

Multi-View Stereo Revisited

Michael Goesele Brian Curless Steven M. Seitz
University of Washington

Abstract

We present an extremely simple yet robust multi-view stereo algorithm and analyze its properties. The algorithm first computes individual depth maps using a window-based voting approach that returns only good matches. The depth maps are then merged into a single mesh using a straightforward volumetric approach. We show results for several datasets, showing accuracy comparable to the best of the current state of the art techniques and rivaling more complex algorithms.

1. Introduction

In the last decade, multi-view stereo algorithms have evolved significantly in both their sophistication and quality of results. While the early algorithms were simple extensions of window-based binocular stereo matchers, the best currently available methods employ powerful nonlinear energy minimization methods (e.g., graph cuts [8, 16, 18, 15], level set evolution [3, 13], mesh evolution [7]) often incorporating careful treatment of visibility conditions and silhouette constraints.

In this paper, we take a big step *backward* and argue that some simple modifications of the original window-based algorithms can produce results with accuracy on par with the best current methods. Our algorithm consists of two steps: In the first step we reconstruct a depth map for each input view, using a robust version of window-matching with a small number of neighboring views. The second step is to merge the resulting depth maps using a standard volume merging method [2].

The key new idea underlying this work is to attempt to reconstruct only the portion of the scene that can be matched with high confidence in each input view. Consequently, each individual depth map may contain numerous holes, e.g., near silhouettes, oblique surfaces, occlusions, highlights, low-textured regions and so forth. Because most of these effects (with the exception of low-textured regions) occur in different image regions in different views, the merging step fills in most of these missing areas, and improves accuracy in regions that are reconstructed multiple



Figure 1. Reconstruction (*right*) from the *templeFull* dataset and an input image (*left*) for comparison. Despite the simplicity of the proposed algorithm, it is able to estimate the object's shape to sub-millimeter accuracy.

times. While this is a simple idea, it is a departure from most modern binocular and multi-view stereo algorithms, which seek models that are as complete as possible, often using regularization to help fill in uncertain regions. While we do not optimize or enforce completeness, our algorithm nonetheless reconstructs dense shape models with few holes and gaps, given a sufficient number of input views.

The benefits of our algorithm are as follows:

- The algorithm outputs dense and very accurate shape estimates, on par with the current best performing methods (as reported in a new evaluation of multi-view stereo methods [14]).
- Confidence is provided for each reconstructed point.
- The performance is easy to understand, analyze, and predict, due to the algorithm's simplicity.
- It is unusually easy to implement.

Some disadvantages of our approach are

- Holes can occur in areas with insufficient texture.
- Each surface point must be seen in at least three views to be reconstructed.

The technical tools used in this algorithm are not particularly new — indeed, they are directly inspired by prior work, in particular [9, 7]. However, they have not been used previously in this combination, and we argue that it is this particular synthesis of existing ideas which is the key to its success. A second contribution of this paper is an analysis of *why* the algorithm performs as well as it does, and under which conditions it fails.

The remainder of this paper is structured as follows: We first give an overview of related work before we describe our algorithm in detail (Section 2). We then describe the datasets used for our reconstructions (Section 3) and show results for a large number of parameter settings (Section 4) before we conclude.

1.1. Related Work

While there is a large body of prior work on multi-view stereo algorithms, the three papers that are most closely related are Narayanan et al.’s *Virtualized Reality* technique [9], Pollefeys et al.’s *visual modeling* system [11, 12], and the *multi-stereo* approach of Hernández and Schmitt [6, 7].

Narayanan et al. [9] proposed the idea of creating dense shape models by volumetric merging of depth maps. The key difference between their work and ours is the method used to reconstruct depth maps. Narayanan et al. used a traditional multi-baseline stereo matcher [10] that seeks to estimate a *complete* depth map. As the authors point out, this method produces noisy results with many outliers that lead to problems with merging and errors in the resulting shapes. In contrast, we devise a specialized matcher that computes depth only at high confidence points, simplifying the merging step and leading to much higher quality reconstructions.

Pollefeys et al. [11, 12] use a three step technique. They first perform a pair-wise disparity estimation for directly adjacent views using a modification of Cox et al.’s dynamic programming scheme [1] which yields dense but incomplete depth maps by enforcing various constraints on the solution. An optimal joint estimate for each view is then computed by adding corresponding disparity estimates from gradually farther away views on a per-pixel basis as long as they are not classified as outliers relative to the current depth estimate for the pixel under consideration. The fused depth maps are then combined using a volumetric merging approach. Compared to Pollefeys et al., our system reconstructs depth maps in a single pass with a much simpler approach yielding potentially less complete depth maps. In addition, we only use the quality of a match between a fixed number of neighboring views as the acceptance criterion instead of performing an outlier classification based on reconstructed depth values.

