

DAISY: An Efficient Dense Descriptor Applied to Wide-Baseline Stereo

Engin Tola, Vincent Lepetit, and Pascal Fua, *Senior Member, IEEE*

Abstract—In this paper, we introduce a local image descriptor, DAISY, which is very efficient to compute densely. We also present an EM-based algorithm to compute dense depth and occlusion maps from wide-baseline image pairs using this descriptor. This yields much better results in wide-baseline situations than the pixel and correlation-based algorithms that are commonly used in narrow-baseline stereo. Also, using a descriptor makes our algorithm robust against many photometric and geometric transformations. Our descriptor is inspired from earlier ones such as SIFT and GLOH but can be computed much faster for our purposes. Unlike SURF, which can also be computed efficiently at every pixel, it does not introduce artifacts that degrade the matching performance when used densely. It is important to note that our approach is the first algorithm that attempts to estimate dense depth maps from *wide-baseline image pairs*, and we show that it is a good one at that with many experiments for depth estimation accuracy, occlusion detection, and comparing it against other descriptors on laser-scanned ground truth scenes. We also tested our approach on a variety of indoor and outdoor scenes with different photometric and geometric transformations and our experiments support our claim to being robust against these.

Index Terms—Image processing and computer vision, dense depth map estimation, local descriptors.

1 INTRODUCTION

THOUGH dense short-baseline stereo matching is well understood [9], [25], its wide-baseline counterpart is, in contrast, much more challenging due to large perspective distortions and increased occluded areas. It is nevertheless worth addressing because it can yield more accurate depth estimates while requiring fewer images to reconstruct a complete scene. Also, it may be necessary to compute depth from two widely separated cameras such as a surveillance application in which installing cameras side-by-side is not feasible.

Large correlation windows are not appropriate for wide-baseline matching because they are not robust to perspective distortions and tend to straddle areas of different depths or partial occlusions. Thus, most researchers favor simple pixel differencing [6], [16], [24] or correlation over very small windows [26]. They then rely on optimization techniques such as graph-cuts [16] or PDE-based diffusion operators [27] to enforce spatial consistency. The drawback of using small image patches is that reliable image information can only be obtained where the image texture is of sufficient quality. Furthermore, the matching becomes very sensitive to illumination changes and repetitive patterns.

An alternative to performing dense wide-baseline matching is to first match a few feature points, triangulate them, and then locally rectify the images. This approach, however,

potentially is not without problems. If some matches are wrong and are not detected as such, gross reconstruction errors will occur. Furthermore, image rectification in the triangles may not be sufficient if the scene within cannot be treated as locally planar.

We instead advocate replacing correlation windows with local region descriptors, which lets us take advantage of powerful global optimization schemes such as graph-cuts to force spatial consistency. Existing local region descriptors such as SIFT [19] or GLOH [21] have been designed for robustness to perspective and lighting changes and have proven successful for sparse wide-baseline matching. However, they are much more computationally demanding than simple correlation. Thus, for dense wide-baseline matching purposes, local region descriptors have so far only been used to match a few seed points [33] or to provide constraints on the reconstruction [27].

We therefore introduce a new descriptor that retains the robustness of SIFT and GLOH and can be computed quickly at every single image pixel. Its shape is closely related to that of [32], which has been shown to be optimal for sparse matching but is not designed for efficiency. We use our descriptor for dense matching and view-based synthesis using stereo pairs having various image transforms or for pairs with too large a baseline for standard correlation-based techniques to work, as shown in Figs. 1, 2, 3, and 4. For example, on a standard laptop, it takes less than 4 seconds to perform the computations using our descriptor over all the pixels of an 800×600 image, whereas it takes over 250 seconds using SIFT. Furthermore, it gives better results than SIFT, SURF, NCC, and pixel differencing, as will be shown by comparing the resulting depth maps to laser-scanned data.

To be specific, SIFT and GLOH owe much of their strength to the use of gradient orientation histograms, which are

• The authors are with the Computer Vision Laboratory, Ecole Polytechnique Fédérale de Lausanne (EPFL), EPFL/IC/ISIM/CVLab, Station 14, CH-1015 Lausanne, Switzerland.

E-mail: {engin.tola, vincent.lepetit, pascal.fua}@epfl.ch.

Manuscript received 11 Aug. 2008; revised 27 Jan. 2009; accepted 27 Mar. 2009; published online 31 Mar. 2009.

Recommended for acceptance by S.B. Kang.

For information on obtaining reprints of this article, please send e-mail to: tpami@computer.org, and reference IEEECS Log Number TPAMI-2008-08-0482.

Digital Object Identifier no. 10.1109/TPAMI.2009.77.



Fig. 1. Contrast Change: The first two images are used as input. We manually increased the contrast of the first image and tried to estimate the depth from the second image's point of view. We used NCC, SIFT, and DAISY, and their reconstructions are displayed in the second row, respectively. We also resynthesized the second image using the depth map of DAISY and the first image's intensities and show it at the end of the first row.

relatively robust to distortions. The more recent **SURF** descriptor [4] approximates them by using integral images to compute the histograms bins. This method is computationally effective with respect to computing the descriptor's value at every pixel, but does away with SIFT's spatial weighting scheme. All gradients contribute equally to their respective bins, which results in damaging artifacts when used for dense computation. The key insight of this paper is that computational efficiency can be achieved without performance loss by convolving orientation maps to compute the bin values using Gaussian kernels. This lets us match relatively large patches— 31×31 —at an acceptable computational cost and improve robustness in unoccluded

areas over techniques that use smaller patches. Using large areas requires handling occlusion boundaries properly, though, and we address this issue by using different masks at each location and selecting the best one by using an **Expectation Maximization (EM) framework**. This is inspired by the earlier works of [13], [14], [15] where multiple or adaptive correlation windows are used.

