. We also implemented a shortest path texturing method that produces seamless textures on planes where multiple head-on images are available. Both of these approaches are highly modular, and easily tunable for similar systems across multiple environments.

### 6.1 Future Work

While our method works well on our datasets, it is heavily reliant on having a large number of input images, and minimal non-2D error. A technique for resolving camera error in 3D that would mesh well with our approach is to perform matching between image lines and geometry in 3D, which can be done reasonably efficiently.<sup>15,16</sup> Using linear features on top of our SIFT features is also likely to result in improved matches, as indoor scenes often have long, unbroken lines spanning multiple images.<sup>17</sup> Finally, our blending approach is quite basic, and applying more sophisticated methods of blending as well as normalization would benefit our final visual quality, and more robustly handle motion-based or parallax errors.

## REFERENCES

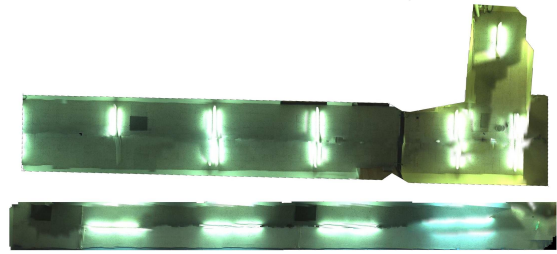
- [1] Liu, T., Carlberg, M., Chen, G., Chen, J., Kua, J., and Zakhor, A., “Indoor localization and visualization using a human-operated backpack system,” in [*Indoor Positioning and Indoor Navigation (IPIN), 2010 International Conference on*], 1–10, IEEE (2010).
- [2] Chen, G., Kua, J., Shum, S., Naikal, N., Carlberg, M., and Zakhor, A., “Indoor localization algorithms for a human-operated backpack system,” in [*Int. Symp. on 3D Data, Processing, Visualization and Transmission (3DPVT)*], (2010).
- [3] Kua, J., Corso, N., and Zakhor, A., “Automatic loop closure detection using multiple cameras for 3d indoor localization,” in [*IS&T/SPIE Electronic Imaging*], (2012).
- [4] Sanchez, V. and Zakhor, A., “Planar 3d modeling of building interiors from point cloud data,” in [*International Conference on Image Processing*], (2012).
- [5] Glassner, A. et al., [*An introduction to ray tracing*], Academic Press (1989).
- [6] Szeliski, R., “Image alignment and stitching: A tutorial,” *Foundations and Trends® in Computer Graphics and Vision* **2**(1), 1–104 (2006).
- [7] Brown, M. and Lowe, D. G., “Automatic panoramic image stitching using invariant features,” *Int. J. Comput. Vision* **74**, 59–73 (Aug. 2007).
- [8] Brown, M. and Lowe, D., “Autostitch.” <http://www.cs.bath.ac.uk/brown/autostitch/autostitch.html>.

---

<sup>†</sup><http://www.eecs.berkeley.edu/~pcheng/indoormapping>



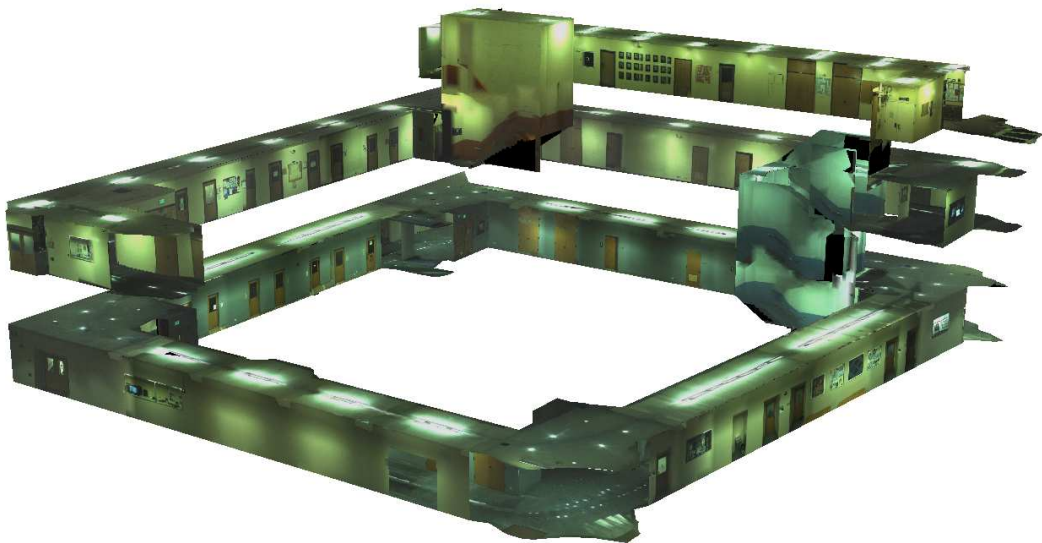
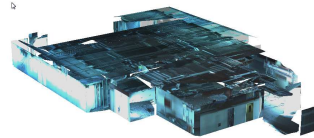
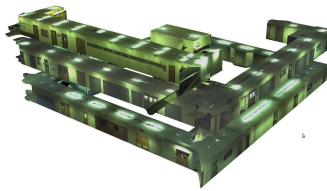
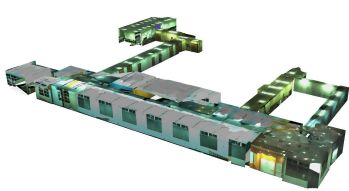
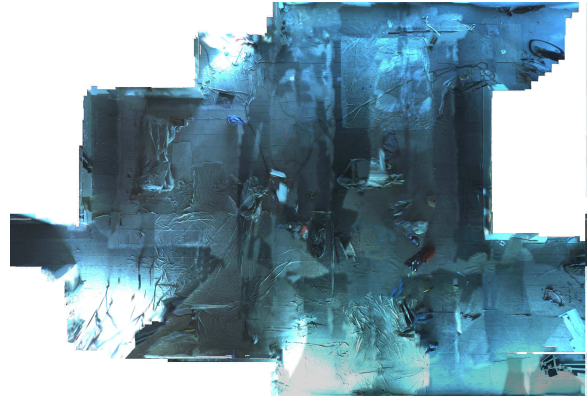
(a)



(b)



(c)



(d)

Figure 8: Examples of our final texture mapping output for (a) walls, (b) ceilings, (c) floors, (d) full models.

- [9] Fischler, M. and Bolles, R., “Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography,” *Communications of the ACM* **24**(6), 381–395 (1981).
- [10] Lowe, D., “Object recognition from local scale-invariant features,” in [*Computer Vision, 1999. The Proceedings of the Seventh IEEE International Conference on*], **2**, 1150–1157, Ieee (1999).
- [11] Mikolajczyk, K. and Schmid, C., “A performance evaluation of local descriptors,” *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **27**(10), 1615–1630 (2005).
- [12] Wu, C., “Siftgpu.” <http://www.cs.unc.edu/~ccwu/siftgpu/>.
- [13] Dijkstra, E., “A note on two problems in connexion graphs,” *Numerische Mathematik* **1**, 269–271 (1959).
- [14] “Openscenegraph.” <http://www.openscenegraph.org/projects/osg>.
- [15] Elqursh, A. and Elgammal, A., “Line-based relative pose estimation,” in [*Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*], 3049–3056 (june 2011).
- [16] Koeck, J. and Zhang, W., “Extraction, matching and pose recovery based on dominant rectangular structures,” (2005).
- [17] Ansar, A. and Daniilidis, K., “Linear pose estimation from points or lines,” *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **25**, 578 – 589 (may 2003).