

Scaled-neighborhood Patches Fusion for Multi-view Stereopsis

Ning An , Yicong He, Hang Dong and Fei Wang

Institute of Artificial Intelligence and Robotics, Xi'an Jiaotong University, 710049, P. R. China

Abstract. In this paper, we present a multi-view stereo reconstruction approach which fuses scaled-neighborhood information. PMVS proposed by Furukawa is one of the most excellent algorithms, and it has a good performance on many datasets both the accuracy and the completeness. However, there are still further improvements on this algorithm. PMVS cannot perform well in the presence of slanted surfaces, which are usually imaged at oblique angles. According to these aspects, on the one hand we propose to estimate the initial normal of every seed patch via fitting quadrics with scaled-neighborhood patches, which greatly improves the accuracy of the normal. On the other hand, we present to compute scaled-window for the further optimization based on texture. And it has been tested that employing the scaled-window will dramatically smooth the surfaces and enhance the reconstruction precision.

1 Introduction

Multi-view stereo (MVS) is to reconstruct a complete 3D model from a collection of images taken from a scene or an object. Nowadays, huge numbers of images could be easily available with rapidly development of the modern digital camera and the Internet. It has seen a surge of interest in computer vision that how to exploit the largest and most diverse images collection ever assembled to reconstruct the 3D model about a scene [1, 2]. Much progress has been made, both in terms of precision and performance. Snavely et al. [2] proposed a method to automatically rebuild 3D large-scale scenes based on massive photographs from the Internet imagery, which realizes sparse 3D scene modeling from multi-view.

In the research of dense reconstruction, Furukawa[3] presented a novel algorithm, PMVS, based on patches which has been tested on various datasets. It has been proved to perform well on completeness and accuracy. In this paper, we revisit PMVS and show that it appears some limitations when the 3D modeling encounters slanted surfaces. Particularly, PMVS cannot achieve good result where the scene has both plain regions and curved regions. We will show that, with special design, this classical approach can yield some of the higher quality reconstructions.

The rest of this paper is organized as follows: we will first review related work (Section 2), and provide a more detailed overview of our method (Section 3). Section 3.1 briefly describes Furukawa's patch-based MVS algorithm. We then present the individual stages of our method, including initial estimation of patch normal via fitting quadrics with scaled-neighborhood patches (Section 3.2) and scaled-window for further optimization based on texture (Section 3.3). Experimental results and

discussions are given in Section 4. We conclude with results in Section 5.

2 Related work

In the process of MVS reconstruction, there are many factors influencing results such as slanted surfaces, low-texture regions, large concavities, which increase the difficulty for features matching across images. And, it is more difficult to reconstruct a 3D model where the local curvature is too high. And these factors have been always affecting the research of MVS.

Recently, the MVS problem has achieved a great development, yielding a variety of reconstruction algorithms. Most of methods need some prior information as auxiliary to get more precise results. According to the taxonomy of Seitz et al. [4], MVS algorithms can be divided into four categories: 3D volumetric approaches [5, 6], surface evolution techniques [7, 8], algorithms that compute and merge depth maps [9, 10], and techniques that grow regions or surfaces starting from a set of extracted features or seed points [3, 11], our algorithm apparently falls into the last category.

A multi-view framework for computing dense depth estimations was first proposed by Szeliski [12]. Esteban et al. [7] proposed a method based on texture and silhouette information, and fused the silhouette force into the snake framework. Based on these methods, Furukawa et al. [3] presented quite an accurate Patch-based MVS (PMVS) approach that starts from a sparse set of matched key points, and repeatedly expands these till visibility constraints are invoked to filter away false matches. The authors tested their method on the datasets provided by Strecha et al. [13], the available evaluation, and their

results were significantly more accurate and complete than the few other submitted ones.

Hiep et al. [14] proposed a minimum s-t cut based global optimization that transforms a dense point cloud into a visibility consistent mesh, followed by a mesh-based variational refinement that captures small details, smartly handling photo-consistency, regularization and adaptive resolution. In Bleyer’s [15] approach, a 3D scene is represented as a collection of visually distinct and spatially coherent objects. Inspired by Markov Random Field models of image segmentation, they employed object-level color models as a soft constraint, which can improve depth estimation in powerful ways. Jancosek et al. [16] augmented the existing Labatut CGF 2009 method with the ability to cope with these difficult surfaces just by changing the t-edge weights in the construction of surfaces by a minimal s-t cut. Their method uses Visual-Hull to reconstruct the difficult surfaces which are not sampled densely enough by the input 3D point cloud. They proved that the method can considerably better reconstruct difficult surfaces while preserving thin structures and details in the same quality. Qi Shan [17] leveraged occluding contours to improve the performance of multi-view stereo methods. This proposed approach outperforms state of the art MVS techniques for challenging Internet datasets, yielding dramatic quality improvements both around object contours and in surface detail.