The first step of our approach (estimating depth maps) is inspired by the work of Hernández and Schmitt [7], who

also use robust window-based matching to compute reliable depth estimates. While the high level ideas are similar, many of the details are quite different from what we do. First, we are using a simpler and more conservative correlation criterion. Hernández [6] computes the local maxima of the correlation curves between the reference view and the nearby images. These are used to vote for a depth range within which the global maximum is determined from all views that pass a threshold-based selection criterion. In contrast, we require that at least two views pass a threshold-based selection criterion at each candidate depth value. The other important way in which our method differs from [7] is the method of combining depth maps into a full 3D model. They use a combination of volume filtering, mesh evolution based on a snakes formulation, and additional silhouette terms to recover a complete model. The resulting approach, while it generates beautiful results, has very different properties and assumptions than our approach. Since it is based on local refinement via snakes, [7] requires a close initial estimate of the shape being estimated, and the topology of the object must be the same as that of its visual hull. They also require that silhouettes are extracted. In contrast our approach does not require an initial surface estimate, and does not place any restriction on the topology of the object. While we do not require silhouettes, our algorithm can take advantage of them, when available. An advantage of [7] is that it produces complete surface models and can fill in holes using regularization and silhouette terms. While our approach can leave holes in low-contrast regions, the lack of a smoothness term has the advantage of avoiding smoothing over sharp features. A final difference is that our approach is very simple to implement and reproduce, in comparison to [7].

2. Algorithm Description

Our algorithm consists of two steps: 1) reconstructing a depth map for each input view, and 2) merging the results into a mesh model. In the first step, depth maps are computed using a very simple but robust version of window-matching with a small number of neighboring views. We also estimate a confidence value for each pixel and only high confidence points are included in the merging step.

In the second step, we merge the resulting set of confidence-weighted depth maps using the volumetric method by Curless and Levoy [2]. The result of the second step is a triangular mesh with per-vertex confidence values.

The following sections describe both steps of the algorithm in more detail.

2.1. Step 1: Depth Map Generation

We assume as input a set of views $\mathbf{V} = \{V_0, \dots, V_{n-1}\}$ of an object along with camera parameters and an ap-

proximate bounding box or volume containing the object. For each view $R \in \mathbf{V}$ (hereforth called a *reference view*) we first select a set of k neighboring views $\mathbf{C} = \{C_0, \dots, C_{k-1}\} \subset \mathbf{V} - R$ against which we correlate R using robust window-matching.

For each pixel p in R , we march along its backprojected ray inside the bounding volume of the object. For each depth value d we reproject the resulting 3D location into all views in \mathbf{C} . We compute the normalized cross-correlation $NCC(R, C_j, d)$ between an $m \times m$ window centered on p and the corresponding windows centered on the projections in each of the views C_j with subpixel accuracy. (Appendix A defines the normalized cross-correlation NCC formally.)

If two views show the same surface area of a textured Lambertian object, we expect to see a high NCC score for some value of d . If, in contrast, there is an occlusion, specular highlight, or other compounding factor, the NCC value will typically be low for all depths. We wish to rely on a depth value only if the window in the reference view correlates well with the corresponding window in multiple views. We therefore define that a depth value d is *valid* if $NCC(R, C_j, d)$ is larger than some threshold *thresh* for at least two views in \mathbf{C} . The set of all views with NCC larger than *thresh* for a given depth d is denoted as $\mathbf{C}_v(d)$.

For a valid depth d we compute a correlation value $corr(d)$ as the mean of the NCC values of all views in $\mathbf{C}_v(d)$:

$$corr(d) = \frac{\sum_{C_j \in \mathbf{C}_v(d)} NCC(R, C_j, d)}{\|\mathbf{C}_v(d)\|}.$$

$\|\mathbf{C}_v(d)\|$ evaluates to the number of elements in $\mathbf{C}_v(d)$. For each pixel p in R , the depth is chosen to be the value of d that maximizes $corr(d)$, or none if no valid d is found.

Note that this approach is extremely simple, and very similar to standard SSSD-style multi-baseline window matching methods [10], with the following modifications: 1) the robust version of NCC effectively minimizes the impact of occlusions and specularities, and 2) we compute depth only at high confidence points in the image.

We also compute a confidence value $conf(d)$ for each recovered depth value as follows:

$$conf(d) = \frac{\sum_{C_j \in \mathbf{C}_v(d)} (NCC(R, C_j, d) - thresh)}{\|\mathbf{C}\|(1 - thresh)}.$$

This confidence function increases with the number of valid views and is used to inform the merging step, described in the next subsection.