After discussing related work in Section 2, we introduce our new local descriptor and present an efficient way to compute it in Section 3. In Section 4, we detail our EM-based occlusion handling framework. Finally, in Section 5, we present results and compare our descriptor to SIFT, SURF, NCC, and pixel differencing.

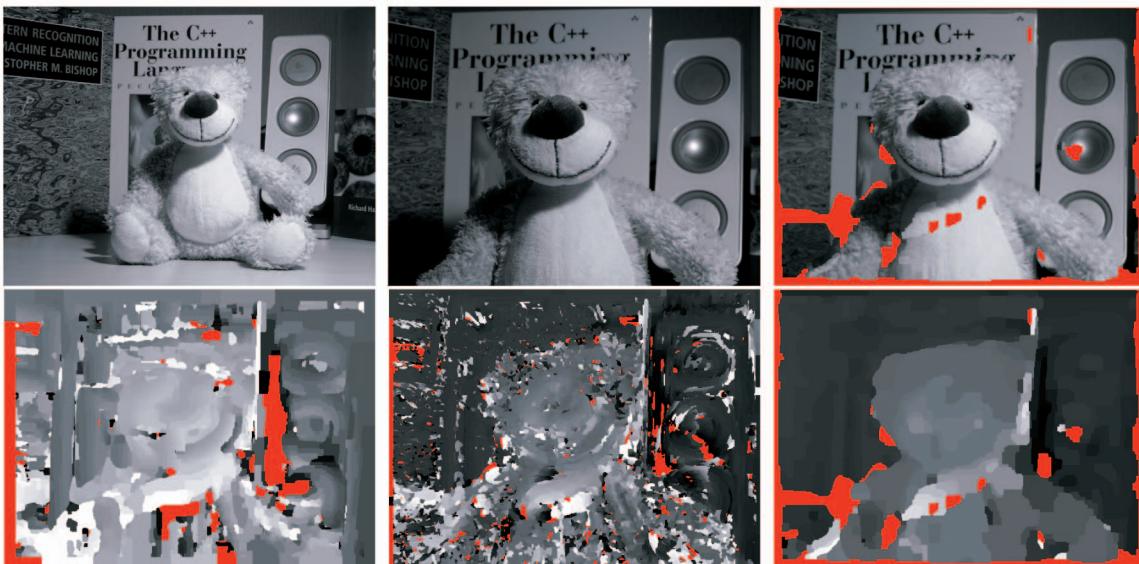


Fig. 2. Scale Change: We used the first two images of the upper row for computing the depth map from the second image's point of view. The depth maps are computed using NCC, SIFT, and DAISY, and they are displayed in the lower row in that order. The last image in the first row shows the resynthesized image using the DAISY's depth estimate. Although scale change is not explicitly addressed in any way and we used the same parameters for the descriptors of two images, we obtain a very acceptable depth map.

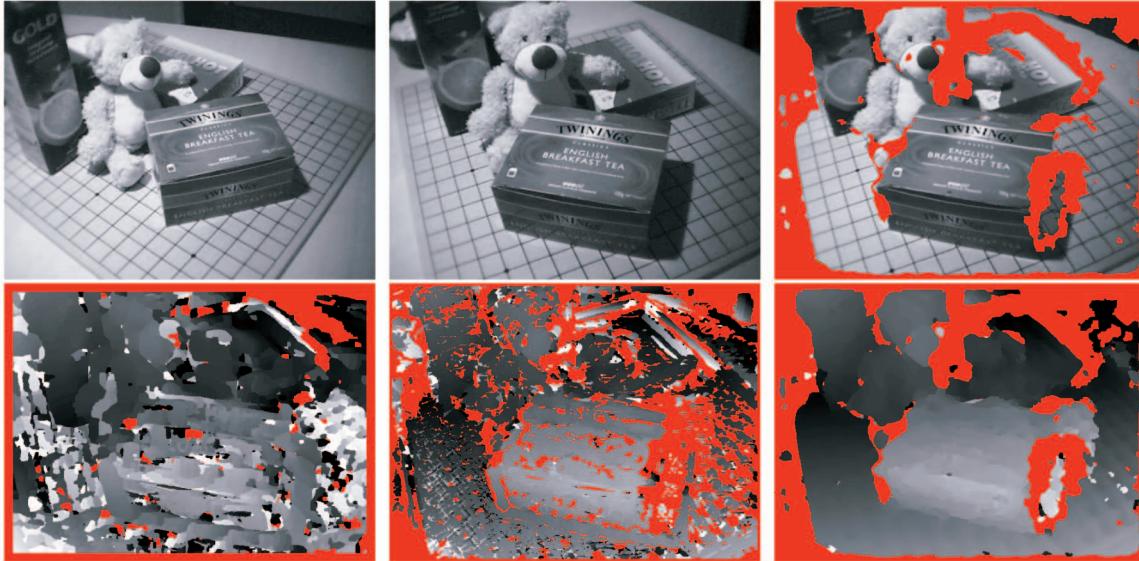


Fig. 3. Image Quality: We used the first two images of the upper row, which are obtained by a webcam, for computing the depth map from the second image's point of view. The depth maps are computed using NCC, SIFT, and DAISY, and they are displayed in the lower row in that order. The last image in the first row shows the resynthesized image using the DAISY's depth estimate. Despite the somewhat blurry, low-quality nature of the images, we can still compute a good depth map.