We focus attentions on the research of improving 3D reconstruction on high curvature regions and thin structures. We look back on the patch-based method, PMVS, and on the framework of iterative optimization for every new seed patch, we propose an improvement approach that is inspired by the concept of multi-scale. It has been showed that, with careful design, this classical method can yield some of the higher quality reconstructions of any multi-view stereo algorithm, remaining highly scalable and offering excellent performance.

3 Algorithm Overview

This paper considers the problem of reconstructing the dense geometry of a 3D model from a number of images, and the camera poses and intrinsic parameters have been previously obtained. Methods about sparse multi-view stereo reconstruction and how to obtain the intrinsic parameters please refer to [1]. The dense MVS is a classic computer vision problem that has been extensively studied and a number of solutions have been published. Our multi-view reconstruction algorithm takes a set of calibrated images as input, which are captured from different viewpoints around the object to be reconstructed.

Decipted in Figure 1, similar to the process of PMVS, features being extracted are first matched across multiple images, yielding a sparse set of patches associated with salient image regions. The pairs of features are constrained to lie the corresponding epipolar lines, and triangulate the 3D points. We compute the normalized cross correlation (NCC) score based on scaled-window

whereas PMVS adopts fixed-window. The seed patches reconstructed by our approach are proved to be more steady, which reduces the impact of outliers in some degree.

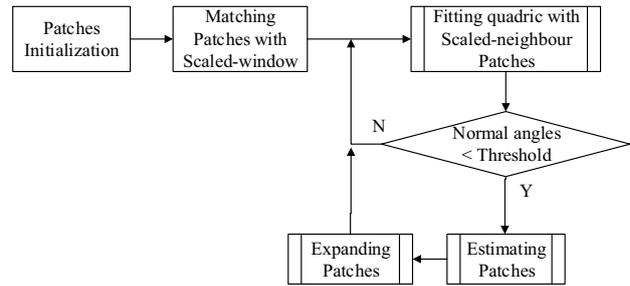


Figure 1. Algorithm framework.

PMVS estimates the surface orientation while enforcing the local photometric consistency, which is significant in practice to obtain accurate models for diversiform datasets. We employ this idea to compute scaled-window for the further optimization based on texture in the next iterations. In the step of optimization, we propose a design to fuse neighborhood priors to initially estimate the patch normal, that multi-scale quadrics are fitted according to K -neighborhood patches for acquiring the precise normal. It can be obviously improper in PMVS that the patch normal is assumed to parallel with the viewing ray passing through the center intersects the plane containing the patch. The application of the theory of multi-scale spatial is effectively improved the reconstruction accuracy with fine surface detail, especially on high curvature parts and thin structure regions.

3.1 Briefly Description of PMVS

PMVS is one of the most classical multi-view stereopsis algorithms. A patch p is essentially a local tangent plane approximation of a surface (Figure 2). Its geometry is fully determined by its center $c(p)$, unit normal vector $n(p)$ oriented toward the cameras observing it, a reference image $R(p)$ where p is visible and visible images $V^*(p)$.

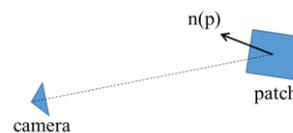


Figure 2. Patch model

It is implemented as a match, expand, and filter procedure. Features found by Harris and difference-of-Gaussians (DoG) operators are first matched across multiple images, yielding a sparse set of seed patches. Expansion is used to spread the initial matches to nearby pixels and obtain a dense set of patches. The geometric parameters, $c(p)$ and $n(p)$, are optimized by minimizing the photometric discrepancy score $g^*(p)$. It is repeatedly expanded before using visibility consistency and regularization constraint to filter away false matches.