There are a number of free parameters in the above description, i.e., the number k and selection of neighboring views, the sampling rate in depth, the window size m , and the threshold *thresh*. We discuss our choice of these parameters in Section 4.

dataset	# views	geometry
templeFull	317	hemisphere
templeRing	47	ring
templeSparseRing	16	ring
dinoFull	363	hemisphere
dinoRing	48	ring
dinoSparseRing	16	ring
nskulla-half	24	8-ring + 16-ring
nskulla-small	24	8-ring + 16-ring

Table 1. Specifications of the datasets. All temple and dino datasets have a resolution of 640×480 pixels. The images of the original nskulla datasets [4] are cropped to different sizes within the dataset. We scaled them down to a resolution of approximately 1000×900 pixels (*nskulla-half*) and 400×360 pixels (*nskulla-small*).

2.2. Step 2: Merging Depth Maps

Step 1 produces a set of incomplete depth maps with confidence values. In Step 2, we merge them into a single surface mesh representation using the freely available implementation of the volumetric method of Curless and Levoy [2, 17]. This approach was originally developed for merging laser range scans. In a nutshell, it converts each depth map into a weighted signed distance volume, takes a sum of these volumes, and extracts a surface at the zero level set. More details can be found in [2].

This merging approach has a number of nice properties that make it particularly appropriate for our algorithm, in particular robustness in the presence of outliers and representation of directional uncertainty. The merging algorithm starts by reconstructing a triangle mesh for each view and downweighting points near depth discontinuities and points seen at grazing angles. These meshes are then scan-converted with per-vertex weights into a volume for merging. Outliers consisting of one or two samples are filtered out automatically, because they cannot form triangles in the first phase of the algorithm. Larger handfulls of outliers will be reconstructed as small disconnected surfaces; these surfaces will have low weight, since all the points are near depth discontinuities and are probably not substantiated by other views. They can be eliminated in a post-processing step by removing low confidence geometry or by extracting the largest connected component. In addition, the approach has been shown to be least squares optimal under certain conditions, particularly assuming uncertainty distributed along sensor lines of sight [2] which by construction applies to the depth maps from Step 1.

3. Description of Datasets

We now describe the datasets used in this paper – the temple, dino, and nskulla datasets (see Table 1 for their specifications). The temple object is a 159.6 mm tall plaster

reproduction of an ancient temple. It is quite diffuse and contains lots of geometric structure and texture. The temple was illuminated by multiple light sources and captured with a camera mounted on a calibrated spherical gantry. Images with cast shadows where the camera or the gantry were in front of a light source were removed from the dataset. *templeFull* is the full dataset with 317 images. *templeRing* contains only 47 views on a ring around the object, *templeSparseRing* is a more sparse version of the *templeRing* dataset with 16 views on a ring around the object. All images have a resolution of 640×480 pixels.

The dino object is a 87.1 mm tall, white, strongly diffuse plaster dinosaur model. It was captured in the same way as the temple and consists of three sets of images: *dinoFull* (sampled along the hemisphere), *dinoRing* (sampled in a ring around the object), and *dinoSparseRing* (sparsely sampled in a ring around the object). A more detailed description of the dino and temple datasets can be found in [14].

The nskulla object is a plaster cast of a human skull. It was lit by several light sources with diffusers to approximate diffuse lighting. The skull was rotated on a turntable while cameras and lights remained fixed so that the lighting conditions are different for each image. Moving highlights are clearly visible on the object's moderately specular surface. The nskulla dataset contains 16 images captured on a ring around the object plus an additional 8 images captured on a sparser ring at higher elevation angles. The datasets differ only in resolution and were downsampled from their original resolution to approximately 1000×900 pixels (*nskulla-half*) and 400×360 pixels (*nskulla-small*). A more detailed description of the nskulla dataset can be found in [4].

4. Results

The description of the algorithm in Section 2 contains several parameters. In this section, we briefly describe how each parameter was chosen in our reconstructions and show results for other parameter choices. We also discuss the influence of other factors such as the reflectance properties of the object. Finally, we report some results of an evaluation of the reconstructed models against ground truth.

4.1. Implementation Notes

We generally set the number of neighboring views $k = 4$. A larger k reduces occlusions but does not significantly improve the results. Due to the arrangement of the camera positions around a ring or distributed on some portion of a hemisphere, we selected neighboring views based on angular distance between the optical axes. For a given reference view, the k closest views were chosen as neighboring views unless the angular distance between a view and the reference view or any other neighboring view was less than 4

degrees. We used a fixed sampling rate Δd in depth to find an initial depth estimate d' . The final depth d was computed by re-running the algorithm with step size $\frac{\Delta d}{10}$ in the interval $(d' - \Delta d, d' + \Delta d)$. We selected $\Delta d = 2.5$ mm for the temple and dino dataset and $\Delta d = 0.2$ for the nskulla datasets. The default value for the NCC threshold is $thresh = 0.6$ except for the nskulla dataset, as discussed in Section 4.5.

4.2. Window Size

Figure 2 shows a comparison of rendered depth maps reconstructed for the same reference view R with two different window sizes. The center row shows the correlation value $corr(d)$ and the bottom row displays the confidence value $conf(d)$.