2 RELATED WORK

Even though multiview 3D surface reconstruction has been investigated for many decades [9], [25], it is still far from being completely solved due to many sources of errors, such as perspective distortion, occlusions, and textureless areas. Most state-of-the-art methods rely on first using local measures to estimate the similarity of pixels across images and then on imposing global shape constraints using dynamic programming [3], level sets [11], space carving [17], graph-cuts [16], [24], [8], PDE [1], [27], or EM [26]. In this paper, we do not focus on the method used to impose the global constraints and use a standard one [8]. Instead,

we concentrate on the similarity measure all of these algorithms rely on.

In a short-baseline setup, reconstructed surfaces are often assumed near frontoparallel, so the similarity between pixels can be measured by cross-correlating square windows. This is less prone to errors compared to pixel differencing and allows normalization against illumination changes.

In a wide-baseline setup, however, large correlation windows are especially affected by perspective distortions and occlusions. Thus, wide-baseline methods [1], [16], [26], [27] tend to rely on very small correlation windows or revert to pointwise similarity measures, which loose the discriminative power larger windows could provide. This

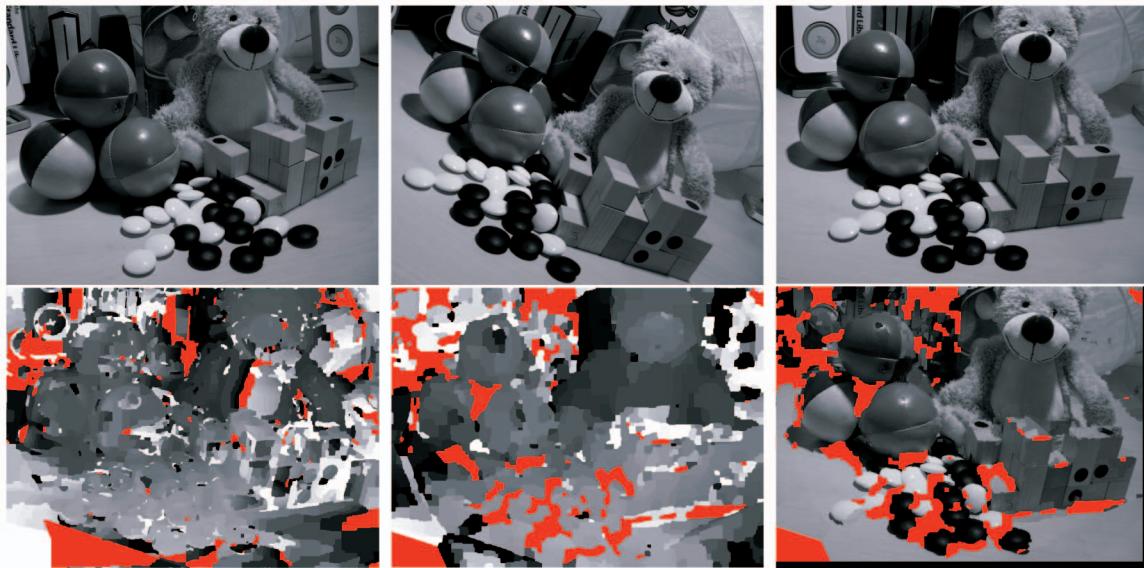


Fig. 4. Rotation around a point: We used the first two images of the upper row for computing the depth map from the third image's point of view. The depth maps are computed using NCC and DAISY, and they are displayed in the lower row in that order. The last image in the second row shows the resynthesized image using the DAISY's depth estimate.

