



Mesh-based Camera Pairs Selection and Occlusion-Aware Masking for Mesh Refinement

Andrea Romanoni^{a,**}, Matteo Matteucci^a

^aPolitecnico di Milano, Italy

ABSTRACT

Many Multi-View-Stereo algorithms extract a 3D mesh model of a scene, after fusing depth maps into a volumetric representation of the space. Due to the limited scalability of such representations, the estimated model does not capture fine details of the scene. Therefore a mesh refinement algorithm is usually applied; it improves the mesh resolution and accuracy by minimizing the photometric error induced by the 3D model into pairs of cameras. The choice of these pairs significantly affects the quality of the refinement and usually relies on sparse 3D points belonging to the surface. Instead, in this paper, to increase the quality of pairs selection, we exploit the 3D model (before the refinement) to compute five metrics: scene coverage, mutual image overlap, image resolution, camera parallax, and a new symmetry term. To improve the refinement robustness, we also propose an explicit method to manage occlusions, which may negatively affect the computation of the photometric error. The proposed method takes into account the depth of the model while computing the similarity measure and its gradient. We quantitatively and qualitatively validated our approach on publicly available datasets against state of the art reconstruction methods.

© 2019 Elsevier Ltd. All rights reserved.

1. Introduction

Multi-View Stereo (MVS) algorithms recover the 3D model of a scene captured by a set of images. Boosted by the benchmarks proposed in [1, 2, 3] and the enhancements in hardware capabilities, several works proposed several accurate and efficient MVS methods.

A well established MVS pipeline, first proposed by Vu *et al.* [4], estimates the camera positions with Structure-from-Motion (SfM) [5, 6]; then it applies plane sweeping [7] or depth map fusion [8] to recover a dense point cloud representation of the scene. This pipeline builds a volumetric partitioning of the space in which the camera to point visibility rays are exploited to estimate free and occupied space (or an implicit representation of the model, such as Truncated Signed Distance Function); the free-occupied space boundary (or the zero crossing surface) is the model of the observed scene, usually represented by means of a mesh. To obtain a very accurate reconstruction,

the last step of the pipeline is a refinement algorithm; it minimizes the photometric error induced by such mesh in pairs of cameras.

A fundamental aspect of a mesh refinement algorithm, and in general a Multi-View Stereo method, is the choice of the camera pairs used to compute and minimize the photometric error. It is well known that too narrow cameras imply noisy reconstruction results. On the other hand, images captured by cameras too far from each other could have limited overlap, *i.e.*, the region of the scene perceived by both cameras is small. The right choice of these pairs leads to a coherent computation of the photometric error and, as a consequence, an effective gradient descent minimization. The most widespread Multi-View Stereo methods select for each camera a pairing view among the others by relying on several factors. Li *et al.* [9] and Ebner *et al.* [10] evaluate the baseline and the angle of the principal viewing direction between the cameras; instead, other methods [11, 12, 13, 14] consider the SfM 3D points and take into account the average angle between the camera-to-point viewing rays, the baseline among views and the scale. Vogiatzis *et al.* [15] leverage on similar metrics to filter out unreliable photometric measures, adopted to estimate the 3D model.

^{**}Corresponding author:
e-mail: andrea.romanoni@polimi.it (Andrea Romanoni)

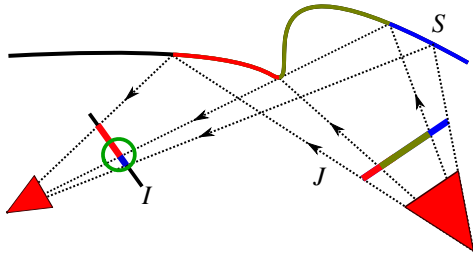


Fig. 1: Effect of the occlusion in the error computation

While pixel-wise camera pair selection has been addressed to estimate accurate point clouds in [16, 17, 18], mesh refinement literature has always limited the choice to the pairs of cameras sharing the visibility of the highest number of 3D points with sufficient parallax.

A second issue which affects mesh refinement and we address in this paper is related to model occlusions. For instance, in Fig. 1 image J projects in I through S ; the patch in the green circle contains a discontinuity. In this case, while computing the projection error in the green patch corresponding to a pixel in the red region, state-of-the-art methods consider both the information from the blue and red (non-adjacent) parts of S . This issue has been considered only when dealing with generative methods [19] or with a simple heuristic avoiding to evolve the mesh just in correspondence of edge segments joining visible and non-visible facets ([20]).

For these reasons, in this paper we propose three contributions:

- A pairwise camera selection method exploiting the knowledge of an approximate model of the scene. It minimizes an energy function defined over the surface instead of just relying on camera poses or 3D points viewing angles (Section 3).
- An energy term to favor symmetric pairs and better compensate image noise while computing the gradient flow (Section 3).
- An occlusion-aware mask term to explicitly identify, for each pixel, which part of the neighborhood has to be considered, during mesh refinement, to compute the similarity measure and its gradient (Section 4).

2. Related works

Mesh refinement is a case of surface evolution. Surface evolution methods are framed into the variational framework formalized by Hermosillo *et al.* [21]. Early works represent the model by level set, *i.e.*, as the zero crossing of a function $f : \mathbf{R}^3 \rightarrow \mathbf{R}$ [22, 23, 24, 25, 26, 27, 28]. Faugeras and Keriven *et al.* [22] define the level set by means of a partial differential equation of f such that a point on the surface moves proportionally to the photo-consistency of its neighborhood. Jin *et al.* [23] and Yoon *et al.* [27] extend this approach to cope with specular reflection. While these methods integrate the photometric measure in the 3D domain, Yezzi and Soatto [24] show that integrating this measure on the image domain yields to more accurate

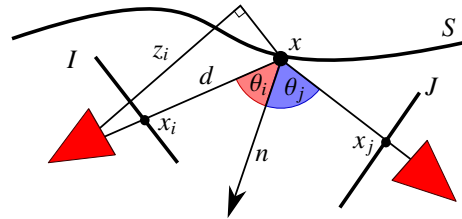


Fig. 2: Variables involved in the photometric refinement process

results. Solem *et al.* [26] and Pons *et al.* [29, 28] replace the partial differential equation with a gradient leading to more robust mesh evolution. Finally, Fuhrmann *et al.* [25] adapts the considered neighborhood around surface points according to their scale. Even if level set methods achieve accurate results, the evolution process is not always easy to track and understand.