Overall, the behavior is as normally expected: A larger window size leads to smoother depth maps and the removal of lower confidence values from the reconstruction (e.g., noise in the background, fine structures in the columns). Note that most of the background noise in the 5×5 dataset is not contained in the depth map generated by [17] due to the inherent outlier filtering (see Section 2.2). It is also assigned a much lower confidence value and can therefore be easily removed. The examples show also that the algorithm detects occlusions reliably for all window sizes without compromising the correlation values — only the confidence value is scaled according to the number of valid views. This is visible in the confidence image as a dark vertical stripe on the left side of the temple where some of the input views are occluded. In our experiments, we found a 5×5 window gave good results, and we used this size for all of our reconstructions.

4.3. Density of Views

Our algorithm requires that each surface point (and the window surrounding it) is seen in at least three views (a reference view and at least two neighboring views) and furthermore yields a high correlation value. The reconstructions of the *templeSparseRing* dataset are therefore incomplete even in low-occlusion areas (see Figures 3 and 7). Some high-occlusion regions such as the temple roof are missing almost completely. The results for the *templeRing* and the *templeFull* dataset show however that adding more views drastically improves the results yielding almost complete surface coverage. The *templeFull* reconstruction contains holes almost exclusively in areas that need to be observed from below. Such views are not included in the dataset.

4.4. Surfaces without Texture

The dino plaster cast has a white, Lambertian surface without obvious texture. Due to the lack of structure, stereo reconstruction using window-matching is extremely difficult. Nevertheless, the algorithm reconstructs geometry for

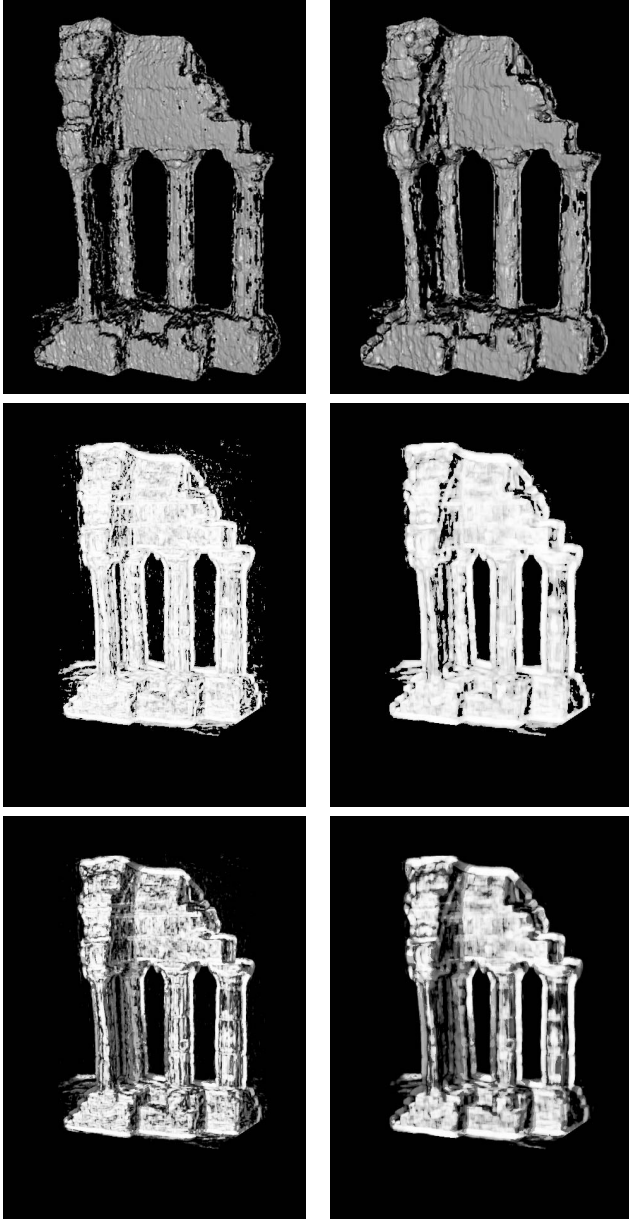


Figure 2. *Top to bottom*: Reconstructed depth map (rendered as triangle mesh generated by [17]), correlation values, and confidence values for a view from the *templeRing* dataset. *Left*: Reconstruction with window size 5×5 . *Right*: Reconstruction with window size 9×9 . The “glow” around the silhouettes of the temple is discussed in Section 4.6.

a large portion of the surface. The input images were captured under fixed but not completely diffuse illumination so that the surface shading is stationary. Regions in the vicinity of shadow boundaries and geometric features are therefore reconstructed. In addition, the algorithm reconstructs geometry in the neighborhood of dust specks on the plaster surface.