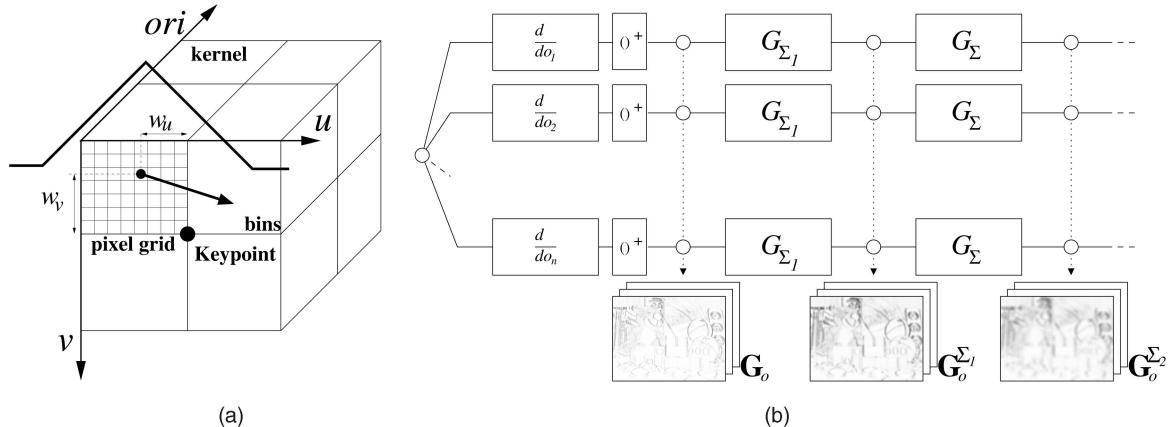


Fig. 5. Relationship between SIFT and DAISY: (a) SIFT is a 3D histogram computed over a local area where each pixel location contributes to bins depending on its location and the orientation of its image gradient, the importance of the contribution being proportional to the norm of the gradient. Each gradient vector is spread over $2 \times 2 \times 2 = 8$ bins to avoid boundary effects, and its contribution to each bin is weighted by the distances between the pixel location and the bin boundaries. (b) DAISY computes similar values but in a dense way. Each gradient vector also contributes to several of the elements of the description vector, but the sum of the weighted contributions is computed by convolution for better computation times. We first compute orientation maps from the original images, which are then convolved to obtain the convolved orientation maps $\mathbf{G}_o^{\Sigma_i}$. The values of the $\mathbf{G}_o^{\Sigma_i}$ correspond to the values in the SIFT bins, and will be used to build DAISY. By chaining the convolutions, the $\mathbf{G}_o^{\Sigma_i}$ can be obtained very efficiently.

loss can be compensated for by using multiple [2], [27] or high-resolution [27] images. The latter is particularly effective because areas that appear uniform at a small scale are often quite textured when imaged at a larger one. However, even then, lighting changes remain difficult to handle. For example, Strecha et al. [27] show results either for wide baseline without light changes, or with light changes but under a shorter baseline.

As we shall see, our feature descriptor reduces the need for higher resolution images and achieve comparable results using fewer number of images. It does so by considering large image patches while remaining stable under perspective distortions. Earlier approaches to this problem relied on warping the correlation windows [10]. However, the warps were estimated from a first reconstruction obtained using classical windows, which is usually not practical in wide-baseline situations. In contrast, our method does not require an initial reconstruction. Additionally, in a recent publication [32], a descriptor which is very similar to ours in shape has been shown to outperform many state-of-the-art feature descriptors for sparse point matching. However, unlike this descriptor, ours is designed for fast and efficient computation at every pixel in the image.

Local image descriptors have already been used in dense matching, though in a more traditional way, to match only sparse pixels that are feature points [31], [19]. In [27], [33], these matched points are used as anchors for computing the full reconstruction. Yao and Cham [33] propagate the disparities of the matched feature points to their neighbors, while Strecha et al. [27] use them to initialize an iterative estimation of the depth maps.

To summarize, local descriptors have already proved their worth for dense wide-baseline matching, but only in a limited way. This is due in part to their high computational cost and in part to their sensitivity to occlusions. The technique we propose addresses both issues.

3 OUR LOCAL DESCRIPTOR

In this section, we briefly describe SIFT [19] and GLOH [21] and then introduce our DAISY descriptor. We discuss both

its relationship with them and its greater effectiveness for dense computations.

The SIFT and GLOH descriptors involve 3D histograms in which two dimensions correspond to image spatial dimensions and the additional dimension to the image gradient direction. They are computed over local regions, usually centered on feature points but sometimes also densely sampled for object recognition tasks [12], [18]. Each pixel belonging to the local region contributes to the histogram depending on its location in the local region, and on the orientation and the norm of the image gradient at its location. As depicted by Fig. 5a, when an image gradient vector computed at a pixel location is integrated to the 3D histogram, its contribution is spread over $2 \times 2 \times 2 = 8$ bins to avoid boundary effects. More precisely, each bin is incremented by the value of the gradient norm multiplied by a weight inversely related to the distances (i.e., as the distance increases, weight decreases) between the pixel location and the bin boundaries, and also to the distance between the pixel location and the one of the key point. As a result, each bin contains a weighted sum of the norms of the image gradients around its center, where the weights roughly depend on the distance to the bin center.

In this work, our goal is to reformulate these descriptors so that they can be efficiently computed at every pixel location. Intuitively, this means computing the histograms only once per region and reusing them for all neighboring pixels.

To this end, we replace the weighted sums of gradient norms by convolutions of the gradients in specific directions with several Gaussian filters. We will see that this gives the same kind of invariance as the SIFT and GLOH histogram building, but is much faster for dense-matching purposes and allows the computation of the descriptors in all directions with little overhead.

Even though SIFT, GLOH, and DAISY involve different weighting schemes for the orientation gradients, the computed histograms can be expected to be very similar which gives a new insight on what makes SIFT work: Convolving with a kernel simultaneously dampens the noise and gives a measure of invariance to translation. This