Differently from level set methods, mesh refinement algorithms directly represent the surface as a 2D mesh embedded into a 3D space [8, 30, 19, 31, 32, 33]: given an initial rough mesh of the scene, they move the position of its vertices to obtain a more faithful model. Vu *et al.* [8] discretized the continuous level set formulation of [28] to work directly with meshes. Delaunoy *et al.* [32] extended this method to take into account occlusions and in [34] they jointly optimize the surface and the camera in a bundle adjustment fashion. Li *et al.* [35] proposed an improved smoothness term of the energy function to output very smooth surfaces keeping the details of the scene. Recently, Li *et al.* [20] simplify the mesh, decreasing the resolution where few vertices are sufficient to capture the structure of the scene, without affecting the accuracy significantly. In [36], photometric refinement is coupled with a moving object detection method to avoid using their image areas to refine the static model of the scene. Finally, two mesh refinement approaches [37, 38] exploit the semantic 2D segmentation of the images to improve the robustness of the refinement process, especially where two objects of different classes are adjacent.

Mesh Refinement. Mesh refinement takes as input an initial mesh which is a rough model of the scene and a set of images capturing the scene. The most popular approach, proposed in [4], minimizes an energy function E :

$$E = E_{\text{photo}} + E_{\text{smooth}}, \quad (1)$$

where E_{photo} represents the photo-consistency error of the model with respect to the images, and E_{smooth} enforces the smoothness of the surface.

To minimize the term E_{photo} , the mesh refinement procedure applies gradient descent. Let's consider two images I and J , and a point x belonging to the surface S (Fig. 2); we define the error function $err_{A,B}(x)$ that decreases if the similarity between the patch around the projection of x in A and B increases, *e.g.*, in our case, the negative ZNCC of the 5×5 pixels neighborhood. The energy E_{photo} in Equation (1) is defined as:

$$E_{\text{photo}} = E(S) = \sum_{i,j} \int_{\Omega_{i,j}^S} err_{I,J}(x_i) dx_i = \sum_{i,j} \mathcal{E}_{i,j}^{im}(x), \quad (2)$$

where $I_{i,j}^S$ represents the reprojection of the image from the j -th camera onto image I through the surface S , and $\Omega_{i,j}^S$ is the image

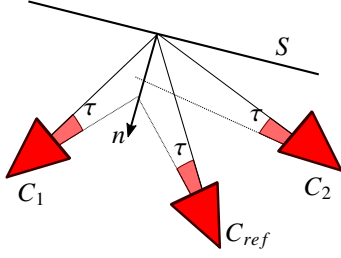


Fig. 3: Example of symmetric (C_{ref} and C_1) and non symmetric cameras (C_{ref} and C_2) with respect to the surface normal n .

region where the reprojection is defined. Now, let $X_i \in \mathbb{R}^3$ be a vertex of the surface mesh S , and $\phi_i(x)$ be the barycentric coordinates of a surface point x in the triangle with vertex X_i . The discrete gradient of $E_{photo} = E(S)$ computed for a vertex X_i is:

$$\begin{aligned} \frac{dE(S)}{dX_i} &= \int_S \phi_i(x) \nabla E(S) dx = \int_S \phi_i(x) \nabla \left(\sum_{i,j} \mathcal{E}_{ij}^{im}(x) \right) dx = \\ &= \int_S \phi_i(x) \sum_{i,j} \nabla \mathcal{E}_{ij}^{im}(x) dx. \end{aligned} \quad (3)$$

By changing the variable of integration from x to x_i it is possible to integrate the energy over the image I instead of over the surface S . Let define \vec{n} as the normal at x pointing outward the surface S , x_i the projection of x into the I image, \mathbf{d}_i as the vector from camera i to x , z_i as the depth of x in camera i (see Fig. 2); with the change of variable $dx_i = -\vec{n}^T \mathbf{d}_i dx / z_i^3$ [28] we obtain:

$$\frac{dE(S)}{dX_i} = \sum_{i,j} \int_{\Omega_{ij}^S} \phi_i(x) \nabla \mathcal{E}_{ij}^{im}(x) \frac{z_i^3}{\vec{n}^T \mathbf{d}_i} dx_i. \quad (4)$$

To define which pairs (i, j) are adopted to compute the gradients, the most widespread methods [12, 4, 39] leverage on 3D points correspondences estimated by the Structure from Motion method adopted to calibrate the cameras. For each camera i it chooses the camera j with the highest number of common 3D points with a reasonable parallax (e.g., in [12] between 10° and 30°).

Finally, the evolution process is complemented by the umbrella operator [40], which moves each vertex in the mean position of its neighbors; this approximates the Laplace-Beltrami operator, and it minimizes the energy term E_{smooth} .

3. Model-based Camera Pairs Selection

One of the most relevant aspect to effectively minimize E_{photo} is the choice of the camera pairs (i, j) . Instead of basing this choice on just 3D sparse points estimated by Structure from Motion, we propose to exploit the knowledge of the initial 3D mesh to find the camera pairs having a good trade-off between reasonable parallax and image overlap.

We first define a term $E_p^{i,j}$ to represent the quality of the parallax between camera i and j . Given the initial surface S , the

camera centers, c_i and c_j , and a point $x \in S$, let's define the parallax as the angle $\theta_{i,x,j} = \angle(c_i x, c_j x)$. We compute the average parallax as:

$$P_{i,j} = \frac{1}{A(\Omega_{ij}^S)} \int_{\Omega_{ij}^S} \theta_{i,\Pi_i^{-1}(x_i),j} dx_i \quad (5)$$

where $A(\cdot)$ represents the area of an image region, in this case the region Ω_{ij}^S . To define the reference parallax, we recall that small angles induce good overlap between patches, while larger angles induce more stability in the refinement process. Tola *et al.* [12] choose a small parallax between 10° and 30° to avoid erroneous image warping caused by occlusions. In our case, however, we know the geometry, and we explicitly handle occlusions, therefore we prefer a larger reference angle. In [41] the parallax ranges around 40° to 70° ; in [42] and in [43] the convergent angle, which is strictly related to the parallax, is chosen respectively as 50° and 45° . Among these values, we experimentally choose the reference parallax to be 50° . Therefore, we define:

$$E_p^{i,j} = -e^{-\left(\frac{P_{i,j}-50^\circ}{2\sigma_p}\right)^2}, \quad (6)$$

where we put the variance $\sigma_p = 45^\circ$.