$$c(p) \leftarrow \{Triangulation \text{ from } f \text{ and } f^i\} \quad (1)$$

$$n(p) \leftarrow \frac{\overline{c(p)O(I_i)}}{|\overline{c(p)O(I_i)}|} \quad (2)$$

$$R(p) \leftarrow I_i \quad (3)$$

$$V(p) \leftarrow \{I \mid n(p) \cdot \overline{c(p)O(I)} / |\overline{c(p)O(I)}| > \cos(\theta)\} \quad (4)$$

$$V^*(p) = \{I \mid I \in V(p), h(p, I, R(p)) \leq \alpha\} \quad (5)$$

$$g^*(p) = \frac{1}{|V^*(p) \setminus R(p)|} \sum_{I \in V^*(p) \setminus R(p)} h(p, I, R(p)) \quad (6)$$

To simplify computations, PMVS constrains $c(p)$ to lie on a ray such that its image projection in one of the visible images and $n(p)$ is assumed to parallel with the viewing ray passing through the center which intersects the plane containing the patch. Considering the reconstruction of high curvature and slanted surfaces, the initial estimation of patch normal determined by any visible image diverges from the true normal. It cannot make any sense by minimizing the photometric discrepancy score for the optimization of the geometric parameters. In addition, there are both plain regions and fold regions existing, and PMVS employs the fixed-window to compute photometric discrepancy scores for the further optimization, which is unable to perform well.

3.2 Normal Estimation by Scaled-Quadric

Our algorithm follows the framework of patch-based PMVS, however, some key designs are proposed. We compensate for estimating normal of seed patch by employing fitting quadrics with scaled-neighborhood patches, which improves the quality especially in high curvature regions and slanted surfaces. From the Figure 3, the dashed lines are the initial estimation of patch normal in PMVS. It can be assumed almost true when the reconstructed surface exactly faces towards cameras, while some high curvature regions exist unavoidable.

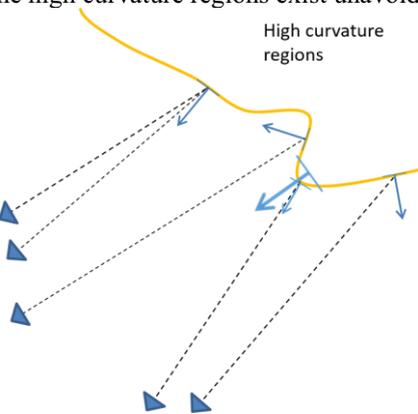


Figure 3. Cameras and the surface of 3D scenes

Shi Li-Min[18] has found the deficiency in PMVS, and proposed a method that K neighborhoods around the seed patch are selected to fit quadrics. The normals of fitted quadrics are computed to be as the initial estimation

normal of the seed patch. This method rectifies the estimation of patch normal but still has some limitations. The value of K is set beforehand, and if the reconstructed region is plain, some outliers will lead to the deviation when the quadrics are fitted without enough patches. The large K neighborhood patches has certain endurance to fault tolerance when the plain regions are reconstructed. From the Figure 3, the normal computed with the large K neighborhood patches has some deviations while the smaller K will be more accurate.

We propose an approach to search multi-scale K neighborhoods around the seed patch that more effective patches are acquired to fit quadrics. It can be found more appropriate quadrics to estimate the normal when the high curvatures are reconstructed. Firstly, we search $\sigma \times K$ neighborhood patches to fit a quadric and search $\lambda \sigma \times K$ neighborhood patches to fit another quadric in a lower scale (Figure 4). Then we get two normals $V1, V2$. When the two normals satisfy the formula (7), the normals are regarded as the same direction, and the normal computed in higher scale is set to be the initial estimation of the seed patch. Otherwise, $\lambda^2 \sigma \times K$ neighborhood patches are searched to compute the normal $V3$, checking whether $V2$ and $V3$ meet the threshold. And repeat the above process. If the two normals still have comparatively large deviation, the normal computed in smaller scale is used to as the initial estimation of the seed patch. Right now the new seed patch already has a relatively accurate initial normal, which can dramatically improves the further optimization of the geometric parameters, $c(p)$ and $n(p)$.

$$\cos^{-1}(\text{normal1}, \text{normal2}) < \text{AngleThreshold} \quad (7)$$

3.3 Scaled-Window for Further Optimization

As mentioned above section, it is imperative that the idea of adaptive multi-scale is introduced in the complex 3D reconstruction. During the process of optimization based on texture, it obtains amazing results when we employ this concept to rectify the texture window. $c(p)$ is constrained to lie on a ray such that its image projection in one of visible images does not change, reducing its number of degrees of freedom to one and solving only for a depth. $n(p)$ is parameterized by Euler angles, yaw and pitch. We use normalized cross correlation (NCC) on the square texture region as a metric for minimizing the photometric discrepancy.