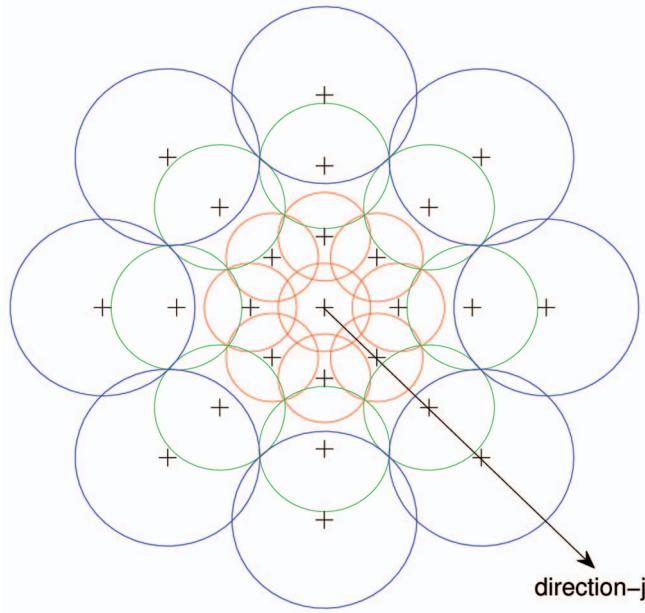


Fig. 6. The DAISY descriptor: Each circle represents a region where the radius is proportional to the standard deviations of the Gaussian kernels and the “+” sign represents the locations where we sample the convolved orientation maps center being a pixel location where we compute the descriptor. By overlapping the regions, we achieve smooth transitions between the regions and a degree of rotational robustness. The radii of the outer regions are increased to have an equal sampling of the rotational axis, which is necessary for robustness against rotation.

is also better than integral image-like computations of histograms [22] in which all gradient vectors have the same contribution. We can very efficiently reduce the influence of gradient norms from distant locations.

Fig. 6 depicts the resulting descriptor. Note that its shape resembles that of a descriptor [32] that has been shown to outperform many state-of-the-art ones. However, unlike that descriptor, DAISY is also designed for effective dense computation. The parameters that control its shape are listed in Table 1. We will discuss in Section 5 how they should be chosen.

There is a strong connection between DAISY and geometric blur [5]. In this work, the authors recommended using smaller blur kernels near the center and larger away from it and reported successful results using oriented edge filter responses. DAISY follows this recommendation by using larger Gaussian kernels in its outer rings but replaces the edge filters by simple convolutions for the sake of efficiency.

3.1 The DAISY Descriptor

We now give a more formal definition of our DAISY descriptor. For a given input image, we first compute H number of *orientation maps*, \mathbf{G}_i , $1 \leq i \leq H$, one for each quantized direction, where $\mathbf{G}_o(u, v)$ equals the image gradient norm at location (u, v) for direction o if it is bigger than zero, else it is equal to zero. This preserves the polarity of the intensity changes. Formally, orientation maps are written as $\mathbf{G}_o = (\frac{\partial \mathbf{I}}{\partial o})^+$, where \mathbf{I} is the input image, o is the orientation of the derivative, and $(.)^+$ is the operator such that $(a)^+ = \max(a, 0)$.

Each orientation map is then convolved several times with Gaussian kernels of different Σ values to obtain *convolved orientation maps* for different sized regions as $\mathbf{G}_o^\Sigma = G_\Sigma * (\frac{\partial \mathbf{I}}{\partial o})^+$ with G_Σ a Gaussian kernel. Different Σ s are used to control the size of the region.

Our primary motivation here is to reduce the computational requirements and convolutions can be implemented very efficiently especially when using Gaussian filters, which are separable. Moreover, we can compute the orientation maps for different sizes at low cost because convolutions with a large Gaussian kernel can be obtained from several consecutive convolutions with smaller kernels. More specifically, given $\mathbf{G}_o^{\Sigma_1}$, we can efficiently compute $\mathbf{G}_o^{\Sigma_2}$ with $\Sigma_2 > \Sigma_1$ as

$$\mathbf{G}_o^{\Sigma_2} = G_{\Sigma_2} * \left(\frac{\partial \mathbf{I}}{\partial o} \right)^+ = G_\Sigma * G_{\Sigma_1} * \left(\frac{\partial \mathbf{I}}{\partial o} \right)^+ = G_\Sigma * \mathbf{G}_o^{\Sigma_1},$$

with $\Sigma = \sqrt{\Sigma_2^2 - \Sigma_1^2}$. This computational flow, the incremental computation of the *convolved orientation maps* from an input image, is summarized in Fig. 5b.

To make the link with SIFT and GLOH, note that each pixel location of the convolved orientation maps contains a value very similar to the value of a bin in SIFT or GLOH that is a weighted sum of gradient norms computed over a small neighborhood. We use a Gaussian kernel whereas SIFT and GLOH rely on a triangular shaped kernel. It can also be linked to tensor voting in [20] by thinking of each location in our orientation maps as a voting component and of our aggregation kernel as the voting weights.

As depicted by Fig. 6, at each pixel location, DAISY consists of a vector made of values from the convolved orientation maps located on concentric circles centered on the location, and where the amount of Gaussian smoothing is proportional to the radii of the circles. As can be seen

TABLE 1
DAISY Parameters

Parameter Name	Symbol	Description and Default Value
Radius	R	Distance from the center pixel to the outer most grid point. (15)
Radius Quantization No.	Q	Number of convolved orientations layers with different Σ 's. (3)
Angular Quantization No.	T	Number of histograms at a single layer. (8)
Histogram Quantization No.	H	Number of bins in the histogram. (8)
Grid Point No.	S	Number of histograms used in the descriptor = $Q * T + 1$.
Descriptor Size	D_s	The total size of the descriptor vector = $S * H$

from the figure, this gives the descriptor the appearance of a flower, hence its name.