To favor camera pairs with a similar resolutions and thus inducing coherent refinement, we introduce the resolution term $E_R^{i,j}$. Let $\rho_i = \frac{\|c_i x_i\|}{f_i}$ and $\rho_j = \frac{\|c_j x_i\|}{f_j}$ be the distances of point x_i from the two cameras respectively, normalized with respect to the corresponding focal length f_i and f_j . We define:

$$\rho_{(i,j)} = \frac{|\rho_i - \rho_j|}{\|c_i x_i\|}, \quad (7)$$

as the normalized discrepancy of the resolutions of the two images with respect to the length of $c_i x_i$ ray. We compute the average of these values as:

$$R_{i,j} = \frac{1}{A(\Omega_{ij}^S)} \int_{\Omega_{ij}^S} \rho_{(i,j)} dx_i. \quad (8)$$

To favor similar resolutions we define:

$$E_r^{i,j} = -e^{-\left(\frac{R_{i,j}-0}{2\sigma_r}\right)^2}, \quad (9)$$

where we put the variance $\sigma_r = 0.25$, which represents a resolution discrepancy of 25%.

Finally, to take into account overlap, we define $E_o^{i,j}$ as:

$$E_o^{i,j} = -\frac{A(\Omega_{ij})}{A(I_i)}. \quad (10)$$

Symmetry Term. In most cases, these two terms provide a fair evaluation of the camera pair quality with respect to the mesh refinement problem. However, in some cases, they are not sufficient to discriminate among different camera pairs. For instance, in Fig. 3, the surface S is perceived by a reference camera C_{ref} and by two other cameras C_1 and C_2 . The surface is entirely visible by the three cameras, *i.e.*, the overlap is 100%, and the baselines $C_{ref} - C_1$ and $C_{ref} - C_2$ have very similar

values; by relying on just E_p and E_o , both pairs C_{ref}, C_1 and C_{ref}, C_2 are considered equally good (or equally bad).

To overcome this issue, we evaluate a third term that favors camera pairs with points of view symmetric with respect to the scene. Intuitively, since the surface S evolves along its normal \vec{n} and assuming the images affected by Gaussian noise on the image plane, if we compute the gradient between I_{ref} captured by C_{ref} and I_2 captured by C_2 , similar noise in I_{ref} and in I_2 (inducing an uncertainty angle τ) translate into significantly different gradient noises along \vec{n} . Instead, if we consider I_{ref} and I_1 the same noise affects similarly the gradients along \vec{n} . Statistically, in the former case, the noise accumulates as the gradients are computed, while, in the latter, they likely compensate each other.

In addition to parallax and overlap we then evaluate a symmetry term $E_s^{i,j}$ with respect to the surface normal. To do so we define the oriented angle difference (OAD): $\delta_{i,x,j} := \text{sign} \cdot \frac{1}{2} [\angle(c_i x, \vec{n}(x)) - \angle(c_j x, \vec{n}(x))]$, where $\vec{n}(x)$ is the normal of the surface S at x and, if the $c_i x$ and $c_j x$ belong to the same half-space defined by the plane parallel to \vec{n} through x , then $\text{sign} = 1$, otherwise $\text{sign} = -1$. The OAD average on the surface is computed as:

$$S_{i,j} = \frac{1}{A(\Omega_{i,j}^S)} \int_{\Omega_{i,j}^S} \delta_{i,\Pi_i^{-1}(x_i),j} dx_i, \quad (11)$$

and, the novel energy term $E_s^{i,j}$ is computed as: $E_s^{i,j} = -e^{-\left(\frac{S_{i,j}-0^\circ}{2\sigma_s}\right)^2}$, where we put experimentally the variance $\sigma_s = \sigma_p = 45^\circ$. This term is combined with $E_p^{i,j}$ and $E_o^{i,j}$ to define the energy function $E_{BPV}(i, j)$ for a pair of cameras C_i and C_j :

$$E_{BPV}^{i,j} = \mu_1 \cdot E_p^{i,j} + \mu_2 \cdot E_o^{i,j} + \mu_3 \cdot E_s^{i,j} + \mu_4 \cdot E_r^{i,j}, \quad (12)$$

where μ_1, μ_2 and μ_3 are three coefficient weighting the contribution of the three energy term (in our case $\mu_1 = 0.25, \mu_2 = 0.25, \mu_3 = 0.5$ and $\mu_4 = 0.25$).

Model Coverage. A second relevant aspect when choosing a camera pair is model coverage. While the overlap term is related to the overlap among the images in the pair, in principle no terms discussed previously avoids all the camera pairs perceive, and therefore refine just a small portion of the mesh. For this reason, we enforce camera pair configurations providing good coverage.

We first initialize the camera pairs as follows. Let \mathbb{C} be the set of cameras and \mathbb{P} the set of camera pairs adopted for the refinement. Our algorithm initializes the initial set \mathbb{P}^0 of best camera pairs computed leveraging on the previous energy as $\mathbb{P}^0 = \{(i, j) \forall i \in \mathbb{C} \text{ s.t. } j = \arg \min_j \{E_{BPV}^{i,j}\}\}$, i.e., with the best pair for each camera.

To enforce model coverage, the idea is to perturb this initial set of pairs \mathbb{P}^0 . In the first step we define the model coverage as the average number of camera pairs in which all the facets are visible. Let \mathbb{F} be the set of facets and let define a visibility function $v_f^{i,j}$ of facet $f \in \mathbb{F}$ with respect to cameras i and j , i.e. $v_f^{i,j} = 1$ if is visible from both cameras, $v_f^{i,j} = 0$ otherwise.

Algorithm 1 Camera Pairs Selection

Input:

\mathbb{C} , set of N_{cam} cameras, \mathcal{M} 3D mesh, \mathbb{P}^0

Output:

\mathbb{P} , set of camera pairs

Algorithm:

Initial Camera Pairs Selection:

- 1: **for** $i \in \mathbb{C}$ **do**
 - 2: $\mathbb{P}^0 \leftarrow (i, \arg \min_j \{E_{BPV}^{i,j}\})$
 - 3: **end for** Camera Pairs Selection:
 - 4: $changed = true$
 - 5: $camUsed = \{\emptyset_1, \dots, \emptyset_{N_{cam}}\}$
 - 6: $\mathbb{P}^{Prev} = \mathbb{P}^0$
 - 7: $E_{ref} = \sum_{(i,j) \in \mathbb{P}^0} E_{BPV}^{i,j}$
 - 8: $E_{new} = E_{ref}$
 - 9: **while** $E_{new} > 0.9 \cdot E_{ref}$ **and** $changed$ **do**
 - 10: $changed = false$
 - 11: $\mu_{prev} = \mu^{\mathbb{P}^0}, \sigma_{prev} = \sigma^{\mathbb{P}^0}$
 - 12: **for** $i \in \mathbb{C}$ **do**
 - 13: $\mathbb{P}^{Cur} = \mathbb{P}^{Prev}$
 - 14: remove (i, j_{old}) from \mathbb{P}^{Cur}
 - 15: $j_{new} = \arg \min_j \{E_{BPV}^{i,j} \text{ s.t. } (i, j) \notin camUsed_i\}$
 - 16: $\mathbb{P}^{Cur} \leftarrow (i, j_{new})$
 - 17: $\mu_{cur} = \mu^{Cur}$
 - 18: $\sigma_{cur} = \sigma^{Cur}$
 - 19: **if** $\mu_{cur} > \mu_{prev}$ **and** $\sigma_{cur} < \sigma_{prev}$ **then**
 - 20: $\mathbb{P}^{Cur} = \mathbb{P}^{Cur}$
 - 21: $camUsed_i = camUsed_i \cup j_{new} 1;$
 - 22: $E_{new} = E_{new} + E_{BPV}^{i,j_{new}} - E_{BPV}^{i,j_{old}}$
 - 23: $changed = true$
 - 24: $\mu_{prev} = \mu_{cur}$
 - 25: $\sigma_{prev} = \sigma_{cur}$
 - 26: **else**
 - 27: $\mathbb{P}^{Cur} = \mathbb{P}^{Prev}$
 - 28: **end if**
 - 29: **end for**
 - 30: **end while**
 - 31: **return** \mathbb{P}^{Cur}
-

Then, the global visibility of f is $V_f = \sum_{i,j} v_f^{i,j}$ and the coverage of the whole mesh is represented by $C_{\mathbb{P}} = \{V_f, \forall f \in F\}$.

The second step aims at replacing camera pair (i, j) with a reasonable pair (i, k) that, even at the cost of a small decrease of energy E_{BPV} it improves the model coverage.

To do so, given the initial set \mathbb{P}^0 , we compute $\mu^{\mathbb{P}^0} = \mathbb{E}(C_{\mathbb{P}^0})$ and $\sigma^{\mathbb{P}^0} = \text{stddev}(C_{\mathbb{P}^0})$. For each camera i we compare the current camera pair $(i, j_1) \in \mathbb{P}^0$ with the second best camera pair (i, j_2) among the pairs (i, \cdot) . If (i, j_2) increases the mean coverage $\mu^{\mathbb{P}^0}$ and decreases $\sigma^{\mathbb{P}^0}$ we try to switch (i, j_1) and (i, j_2) so to obtain a new set \mathbb{P}^1 . If the sum of the energies $\sum_{(i,j) \in \mathbb{P}^1} E_{BPV}^{i,j} < 0.9 \cdot \sum_{(i,j) \in \mathbb{P}^0} E_{BPV}^{i,j}$ then the pair change is successful, i.e., $\mathbb{P}^1 = \mathbb{P}^1$, otherwise $\mathbb{P}^1 = \mathbb{P}^0$. We iterate this process until no further change happens (Algorithm 1).

Let notice that in the previous process we considered only one camera j for each reference camera i . This does not represent a limitation of the algorithm but a choice to have a fair

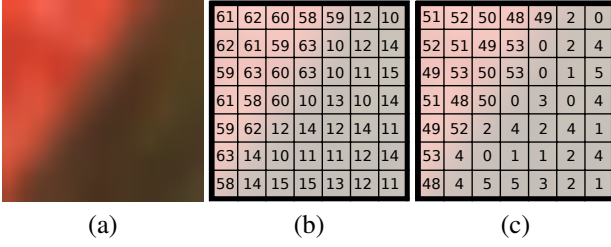


Fig. 4: Patch (a), Depth Patch (b) and Difference of depths patch (c)

comparison with the baselines of Tola *et al.* [12] and our implementation of Vu *et al.* [4] that compare one camera for each reference too. We refer to the experimental section for a more detailed discussion.

4. Occlusion Aware Masking

For each pixel $p(x, y)$, the mesh refinement presented in Section 2 aggregates the gradients of the similarity measure from the neighboring pixels in a squared patch P (in our case with size 5×5 px) (Fig. 4(a)).

In most cases, all the pixels in P contain information useful to refine the position of the 3D point corresponding to $p(x, y)$ by gradient descent. However, in the case of occlusions, the squared patch P contains information from regions of the scene not related to $p(x, y)$, which translates into errors in the gradient computation. To tackle this issue, we rely on the depth map of the current 3D model of the scene. The idea is to find which pixels of P have a depth similar and coherent to the pixel $p(x, y)$, and to use only them during the similarity gradient computation. In the following these pixels are named *valid pixels*.

For each pixel $p(h, k) \in P$, we compute the depth values $\delta(h, k)$ as the camera to model distance (4(b)) and its difference with respect to the depth of pixel $p(x, y)$ as $dd(h, k) = \delta(h, k) - \delta(x, y)$ (4(c)).

Since abrupt depth discontinuities would induce very high variances on the whole patch P , in principle, the standard deviation of the $dd(h, k)$ values are not enough informative to classify the *valid pixels*.

For this reason we propose the following procedure- First, we cluster the pixels between those with a depth closer and those with a depth are farther to the depth of $p(x, y)$. In other words, we cluster all the values $dd(h, k)$ in two sets DD_{min} and DD_{max} , which contain respectively the $dd(h, k)$ values closest to 0 and closest to $\max\{dd(h, k)\}$.

Then, we approximate a robust estimator $\hat{\sigma}$ of the standard deviation of the depths of the valid pixels. Assuming the pixels in DD_{min} to be likely *valid*, and the highest depth variances usually corresponding to the outer part of the patch, we compute $\hat{\sigma}$ as the value $dd(h, k) \in DD_{min}$ of the (spatially) farthest pixel w.r.t. the patch center (x, y) . A pixel $p(h, k)$ is valid when $|dd(h, k) - \hat{\sigma}| < 10|dd(h, k) - \max\{dd(h, k)\}|$

In Fig. 5 we illustrate the first iteration of mesh refinement for DTU sequence 63: the intensity of red represents the number of pixels considered in the gradient computation. It is possible to notice that the number of pixels considered decreases approaching to the mesh discontinuities.