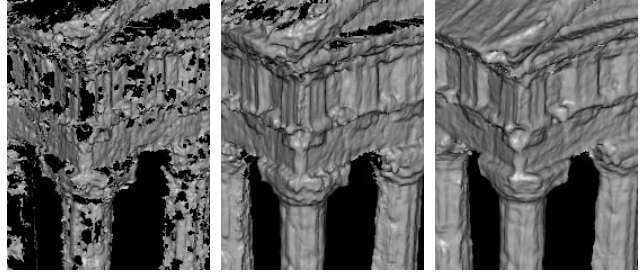


Figure 3. Detail of the temple reconstructed from increasing numbers of input views. *Left to right*: *templeSparseRing* (16 views), *templeRing* (47 views), *templeFull* (317 views). The full version of the datasets is shown in Figure 7.

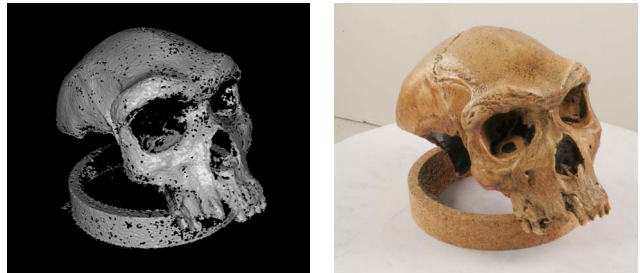


Figure 4. Example view from the *nskulla-half* dataset and the reconstructed mesh. Note the specular reflection on the skull surface and the reconstructed geometry in the eye socket.

4.5. Specular vs. Lambertian Surfaces

A textured, Lambertian surface is an optimal case for NCC-based window matching. The surface of the skull cast in the *nskulla* datasets is however quite specular and the lighting conditions are changing per-view. Figure 4 shows that our algorithm can nevertheless reconstruct a triangular mesh of the unoccluded regions of the skull. The individual depth maps were however quite incomplete so that the mesh contains many small holes.

The reconstruction of the *nskulla-small* dataset with standard parameters (*thresh* = 0.6, 5×5 window) yields comparably incomplete depth maps (see Figure 5). Lowering *thresh* to 0.4 improves coverage but includes also a large number of incorrect samples with high confidence value. All other datasets shown in the paper were therefore reconstructed with *thresh* = 0.6.

4.6. Silhouettes

Figure 6(b) shows a sample depth map from the *templeRing* dataset with spurious geometry around the silhouettes. Spurious geometry can be created in regions with low-contrast background where the windows are still dominated by the silhouette edge although they are centered on the background. The silhouette is clearly visible after the normalization step of the NCC (see Figure 6(d)). The normalized reference window matches the two windows from neighboring views shown on the left of Fig. 6(e) and 6(f)

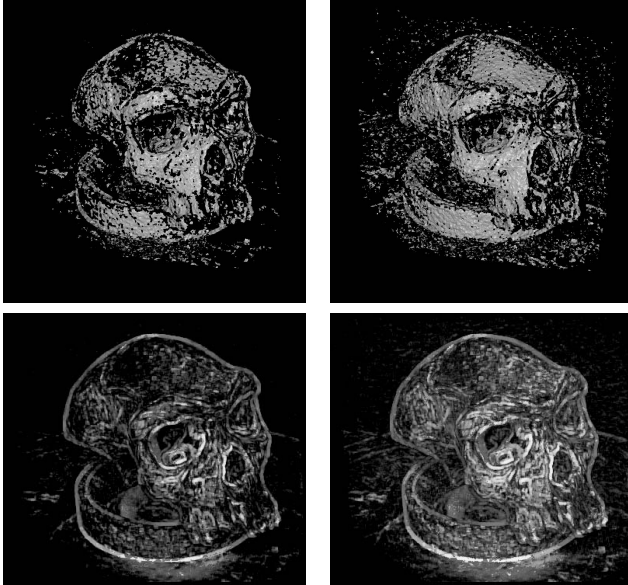


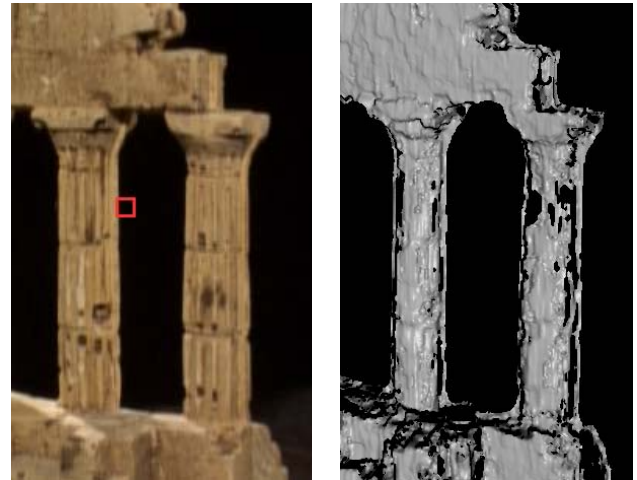
Figure 5. Comparison of different thresholds $thresh$ for a view from the *nskulla-small* dataset. *Top left*: $thresh = 0.6$ leads to large holes in the reconstruction, as shown in this rendering of the depth map. *Top right*: Lowering the threshold to $thresh = 0.4$ fills some of these holes but introduces strong noise. *Bottom*: Confidence values for both thresholds. Lowering the threshold from 0.6 (*left*) to 0.4 (*right*) increases the confidence in the noisy regions.