Let $\mathbf{h}_\Sigma(u, v)$ represent the vector made of the values at location (u, v) in the orientation maps after convolution by a Gaussian kernel of standard deviation Σ .

$$\mathbf{h}_\Sigma(u, v) = [\mathbf{G}_1^\Sigma(u, v), \dots, \mathbf{G}_H^\Sigma(u, v)]^\top, \quad (1)$$

where \mathbf{G}_1^Σ , \mathbf{G}_2^Σ , and \mathbf{G}_H^Σ denote the Σ -convolved orientation maps in different directions. We normalize these vectors to unit norm, and denote the normalized vectors by $\tilde{\mathbf{h}}_\Sigma(u, v)$. The normalization is performed in each histogram independently to be able to represent the pixels near occlusions as correct as possible. If we were to normalize the descriptor as a whole, then the descriptors of the same point that is close to an occlusion would be very different when imaged from different viewpoints.

Problems might arise in homogeneous regions since we normalize each histogram independently. However, consider that we are designing this descriptor for a stereo application. In a worst-case scenario, DAISY will not perform any worse than a standard region-based metric like NCC. However, this will not happen as often because we use a relatively large descriptor. Furthermore, the global optimization algorithm discussed in Section 4 will often fix the resulting errors. If one truly wants the robustness of large regions, one solution to this might be computing unnormalized descriptors and normalizing the visible parts of the descriptor globally before dissimilarity computation. This, however, will increase the computation time of the matching stage with two additional normalization operations for each possible depth per pixel. For applications other than stereo, the normalization should probably be changed depending on the specifics of the application, but this is beyond the scope of this paper.

If Q represents the number of different circular layers, then the full DAISY descriptor $\mathcal{D}(u_0, v_0)$ for location (u_0, v_0) is defined as the concatenation of \mathbf{h} vectors:

$$\begin{aligned} \mathcal{D}(u_0, v_0) = & \\ & [\tilde{\mathbf{h}}_{\Sigma_1}^\top(u_0, v_0), \\ & \tilde{\mathbf{h}}_{\Sigma_1}^\top(\mathbf{l}_1(u_0, v_0, R_1)), \dots, \tilde{\mathbf{h}}_{\Sigma_1}^\top(\mathbf{l}_T(u_0, v_0, R_1)), \\ & \tilde{\mathbf{h}}_{\Sigma_2}^\top(\mathbf{l}_1(u_0, v_0, R_2)), \dots, \tilde{\mathbf{h}}_{\Sigma_2}^\top(\mathbf{l}_T(u_0, v_0, R_2)), \\ & \dots \\ & \tilde{\mathbf{h}}_{\Sigma_Q}^\top(\mathbf{l}_1(u_0, v_0, R_Q)), \dots, \tilde{\mathbf{h}}_{\Sigma_Q}^\top(\mathbf{l}_T(u_0, v_0, R_Q))]^\top, \end{aligned}$$

where $\mathbf{l}_j(u, v, R)$ is the location with distance R from (u, v) in the direction given by j when the directions are quantized into the T values of Table 1.

We use a circular grid instead of SIFT's regular one since it has been shown to have better localization properties [21]. In that sense, our descriptor is closer to GLOH without PCA than to SIFT. Combining an isotropic Gaussian kernel with a circular grid also makes our descriptor naturally resistant to rotational perturbations. The overlapping regions ensure a smoothly changing descriptor along the rotation axis and, by increasing the overlap, we can make it more robust up to the point where the descriptor starts losing its discriminative power.

As mentioned earlier, one important advantage of the circular design and using isotropic kernels is that, when we

want to compute the descriptor in a different orientation, there is no need to recompute the convolved orientation maps; they are still valid and we can recompute the descriptor by simply rotating the sampling grid. The histograms will then also need to be shifted circularly to account for the change in relative gradient orientations but this only represents a very small overhead.

We mentioned earlier that the variance of the Gaussian kernels is chosen to be proportional to the size of the regions in the descriptor. Specifically, they are taken to be

$$\sigma_i = \frac{R(i+1)}{2Q}, \quad (2)$$

where i represents the i th layer in the circular grid (see Fig. 6). Histogram locations are expressed in polar coordinates as

$$\begin{aligned} r_i &= \frac{R(i+1)}{Q}, \\ \theta_j &= \frac{2\pi j}{T}. \end{aligned}$$

3.2 Computational Complexity

The efficiency of DAISY comes from the fact that most computations are separable convolutions and that we avoid computing more than once the histograms common to nearby descriptors. In this section, we give a formal complexity analysis of both DAISY and SIFT and, then, compare them.

3.2.1 Computing DAISY

Recall from Table 1 that DAISY is parameterized with its radius R , number of rings Q , number of histograms in a ring T , and the number of bins in each histogram H . Assuming that the image has P pixels, we begin by computing the orientation layers. In practice, we do not compute gradient norms for each direction separately since they can be computed from the horizontal and vertical ones as

$$\mathbf{G}_\theta = \left(\cos \theta \frac{\partial \mathbf{I}}{\partial x} + \sin \theta \frac{\partial \mathbf{I}}{\partial y} \right)^+. \quad (3)$$

Therefore, for horizontal and vertical gradients, we perform two 1D convolutions with kernels $[1, -1]$ and $[1, -1]^T$, respectively, requiring $2P$ additions to calculate in both directions. Orientation layers are, then, computed from these according to (3) with $2P$ multiplications and P additions for each layer. Then, for each radius quantization level, Q , we perform H convolutions. This is done again as two successive 1D convolutions instead of a single 2D one, thanks to the separability of Gaussian kernels.