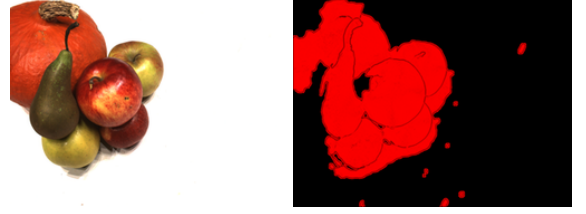
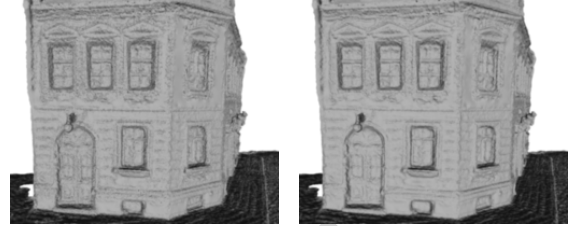


Fig. 5: Number of pixels considered for mesh refinement after occlusion masking



(a) without mask refinement (b) with mask refinement

Fig. 6: The effect of the masking refinement

5. Experiments

We tested our refinement method against 12 sequences of the DTU dataset [46], the fountain-P11 [1] and the Southbuilding [47] sequences. In Table 1 we reported the quantitative results on the DTU dataset against state of the art MVS methods. As mentioned in the Section 3 to have a fair comparison against Tola *et al.* and Vu *et al.* we choose one image for each reference camera. Indeed, both in Tola *et al.* and in our implementation of Vu *et al.* for each reference camera the second camera is chosen by relying on the knowledge of the visibility of the Structure from Motion 3D points. For the implementation of Vu *et al.* we choose the camera with the highest number of common points with a parallax between 20° and 60° , while Tola *et al.* with a parallax between 10° and 30° and at least 0.8 scale difference between the corresponding DAISY descriptors. In Figure 8 we tested our method comparing 1, 2, 4 or 8 cameras with respect to the reference image and the improvement obtained is almost negligible.

According to the procedure described in [46] we compared the distance (in mm) from refined 3D mesh to the ground truth point cloud to compute the accuracy of the model and vice-versa to evaluate its completeness. In the table we list the mean and median values of such distances. In these experiments we used as initial mesh those extracted by the method of Tola *et al.* + Poisson Reconstruction. Our method improves accuracy and completeness with respect to the baseline refinement [4]. In Table 2 we show the ablation study on the DTU dataset and in Fig. 9 we reported some examples of the outcome of our mesh refinement. From the actual proposal (last column of Table 2) we tested the refinement by using just subsets of the energy terms (P = Parallax, S = Symmetry, O = Overlap), with or without the coverage algorithm (represented by C) and with or without the masked refinement (represented by M). We do not mention the resolution prior since, in our scenario, it does not affect the result. This outcome is expected and coherent

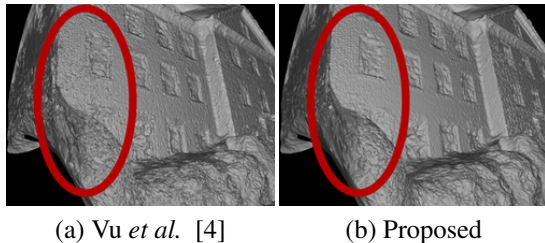


Fig. 7: The effect of the coverage algorithm

Table 1: Accuracy (model-to-GT distance) and Completeness (GT-to-model distance) of the refined model w.r.t. the state-of-the-art on the DTU dataset.

			[44]	[45]	[12]	[4]	Proposed
4	Acc.	Mean	0.7265	0.7947	0.3202	0.3254	0.3171
		Median	0.4014	0.3016	0.1888	0.2028	0.1972
4	Compl.	Mean	0.6054	0.6014	0.7791	0.7988	0.7936
		Median	0.4257	0.3644	0.3245	0.3463	0.3370
6	Acc.	Mean	1.0727	0.6660	0.3681	0.3666	0.3433
		Median	0.4568	0.2628	0.1993	0.2095	0.1990
6	Compl.	Mean	0.4281	0.4739	0.5208	0.5221	0.5095
		Median	0.3454	0.3568	0.3335	0.3427	0.3273
15	Acc.	Mean	0.9829	0.8544	0.4771	0.4614	0.4436
		Median	0.5588	0.4379	0.2866	0.2754	0.2689
15	Compl.	Mean	0.3243	0.5028	0.6906	0.6709	0.6737
		Median	0.2606	0.3876	0.4203	0.3973	0.3981
18	Acc.	Mean	1.3916	0.9665	0.5125	0.4918	0.4696
		Median	0.6793	0.4247	0.2603	0.2592	0.2547
18	Compl.	Mean	0.3640	0.4922	0.8996	0.8763	0.8742
		Median	0.2899	0.3715	0.3937	0.3740	0.3720
24	Acc.	Mean	3.4509	0.7518	0.3941	0.3871	0.3802
		Median	0.5255	0.3141	0.2659	0.2651	0.2632
24	Compl.	Mean	0.4010	0.4827	0.8512	0.8293	0.8300
		Median	0.2954	0.3691	0.4339	0.4119	0.4124
36	Acc.	Mean	0.5972	0.6270	0.3125	0.2859	0.2801
		Median	0.2317	0.2778	0.2007	0.1831	0.1831
36	Compl.	Mean	0.4622	0.6101	1.0331	1.0093	1.0070
		Median	0.2317	0.2930	0.2856	0.2527	0.2533
63	Acc.	Mean	2.4241	2.3992	0.9082	0.8461	0.7836
		Median	0.2782	1.1192	0.2711	0.2495	0.2303
63	Compl.	Mean	0.4730	0.6401	0.7189	0.7159	0.7158
		Median	0.2782	0.3849	0.2985	0.2916	0.2924
106	Acc.	Mean	0.5918	0.7881	0.3028	0.2844	0.2765
		Median	0.2793	0.3028	0.1902	0.1846	0.1821
106	Compl.	Mean	0.6902	0.7004	0.9950	0.9936	0.9935
		Median	0.2793	0.3244	0.3256	0.3220	0.3226
110	Acc.	Mean	3.4509	1.0922	0.7378	0.6867	0.6674
		Median	0.5255	0.3802	0.2314	0.2237	0.2222
110	Compl.	Mean	0.4010	0.5547	0.5675	0.5872	0.5892
		Median	0.2954	0.4041	0.4134	0.4238	0.4272
114	Acc.	Mean	0.6104	0.5789	0.2734	0.2696	0.2714
		Median	0.29700	0.2616	0.1835	0.1789	0.1802
114	Compl.	Mean	0.3544	0.4001	0.3895	0.3872	0.3878
		Median	0.2948	0.3191	0.2999	0.2964	0.2970
118	Acc.	Mean	5.5335	0.65.25	0.3093	0.2982	0.2911
		Median	4.0794	0.2922	0.1946	0.1913	0.1897
118	Compl.	Mean	0.3682	0.5489	0.7389	0.7416	0.7399
		Median	0.2856	0.3296	0.3301	0.3302	0.3290
122	Acc.	Mean	0.6367	0.6621	0.3049	0.2920	0.2884
		Median	0.3276	0.2967	0.1978	0.1920	0.1912
122	Compl.	Mean	0.3682	0.4576	0.6435	0.6441	0.6400
		Median	0.2856	0.3281	0.3278	0.3256	0.3226