creating spurious geometry. The matching windows are dominated by the edge feature so that the NCC value is high despite the presence of background noise. The spurious geometry is therefore assigned a high confidence value and will most likely appear in the final geometry model. This problem arises particularly in the temple dataset due to the placement of the light sources relative to the object. Because silhouettes are easily determined in this dataset, we can reduce artifacts by omitting windows centered on the background.

4.7. Evaluation against Ground Truth

The reconstructed models from the temple and dino datasets were submitted to a multi-view stereo evaluation effort [14] which assessed accuracy and completeness of the reconstructed models. Accuracy is measured as the distance in mm such that 90 % of the points on a mesh are within this distance of the ground truth. The completeness value measures how well the reconstruction covers the ground truth. A more detailed explanation of these measures and the evaluation methodology appears in [14]. The evaluation results are summarized in Table 2. Figure 7 shows a comparison of an input view, our reconstructions for the three dataset sizes, and an image of the ground truth dataset used in the evaluation.

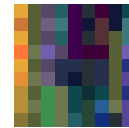
A full evaluation and ranking of the performance on all datasets is available in [14]. In summary, the accuracy of



(a) cropped reference view with approximate location of reference window (b) rendering of the corresponding depth map



(c) reference window



(d) normalized reference window



(e) corresponding windows in neighboring views



(f) normalized corresponding windows in neighboring views

Figure 6. Spurious geometry can occur at silhouettes with low image contrast. Fig. 6(a) and 6(b) show a reference view and the corresponding depth map. The seemingly structureless 9×9 window of the reference view (Fig. 6(c), marked in Fig. 6(a)) actually contains structure which is revealed by the normalization used in the NCC (Fig. 6(d)). The reference window matches the two windows from neighboring views shown on the left of Fig. 6(e) and 6(f) and spurious geometry is created along the edges of the columns.

the proposed algorithm was within a small margin of the best performing algorithm for the temple datasets. Accuracy and completeness clearly improved with the number of input images. For the more difficult *dino* datasets, our algorithm's accuracy actually decreased with the number of input images and achieved the best accuracy out of all methods for the smallest dataset (*dinoSparseRing*). This trend is caused by two factors. First, our algorithm reconstructs depth values for the diffuse, textureless dino only in areas where other features (e.g., geometric structures or shadow

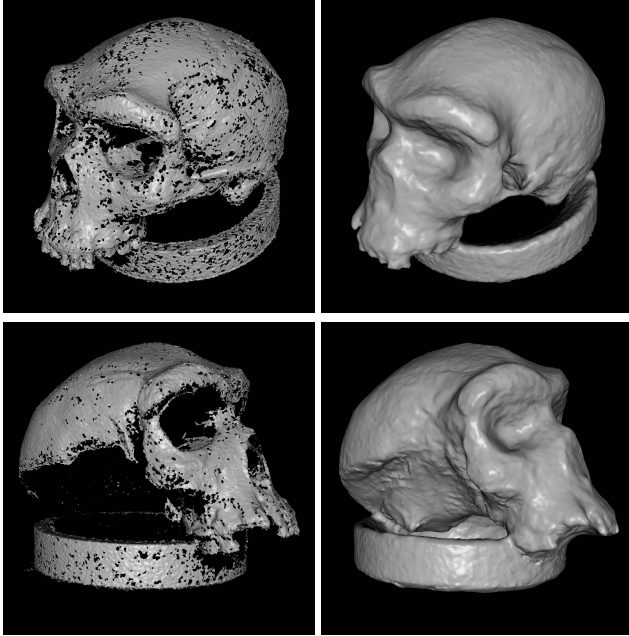


Figure 8. Comparison of the *nskulla-half* dataset reconstructed with our method (*left*) and with one of the currently top-performing multi-view stereo methods [5].

boundaries) allow for a high confidence match. Second, in contrast to the other participants, we do not estimate shape everywhere which avoids introducing inaccurate geometry. Our reconstruction of the *dinoSparseRing* dataset is thus most accurate but least complete.

Our current stereo implementation lacks any optimizations for speed and is consequently the slowest of all participants in [14]. This is, however, not a principle limitation as optimizations such as image rectification or hierarchical evaluation [6] can be easily applied.