Given those gradients, we sample the convolved orientation layers at $Q \times T + 1$ locations for every pixel. For orientations other than 0, an additional shifting operation is required to account for it.

Sampling can be performed by either interpolation or by rounding point locations to the nearest integer. We found that using either method returns roughly equivalent results. Nevertheless, we include both options in the source code we supply [30].

To summarize, computing all the descriptors of an image requires $2H \times Q + 1$ 1D convolutions, $P \times (Q \times T + 1)$ samplings, $2P \times H$ multiplications, and $P \times H$ additions.

TABLE 2
Computation Time in Seconds on an IBM T60 Laptop

Image Size	DAISY	SIFT
800x600	3.8	252
1024x768	6.5	432
1280x960	9.8	651

3.2.2 Computing SIFT

To compute a SIFT descriptor, the image gradient magnitudes and orientations are sampled around the point location at an appropriate scale. Because of the noncircularly symmetric support region and the kernel employed in SIFT, the sampling has to be done at all of the sample locations within the descriptor region. Let us therefore consider the computation of a single descriptor at a single scale where the descriptor is computed over a W_s sample array with S_s histograms of H_s bins.

As does DAISY, SIFT requires image gradients that are computed in the same way and then sampled within the descriptor region at W_s locations. The gradients are then Gaussian smoothed and histograms are formed using trilinear interpolation, that is, each bin is multiplied by a weight of $1 - d$, where d is the distance of the sample from the central value of the bin. The smoothing and interpolation takes $4W_s$ multiplications. Finally, each sample is assigned to a bin and accumulated, which requires W_s multiplications and $S_s \times \binom{2W_s}{S_s} = 2W_s$ summations.

To summarize, the computation of one SIFT descriptor requires W_s samplings, $5W_s$ multiplications, and $2W_s$ summations plus the initial convolution required for gradient computation.

3.2.3 Comparing DAISY and SIFT

For comparison purposes, assume that convolving of a 1D kernel of length N can be done with N multiplications and $N - 1$ summations per pixel. Then, DAISY requires $2H \times Q \times N + 2$ multiplications, $2H \times Q \times N - 1$ summations, and S samplings per pixel.

If we insert the parameters $W_s = 16 \times 16$, $S_s = 4 \times 4$, and $H_s = 8$ for SIFT, as reported in [19], and the parameters we used in this paper for DAISY, $H = 8$, $Q = 3$, and $S = 25$, with an average $N = 5$, we see that SIFT requires 1,280 multiplications, 512 summations, and

256 samplings per pixel, whereas DAISY requires 122 multiplications, 119 summations, and 25 samplings.

Note that computing the descriptor in a different orientation also requires an additional shifting operation in DAISY with $S \times H$ shifts per pixel, whereas this is handled in SIFT during the sampling phase of the gradients with W_s additions.

In any event, a direct comparison of these numbers might be somewhat misleading as one can approximate some of the operations in SIFT in a dense implementation to increase computational efficiency. One might also consider a more efficient implementation of the convolution operation using FFT, which would boost DAISY's performance. The most important speedup difference, however, is due to two facts. First, DAISY descriptors share histograms so that once a histogram is computed for one pixel, it is not computed again for the T other descriptors surrounding this histogram location. Second, the computation pipeline enables a very efficient memory access pattern and the early separation of the histogram layers greatly improves efficiency. This also allows DAISY to be parallelized very easily, and our current implementation [30] allows the use of multiple cores using the OpenMP library. In Table 2, we show typical computation times for various sized images. In Fig. 7, we show the time required to compute DAISY descriptors with varying number of cores on different image sizes. Computation time falls almost linearly with the number of the cores used. The algorithm is also quite suitable for GPU programming which we will pursue in future.

In terms of memory, the precomputed convolved orientation layers require $4Q \times H \times P$ bits and if we also want to precompute all of the descriptors of an image, this will require $4D_s \times P$ bits. In example, for an image of size $1,024 \times 1,024$, these equal to 96 and 800 MB, respectively, for the standard parameter set of $R = 15$, $Q = 3$, $T = 8$, and $H = 8$.

4 DEPTH MAP ESTIMATION

To perform dense matching, we use DAISY to measure similarities across images as shown by Fig. 8. We then feed these measures to a standard graph-cut-based reconstruction algorithm [8]. To properly handle occlusions, we incorporate an occlusion map, which is the counterpart of the visibility maps in other reconstruction algorithms [16]. We do this by introducing an occlusion node with a constant cost in the graph structure. The value of this cost has a significant impact on the proportion of pixels that are

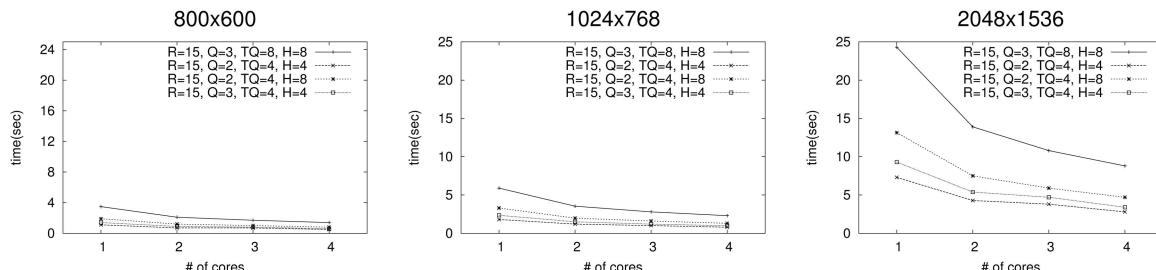


Fig. 7. DAISY Computation Times: Computation times of the DAISY descriptor for all the pixels of an image with various settings on three different sized images. We present the change of the computation time with respect to the number of cores used in parallel.