because the cameras have approximately the same distant from the model and the same image resolution. To have a better understanding of the performance of each version, we reported the average rank and the average error value. To compute the former, we ranked the algorithms for each row, *i.e.*, for each error measure and each sequence, and we reported the average of these rankings. The latter was also used in [48] and it averages all the accuracy and completeness mean and median values for each column. In both cases, it is possible to notice that

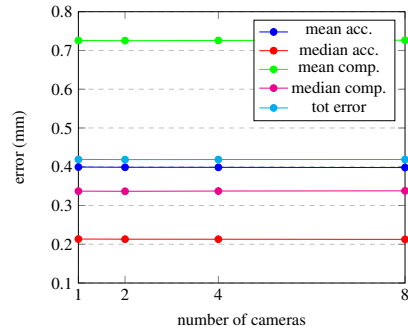


Fig. 8: Errors for different number of camera compared to the reference image

all the terms contribute to achieving a better outcome which is a good trade-off between accuracy and completeness. Fig. 6 shows how the masking refinement produces smoother results, by limiting the effect of occlusions.

Table 3 shows a further quantitative comparison evaluating the accuracy for the fountain-P11 dataset. We computed the Mean Absolute Error (MAE) and the Root Mean Squared Error (RMSE) as proposed in [1], *i.e.*, by comparing the depth maps rendered from the reconstructed model and the GT mesh. Our refinement method improves the accuracy of the classical refinement method in [4].

The Southbuilding dataset does not provide a reference ground truth. However from Fig. 9 and Fig. 7 (and the images reported in the supplementary material with higher resolution) it is possible to notice that our refinement method produces an output with more details than [4], thanks to the effective choice of camera pairs, together with the masked refinement. In particular, Fig. 7 shows that the coverage procedure avoids refinement to focus only on one part of the model producing unbalanced mesh evolution (Fig. 7(a)).

It is worth noticing that the method we use to define the pairs of cameras and to mask the occlusions during the photometric refinement can be easily applied to any mesh refinement method without any relevant modification in the code.

6. Conclusions and future works

In this paper, we addressed two relevant issues of mesh refinement: the choice of camera pairs used to compute similarity gradients and the occlusion management. We defined a model-based energy function to evaluate the overlap, the parallax the resolution and the symmetry among the camera pairs and we proposed a procedure to choose the set of pairs which provides a good trade-off between the defined energy and the model coverage. We also proposed a novel strategy to mask the region of the patch adopted to compute the similarity measures and the gradients such that the influences of model occlusions and discontinuities are neglected or, at least, limited. As future work, we plan to propose a parallel version of the presented refinement method that splits the mesh into several parts that can be processed independently. This allows exploring the level of parallelism depending on both the required mesh resolution and the computing platform. Moreover, we can optimize the

Table 2: Ablation study for the 12 sequences of the DTU dataset. Quantitative comparison according to the overall Average Ranking of the 3D model errors on the 12 sequences and the overall average of accuracy and completeness errors (Average Error Value)

	PC	SC	OC	SPC	OPC	OSC	OSP	OSPC	Proposed
Average Position	6.2	5.2	4.9	4.9	6.1	4.8	5.8	3.3	3.3
Average Error Value	0.4295	0.4258	0.4285	0.4276	0.4258	0.4281	0.4280	0.4223	0.4212

Table 3: Results on the fountain-P11 dataset (depth errors in m).

	Initial Mesh	[4]	Proposed
MAE	0.001513	0.001360	0.001199
RMSE	0.03890	0.03688	0.03462

energy-performance trade-off by leveraging the dynamic reconfiguration support offered in multi-core architectures [49, 50]. We also intend to improve the camera selection exploiting the information recovered by the dense Multi-View Stereo method such as [17], which jointly estimate image depths and pixel-wise camera pairs.