$$NCC(p_0, p_1) = \frac{\sum_i (p_0(i) - \bar{p}_0) \cdot (p_1(i) - \bar{p}_1)}{\sqrt{\sum_i (p_0(i) - \bar{p}_0)^2 \cdot \sum_i (p_1(i) - \bar{p}_1)^2}} \quad (8)$$

It has a certain limitation that PMVS algorithm employs the fixed-window for optimization. The fixed-window is not usually large, or some holes will appear on the reconstructed surfaces because of high curvature regions. However, some plain regions (building wall, ground) exist, and the reconstructed surface will be extremely corrugate if the fixed-window is not large enough. The 3D point cloud looks just like fish scale, which has seriously affected the accuracy of stereo reconstruction.

According to every patch, the scaled-window is selected independently. We choose large texture windows for the plain regions, whose NCC scores could maintain the stability and the 3D model gets more smooth. During the reconstruction of fold regions, our approach adaptively selects smaller texture window for optimization. It maintains the sensibility to high curvature regions, which helps to reconstruct thin structures well. The algorithm is intuitively described in Figure 4.

```

For each new patch candidate  $p'$ 
  Search  $\sigma \times K$  neighborhood Patches;
  Search  $\lambda\sigma \times K$  neighborhood Patches;
  Fit quadratic and compute the normal  $V1, V2$ ;
  If  $\text{acos}(V1, V2) < \text{AngleThreshold1}$ 
     $n(p') = V2$ , window_size =  $WZ$ ;
  else if  $\text{acos}(V1, V2) < \text{AngleThreshold2}$ 
     $n(p') = V1$ , window_size =  $WZ \times \delta$ ;
  else
    Search  $\lambda^2\sigma \times K$  neighborhood Patches
    Fit quadratic and compute the normal  $V3$ ;
    If  $\text{acos}(V1, V3) < \text{AngleThreshold3}$ 
       $n(p') = V1$ , window_size =  $WZ \times \delta$ ;
    else
       $n(p') = V3$ , window_size =  $WZ \times \delta^2$ ;

```

Figure 4. Estimation of the normal via fitting quadrics with scaled-neighborhood patches and optimization based on scaled-window. $K \in [30, 60]$, $\sigma=2$, $\lambda=0.6$, $\delta=0.8$

4 Results

Our algorithm is implemented by VC++ with the CGAL library. We demonstrate our multi-view stereo algorithm with a number of reconstructions. Images are acquired from the real-world scenes during the experimental process, and we have some particular comparison with the result of PMVS. It is mainly analyzed on three 3D models from the reconstruction. The first set of data is taken from a sculpture, and the other two are the open datasets from Strecha [13]. When we reconstruct a scene, it is shot from various angles firstly, yielding a set of pictures. Then structure-from-motion software [19] is applied to restore the position of every camera and get the P matrix for per visible image. We take the P matrix as input of the further dense 3D reconstruction.

The first experiment is carried on a sculpture of typography with nine images (3264×2448 resolution), called “Typography-P9”. Typography-P9’s surface is composed of many slanted planes, and parts of them have characters carved with thin structures. PMVS builds 355142 patches while our algorithm obtains 423148 patches with ten percent more. The denseness reconstructed has improved a lot (Table 1). At the same time, from the comparison of 3D models, our approach performs obviously better than PMVS on both accuracy and completeness on the local region. The regions with red circles (see Figure 5, 6) appear large holes because these areas are slanted and close to the edge of visibility, with small texture for further optimization. And our method conducts an estimation fusing multi-scale neighborhood patches to rectify the normal of a new seed

patch, with adaptive scaled-window for optimization, regarding some new patches supported by small size of texture. It apparently shows that our 3D reconstruction of fold regions can be seen very smooth, with more dense point cloud, in contrast (Figure 5).