5. Discussion

The previous section demonstrates that a conservative window-based multi-stereo algorithm can achieve results on par with the best current methods. Being conservative comes however at the price that the reconstructed models are incomplete if too few input views are available. The reconstructions from smaller datasets are consequently scoring worse in terms of completeness than the reconstructions from full datasets.

There is however also another way to look at the issue: Figure 8 compares the *nskulla-half* dataset reconstructed with our method on the left with one of the currently top-performing multi-view stereo methods [5]. Both reconstructions are (arguably) very good — but they also have very different properties. Our approach leaves holes of various sizes in areas of uncertainty but is able to reconstruct more complex geometry for this model (e.g., at the teeth on the left side of the jaw and in the eye sockets). In con-

trast, most modern multi-view stereo reconstruction methods reconstruct plausible, smooth, and well-behaved geometry even in areas where little or no information is available about the real object. Ultimately, it is a question of the application for which the model is generated and/or of a user's preference whether one or the other reconstruction philosophy is better.

Acknowledgements

We would like to thank Yasutaka Furukawa and Jean Ponce for providing the *nskulla* dataset. This work was supported in part by a Feodor Lynen Fellowship granted by the Alexander von Humboldt Foundation, NSF grant IIS-0413198, the University of Washington Animation Research Labs, the Washington Research Foundation, Adobe, Microsoft, and an endowment by Rob Short and Emer Dooley.

References

- [1] I. Cox, S. Hingorani, and S. Rao. A Maximum Likelihood Stereo Algorithm. *Computer Vision and Image Understanding*, 63(3):542–567, 1996.
- [2] B. Curless and M. Levoy. A volumetric method for building complex models from range images. In *SIGGRAPH*, pages 303–312, 1996.
- [3] O. Faugeras and R. Keriven. Variational principles, surface evolution, PDE's, level set methods and the stereo problem. *IEEE Trans. on Image Processing*, 7(3):336–344, 1998.
- [4] Y. Furukawa and J. Ponce. 3D photography dataset. Available at <http://www-cvr.ai.uiuc.edu/~yfurukaw/research/mview>.
- [5] Y. Furukawa and J. Ponce. Carved Visual Hulls for Image-Based Modeling. In *Proc. ECCV*, 2006.
- [6] C. Hernández. *Stereo and Silhouette Fusion for 3D Object Modeling from Uncalibrated Images Under Circular Motion*. PhD thesis, École Nationale Supérieure des Télécommunications, May 2004.
- [7] C. Hernández and F. Schmitt. Silhouette and stereo fusion for 3D object modeling. *Computer Vision and Image Understanding*, 96(3):367–392, 2004.
- [8] V. Kolmogorov and R. Zabih. Multi-camera scene reconstruction via graph cuts. In *Proc. ECCV*, pages 82–96, 2002.
- [9] P. J. Narayanan, P. Rander, and T. Kanade. Constructing virtual worlds using dense stereo. In *Proc. ICCV*, pages 3–10, 1998.
- [10] M. Okutomi and T. Kanade. A multiple-baseline stereo. *TPAMI*, 15(4):353–363, 1993.
- [11] M. Pollefeys, R. Koch, M. Vergauwen, and L. Van Gool. Metric 3D Surface Reconstruction from Uncalibrated Image Sequences. In *Proc. SMILE Workshop, LNCS 1506*, pages 138–153, 1998.
- [12] M. Pollefeys, L. Van Gool, M. Vergauwen, F. Verbiest, K. Cornelis, J. Tops, and R. Koch. Visual Modeling with a Hand-held Camera. *International Journal of Computer Vision*, 59(3):207–232, 2004.

dataset	accuracy	difference to best method	completeness	run time (hours:minutes)
templeFull	0.42 mm	0.06 mm	98.0 %	<i>200:00</i>
templeRing	0.61 mm	0.09 mm	86.2 %	<i>30:00</i>
templeSparseRing	0.87 mm	0.12 mm	56.5 %	10:06
dinoFull	0.56 mm	0.14 mm	80.0 %	<i>281:00</i>
dinoRing	0.46 mm	0.04 mm	57.8 %	<i>37:00</i>
dinoSparseRing	0.56 mm	—	26.0 %	12:24

Table 2. Results of the evaluation regarding accuracy, completeness, and run time. The third column lists the difference between the accuracy of our method and the accuracy of the best performing method in [14]. The heights of the objects are 159.6 mm (temple) and 87.1 mm (dino). Run times given for a 3.4 GHz Pentium 4 processor. Models whose run times are given in *italics* were computed on a PC cluster and timings were not directly available. The run times were therefore estimated based on the run times for the sparseRing datasets.