Acknowledgements

This work was partially funded by EIT digital

References

- [1] C. Strecha, W. von Hansen, L. Van Gool, P. Fua, U. Thoennessen, On benchmarking camera calibration and multi-view stereo for high resolution imagery, in: *Computer Vision and Pattern Recognition*, IEEE, 2008, pp. 1–8.
- [2] S. M. Seitz, B. Curless, J. Diebel, D. Scharstein, R. Szeliski, A comparison and evaluation of multi-view stereo reconstruction algorithms, in: *Computer vision and pattern recognition*, Vol. 1, IEEE, 2006, pp. 519–528.
- [3] H. Aanæs, R. R. Jensen, G. Vogiatzis, E. Tola, A. B. Dahl, Large-scale data for multiple-view stereopsis, *International Journal of Computer Vision* (2016) 1–16.
- [4] H. H. Vu, P. Labatut, J.-P. Pons, R. Keriven, High accuracy and visibility-consistent dense multiview stereo, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 34 (5) (2012) 889–901.
- [5] P. Moulon, P. Monasse, R. Marlet, Others, *Openmvg. an open multiple view geometry library.*, <https://github.com/openMVG/openMVG>.
- [6] C. Wu, S. Agarwal, B. Curless, S. M. Seitz, Multicore bundle adjustment, in: *Computer Vision and Pattern Recognition*, IEEE, 2011, pp. 3057–3064.
- [7] R. T. Collins, A space-sweep approach to true multi-image matching, in: *Computer Vision and Pattern Recognition*, IEEE, 1996, pp. 358–363.
- [8] H. H. Vu, R. Keriven, P. Labatut, J.-P. Pons, Towards high-resolution large-scale multi-view stereo, in: *Computer Vision and Pattern Recognition*, IEEE, 2009, pp. 1430–1437.
- [9] J. Li, E. Li, Y. Chen, L. Xu, Y. Zhang, Bundled depth-map merging for multi-view stereo, in: *2010 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2010, pp. 2769–2776.
- [10] T. Ebner, O. Schreer, I. Feldmann, Fully automated highly accurate 3d reconstruction from multiple views, in: *International Conference on Image Processing (ICIP)*, IEEE, 2017, pp. 2528–2532.
- [11] M. Goesele, N. Snavely, B. Curless, H. Hoppe, S. M. Seitz, Multi-view stereo for community photo collections, in: *Computer Vision*, 2007. ICCV 2007. IEEE 11th International Conference on, IEEE, 2007, pp. 1–8.
- [12] E. Tola, C. Strecha, P. Fua, Efficient large-scale multi-view stereo for ultra high-resolution image sets, *Machine Vision and Applications* 23 (5) (2012) 903–920.
- [13] S. Shen, Depth-map merging for multi-view stereo with high resolution images, in: *Pattern Recognition (ICPR)*, 2012 21st International Conference on, IEEE, 2012, pp. 788–791.
- [14] L. Lou, Y. Liu, J. Han, J. H. Doonan, Accurate multi-view stereo 3d reconstruction for cost-effective plant phenotyping, in: *International Conference Image Analysis and Recognition*, Springer, 2014, pp. 349–356.
- [15] G. Vogiatzis, P. H. Torr, R. Cipolla, Multi-view stereo via volumetric graph-cuts, in: *Computer Vision and Pattern Recognition*, 2005. CVPR 2005. IEEE Computer Society Conference on, Vol. 2, IEEE, 2005, pp. 391–398.
- [16] E. Zheng, E. Dunn, V. Jovic, J.-M. Frahm, Patchmatch based joint view selection and depthmap estimation, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1510–1517.
- [17] J. L. Schönberger, E. Zheng, J.-M. Frahm, M. Pollefeys, Pixelwise view selection for unstructured multi-view stereo, in: *European Conference on Computer Vision*, Springer, 2016, pp. 501–518.
- [18] A. Romanoni, M. Matteucci, Tapa-mvs: Textureless-aware patchmatch multi-view stereo, *arXiv preprint arXiv:1903.10929*.
- [19] A. Delaunoy, E. Prados, P. G. I. Piracés, J.-P. Pons, P. Sturm, Minimizing the multi-view stereo reprojection error for triangular surface meshes, in: *British Machine Vision Conference, BMVA*, 2008, pp. 1–10.
- [20] S. Li, S. Y. Siu, T. Fang, L. Quan, Efficient multi-view surface refinement with adaptive resolution control, in: *European Conference on Computer Vision*, Springer, 2016, pp. 349–364.
- [21] G. Hermsillo, C. Chefd’Hotel, O. Faugeras, Variational methods for multimodal image matching, *International Journal of Computer Vision* 50 (3) (2002) 329–343.
- [22] O. Faugeras, R. Keriven, Variational principles, surface evolution, pdes, level set methods, and the stereo problem, *IEEE Transactions on Image Processing* 7 (3) (1998) 336–344. doi:10.1109/83.661183.
- [23] H. Jin, A. J. Yezzi, S. Soatto, Variational multiframe stereo in the presence of specular reflections, in: *First International Symposium on 3D Data Processing Visualization and Transmission*, IEEE, 2002, pp. 626–630.
- [24] A. Yezzi, S. Soatto, Stereoscopic segmentation, *International Journal of Computer Vision* 53 (1) (2003) 31–43.
- [25] S. Fuhrmann, M. Goesele, Floating scale surface reconstruction, *ACM Transactions on Graphics (TOG)* 33 (4) (2014) 46.
- [26] J. E. Solem, N. C. Overgaard, A geometric formulation of gradient descent for variational problems with moving surfaces, in: *Scale Space and PDE methods in computer vision*, Springer, 2005, pp. 419–430.
- [27] K.-J. Yoon, E. Prados, P. Sturm, Joint estimation of shape and reflectance using multiple images with known illumination conditions, *International Journal of Computer Vision* 86 (2-3) (2010) 192–210.
- [28] J.-P. Pons, R. Keriven, O. Faugeras, Multi-view stereo reconstruction and scene flow estimation with a global image-based matching score, *International Journal of Computer Vision* 72 (2) (2007) 179–193.
- [29] J.-P. Pons, R. Keriven, O. Faugeras, Modelling dynamic scenes by registering multi-view image sequences, in: *Computer Vision and Pattern Recognition*, Vol. 2, IEEE, 2005, pp. 822–827.
- [30] A. Zaharescu, E. Boyer, R. Horaud, Transformesh: a topology-adaptive mesh-based approach to surface evolution, in: *Asian Conference on Computer Vision*, Springer, 2007, pp. 166–175.
- [31] P. Gargallo, E. Prados, P. Sturm, Minimizing the reprojection error in surface reconstruction from images, in: *International Conference on Computer Vision*, IEEE, 2007, pp. 1–8.
- [32] A. Delaunoy, E. Prados, Gradient flows for optimizing triangular mesh-based surfaces: Applications to 3d reconstruction problems dealing with visibility, *International journal of computer vision* 95 (2) (2011) 100–123.
- [33] H. H. Vu, Large-scale and high-quality multi-view stereo, Ph.D. thesis, Paris Est (2011).
- [34] A. Delaunoy, M. Pollefeys, Photometric bundle adjustment for dense multi-view 3d modeling, in: *Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2014, pp. 1486–1493.
- [35] Z. Li, K. Wang, W. Zuo, D. Meng, L. Zhang, Detail-preserving and content-aware variational multi-view stereo reconstruction, *IEEE Transactions on Image Processing* 25 (2) (2016) 864–877.