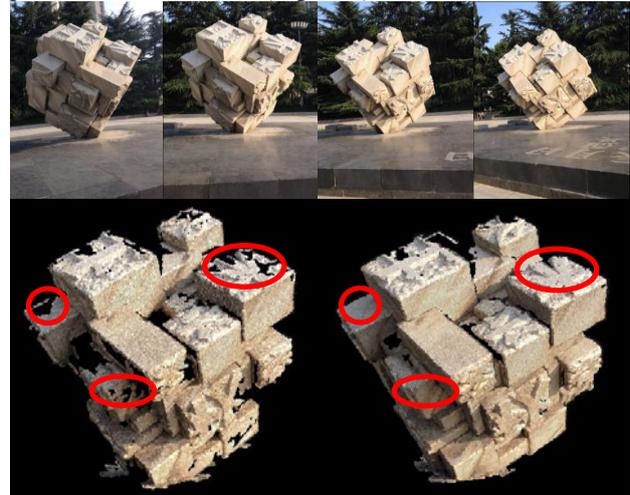


Figure 5. First row : some example images of Typography-P9. 3D point cloud in the last row are computed using PMVS and our algorithm.

The second experiment is on the dataset Herz-Jesu-P8 [13]. The eight images have been precisely calibrated corresponding cameras for their intrinsic and extrinsic parameters. We can directly obtain the P matrix for every image. Our method produces 1565694 patches with 25.7 percent more on the denseness of 3D point cloud (Figure 6, Table 1). This scene has a large flat building wall, which leads to the fish-like scale of the 3D model reconstructed by PMVS with small fixed-window, while large window may lead to some holes on the thin structure surface. Our adaptive approach will choose the large texture window for optimization when the scene has large plain regions. It is the careful design that helps to achieve the fine rebuilding result.

We adopt the dataset of Castle-P10 [13] to make another experiment. This scene has an obvious character of bending building wall, which is more apparently proved the limitation of PMVS. From the 3D modeling, the PMVS’s reconstruction of the front building wall facing towards cameras is just similar to ours, only deficient on some regions with thin structures. However, the bending building wall is almost impossible to rebuild by PMVS. By comparison, our method can not only build the bending building wall, but also have an accurate recovery of carvings, embossments and windows on details (Figure 6).

5 Conclusion

The experiments are proved that the method proposed by us has a better performance on the high curvature regions and slanted surfaces. In comparison to PMVS, the 3D models reconstructed by our algorithm have a great improvement on the accuracy, denseness and thin

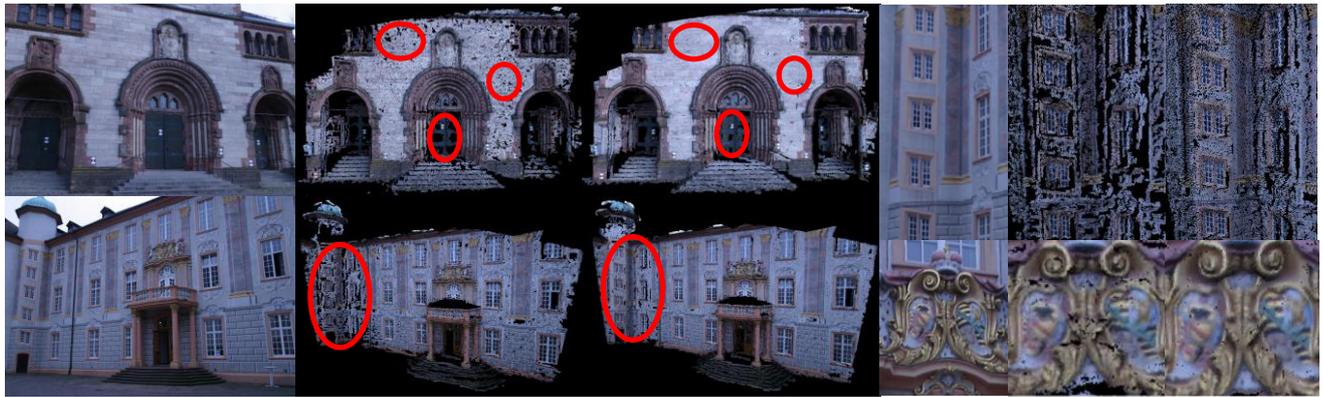


Figure 6. Datasets from Strecha [13]. Top row: from left to right, an image from Herz-Jesu-P8, the 3D model reconstructed by PMVS, by Ours, the bending building wall and the 3D model reconstructed by PMVS, by Ours. Bottom row: from left to right, an image from Castle-P10, the 3D model reconstructed by PMVS, by Ours, the embossment and the 3D model reconstructed by PMVS, by Ours.

structures. The estimation of patch normal via fitting quadrics with scaled-neighborhood patches enhances the sensibility to the high curvature regions. And the application of multi-scale also improves the stability of the normal estimation. It is presented to be more efficient to rebuild a scene existing both plain regions and fold regions with scaled-window for the further optimization based on texture. The adaptiveness guarantees the patches to be optimal during iterations, which leads to a more accurate and complete 3D model.