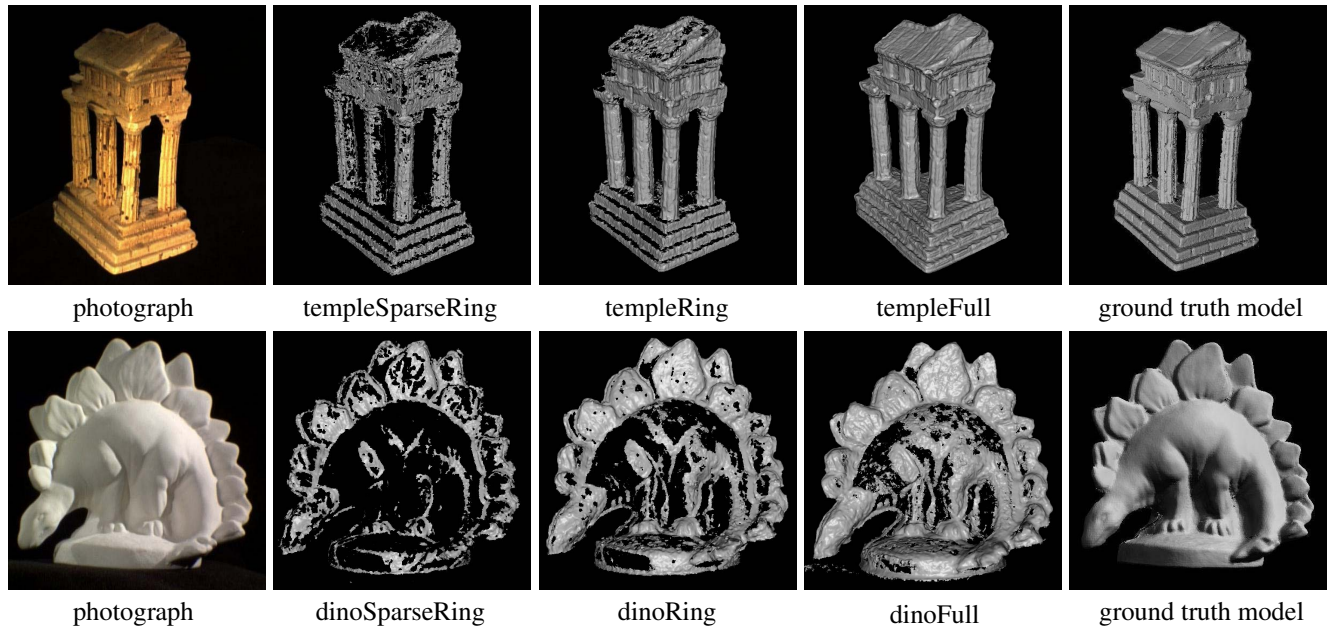


Figure 7. Models evaluated in the multi-view stereo evaluation [14]. *Left to right*: Photograph of the object, reconstruction from the sparseRing, ring, and full dataset, ground truth model used in the evaluation.

- [13] J.-P. Pons, R. Keriven, and O. Faugeras. Modelling dynamic scenes by registering multi-view image sequences. In *Proc. CVPR*, pages 822–827, 2005.
- [14] S. M. Seitz, B. Curless, J. Diebel, D. Scharstein, and R. Szeliski. A Comparison and Evaluation of Multi-View Stereo Reconstruction Algorithms. In *Proc. CVPR*, 2006.
- [15] S. Sinha and M. Pollefeys. Multi-view reconstruction using photo-consistency and exact silhouette constraints: A maximum-flow formulation. In *Proc. ICCV*, pages 349–356, 2005.
- [16] G. Vogiatzis, P. Torr, and R. Cipolla. Multi-view stereo via volumetric graph-cuts. In *Proc. CVPR*, pages 391–398, 2005.
- [17] Vripack: Volumetric range image processing package. Available at <http://grail.cs.washington.edu/software-data/vripack/>.
- [18] G. Zeng, S. Paris, L. Quan, and F. Sillion. Progressive surface reconstruction from images using a local prior. In *Proc. ICCV*, pages 1230–1237, 2005.

A. Normalized Cross-Correlation

We use a version of NCC for n -dimensional RGB color vectors v_0, v_1 with normalization per color channel. The vectors v_i correspond in our application to color values in an $n = m \times m$ window around a pixel position in a view V . To compute the NCC between two vectors v_0 and v_1 we first compute the average color value \bar{v}_i for each vector with $i \in \{0, 1\}$. We can then compute the NCC in a standard way as

$$NCC(v_0, v_1) = \frac{\sum_{j=0}^{n-1} (v_0(j) - \bar{v}_0) \cdot (v_1(j) - \bar{v}_1)}{\sqrt{\sum_{j=0}^{n-1} (v_0(j) - \bar{v}_0)^2 \cdot \sum_{j=0}^{n-1} (v_1(j) - \bar{v}_1)^2}}.$$

A multiplication between two color vectors is evaluated as dot product. The NCC returns a single value in the interval $[-1, 1]$ where 1 means highest correlation.