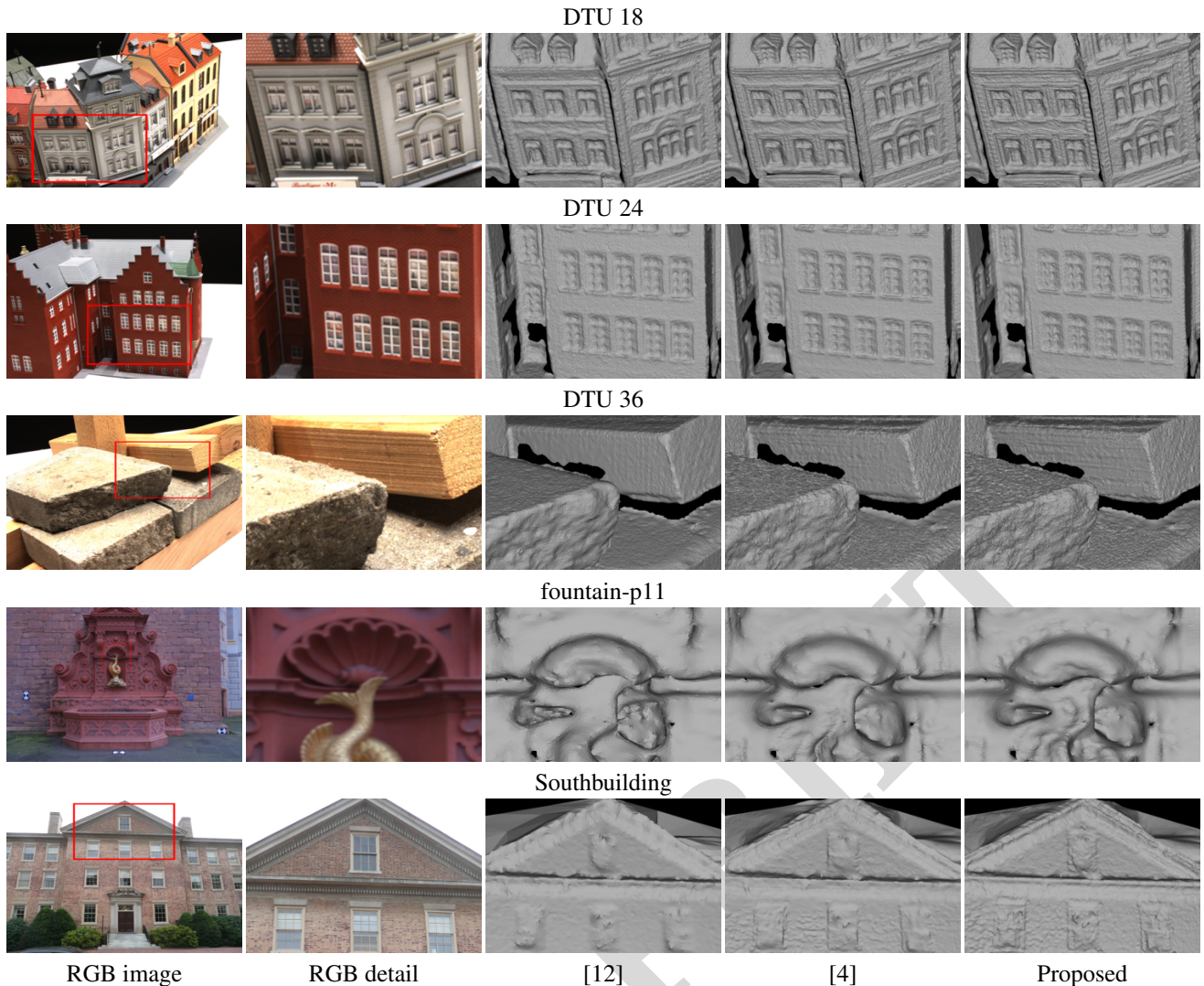


Fig. 9: Examples of refinement results with respect to the initial 3D mesh and the baseline refinement method

- [36] A. Romanoni, D. Fiorenti, M. Matteucci, Mesh-based 3d textured urban mapping, in: 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), IEEE, 2017, pp. 3460–3466.
- [37] M. Blaha, M. Rothermel, M. R. Oswald, T. Sattler, A. Richard, J. D. Wegner, M. Pollefeys, K. Schindler, Semantically informed multiview surface refinement, *International Journal of Computer Vision*.
- [38] A. Romanoni, M. Ciccone, F. Visin, M. Matteucci, Multi-view stereo with single-view semantic mesh refinement, in: *IEEE International Conference on Computer Vision Workshops*, 2017, pp. 706–715.
- [39] H. H. Vu, *Streo multi-vues a grande chelleet de haute qualit*, Ph.D. thesis, Ecole des ponts Paristech (Dec 2011).
- [40] M. Wardetzky, S. Mathur, F. Kälberer, E. Grinspun, Discrete laplace operators: no free lunch, in: *Symposium on Geometry processing*, 2007, pp. 33–37.
- [41] L. Morreale, A. Romanoni, M. Matteucci, Predicting the next best view for 3d mesh refinement, in: *International Conference on Intelligent Autonomous Systems*, Springer, 2018, pp. 760–772.
- [42] E. Nocerino, F. Menna, F. Remondino, Accuracy of typical photogrammetric networks in cultural heritage 3d modeling projects., *International Archives of the Photogrammetry, Remote Sensing & Spatial Information Sciences* 45.
- [43] R. Wackrow, J. H. Chandler, Minimising systematic error surfaces in digital elevation models using oblique convergent imagery, *The Photogrammetric Record* 26 (133) (2011) 16–31.
- [44] N. Campbell, G. Vogiatzis, C. Hernández, R. Cipolla, Using multiple hypotheses to improve depth-maps for multi-view stereo, *European Conference on Computer Vision* (2008) 766–779.
- [45] Y. Furukawa, J. Ponce, Accurate, dense, and robust multiview stereopsis, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 32 (8) (2010) 1362–1376.
- [46] R. Jensen, A. Dahl, G. Vogiatzis, E. Tola, H. Aanæs, Large scale multiview stereopsis evaluation, in: *2014 IEEE Conference on Computer Vision and Pattern Recognition*, IEEE, 2014, pp. 406–413.
- [47] C. Häne, C. Zach, A. Cohen, R. Angst, M. Pollefeys, Joint 3d scene reconstruction and class segmentation, in: *Computer Vision and Pattern Recognition (CVPR)*, 2013 IEEE Conference on, IEEE, 2013, pp. 97–104.
- [48] A. Poms, C. Wu, S.-I. Yu, Y. Sheikh, Learning patch reconstructability for accelerating multi-view stereo, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3041–3050.
- [49] D. Zoni, A. Canidio, W. Fornaciari, P. Englezakis, C. Nicopoulos, Y. Sazeides, Blackout: Enabling fine-grained power gating of buffers in network-on-chip routers, *Journal of Parallel and Distributed Computing* 104 (2017) 130 – 145. doi:<https://doi.org/10.1016/j.jpdc.2017.01.016>.
- [50] D. Zoni, L. Colombo, W. Fornaciari, Darkcache: Energy-performance optimization of tiled multi-cores by adaptively power-gating llc banks, *ACM Trans. Archit. Code Optim.* 15 (2) (2018) 21:1–21:26. doi:10.1145/3186895.