Table 1. Comparisons of denseness

Datasets	Resolution	Number of patches (PMVS)	Number of patches (Ours)	Percent
Typography-P9	3264×2448	355142	423148	19.2%
Herz-Jesu-P8	3072×2048	1246046	1565694	25.7%
Castle-P10	3072×2048	1540772	1863400	20.9%

Acknowledgment

This work was supported by the Natural Science Foundation of China (No61273366), National High Technology Research and Development Program (No2015BAH31F01).

References

1. N. Snavely, S.M. Seitz, R. Szeliski S. Modeling the world from internet photo collections. *International Journal of Computer Vision* **80**(2) :189-210. (2008).
2. S. Agarwal, Y. Furukawa, N. Snavely, et al. Building rome in a day. *Communications of the ACM* **54**(10). (2011).
3. Y. Furukawa, and J. Ponce. Accurate, dense, and robust multiview stereopsis. *Pattern Analysis and Machine Intelligence*. **32**(8) :1362-1376. (2010).
4. S.M. Seitz, B. Curless, et al. A comparison and evaluation of multi-view stereo reconstruction algorithms. *Computer Vision and Pattern Recognition*. Vol. **1** :519-528. (2006).
5. M. Sormann, C. Zach, J. Bauer, et al. Watertight multi-view reconstruction based on volumetric graph-cuts. *Image analysis*. 393-402. (2007).
6. G Vogiatzis, et al. Multiview stereo via volumetric graph-cuts and occlusion robust photo-consistency. *Pattern Analysis and Machine Intelligence*, **29**(12) :2241-2246. (2007).
7. C.H. Esteban, F. Schmitt. Silhouette and stereo fusion for 3D object modeling. *Computer Vision and Image Understanding*. **96**(3) :367-392. (2004).
8. A. Zaharescu, E. Boyer, R. Horaud. Transformesh: a topology-adaptive mesh-based approach to surface evolution. *Computer Vision-ACCV* 166-175. (2007).
9. C. Strecha, W. von Hansen, L. Van Gool, P. Fua, U. Thoennessen Combined depth and outlier estimation in multi-view stereo. *Computer Vision and Pattern Recognition*. Vol. **2** :2394-2401. (2006).
10. Shen, Shuhan, and Z.Y. Hu. How to Select Good Neighboring Images in Depth-Map Merging Based 3D Modeling. *Image Processing*. **23**(1). (2014).
11. D. Bradley, T. Boubekeur, and W. Heidrich. Accurate multi-view reconstruction using robust binocular stereo and surface meshing. *Computer Vision and Pattern Recognition*. 1-8. (2008).
12. R. Szeliski. A multi-view approach to motion and stereo. *Computer Vision and Pattern Recognition*. Vol. **1**.(1999).
13. C. Strecha, W. von Hansen, L. Van Gool. On benchmarking camera calibration and multi-view stereo for high resolution imagery. *Computer Vision and Pattern Recognition*. 1-8. (2008).
14. V.H. Hiep, R. Keriven, P. Labatut. Towards high-resolution large-scale multi-view stereo. *Computer Vision and Pattern Recognition*. 1430-1437. (2009).
15. M. Bleyer, C. Rother, P. Kohli. Object stereo-joint stereo matching and object segmentation. *Computer Vision and Pattern Recognition*. 3081-3088. (2011).
16. J. Michal, and T. Pajdla. Multi-view reconstruction preserving weakly-supported surfaces. *Computer Vision and Pattern Recognition*. 3121-3128. (2011).
17. Q Shan, B. Curless, Y. Furukawa, C. Hernandez, S.M. Seitz. Occluding contours for multi-view stereo. *Computer Vision and Pattern Recognition*. 4002-4009. (2014).

18. L.M. Shi, F.S. Guo, Z.Y. Hu. An Improved PMVS through Scene Geometric Information. *Acta Automatica Sinica*, **37.5**(2011).
19. C. Wu. Towards linear-time incremental structure from motion. *3D Vision*, 127-134. (2013).