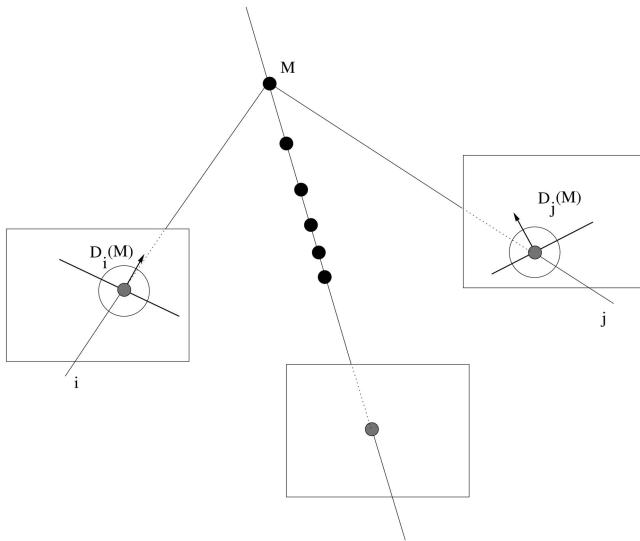


Fig. 8. Problem Definition: The inverse depth is discretized uniformly and the matching score is computed using neighborhoods centered around the projected locations. These neighborhoods are rotated to conform to the orientation of the epipolar lines.

labeled as occluded. As will be discussed in Section 5, we use one image data set to set it to a reasonable value and retain this value for all other experiments shown in this paper. The depth and occlusion maps are estimated by EM, which we formalize below.

We compute the descriptor of every point from its neighborhood. However, for pixels that are close to an occluding boundary, part of the neighborhoods, thereby part of the descriptors, will be different when captured from different viewpoints. To handle this, we exploit the occlusion map and define *binary masks* over our descriptors. We introduce predefined masks that enforce the spatial coherence of the occlusion map, and show that they allow for proper handling of occlusions.

In practice, we assume that we are given at least two calibrated gray-scale images, and we compute the dense depth map of the scene with respect to a particular viewpoint which can either be equal to one of the input view points or it can be a completely different virtual position. We use the calibration information to discretize the 3D space and compute descriptors that take into account the orientation of the epipolar lines, as shown in Fig. 8. In this way, we do not require a rotationally invariant descriptor and take advantage of the fact that DAISY descriptor is very easy to rotate, as described in the previous section.

4.1 Formalization

Given a set of N calibrated images of the scene, we denote their descriptors by $\mathbf{D}_{1:N}$. We estimate the dense depth map \mathbf{Z} for a given viewpoint by maximizing:

$$p(\mathbf{Z}, \mathbf{O} | \mathbf{D}_{1:N}) \propto p(\mathbf{D}_{1:N} | \mathbf{Z}, \mathbf{O})p(\mathbf{Z}, \mathbf{O}), \quad (4)$$

where we also introduced an occlusion map term \mathbf{O} that will be exploited below to estimate the similarities between image locations. As in [8], we enforce piecewise smoothness on the depth map as well as our occlusion map by

penalizing nearby different labels using the Potts model, i.e., $V = \delta(q_i \neq q_j)$ with q_i and q_j are labels of nearby pixels.

For the data-driven posterior, we also assume independence between pixel locations:

$$p(\mathbf{D}_{1:N} | \mathbf{Z}, \mathbf{O}) = \prod_{\mathbf{x}} p(\mathbf{D}_{1:N}(\mathbf{x}) | \mathbf{Z}, \mathbf{O}). \quad (5)$$

Each term $p(\mathbf{D}_{1:N}(\mathbf{x}) | \mathbf{Z}, \mathbf{O})$ of (5) is estimated using our descriptor. Because the descriptor considers relatively large regions, we introduce binary masks computed from the occlusion map \mathbf{O} , as explained in the next section, to avoid including occluded parts into our similarity score.

4.2 Using Masks over the Descriptor

Given the descriptors, we can take the $p(\mathbf{D}_{1:N}(\mathbf{x}) | \mathbf{Z}, \mathbf{O})$ probability to be inversely related to the dissimilarity function

$$D'(\mathbf{D}_i(\mathbf{X}), \mathbf{D}_j(\mathbf{X})) = \frac{1}{S} \sum_{k=1}^S \|\mathbf{D}_i^{[k]}(\mathbf{x}) - \mathbf{D}_j^{[k]}(\mathbf{x})\|_2, \quad (6)$$

where $\mathbf{D}_i(\mathbf{X})$ and $\mathbf{D}_j(\mathbf{X})$ are the descriptors at locations obtained by projecting the 3D point \mathbf{X} (defined by location \mathbf{x} and its depth $\mathbf{Z}(\mathbf{x})$ in the virtual view) onto image i and j , $\mathbf{D}_i^{[k]}(\mathbf{x})$ is the k th histogram $\tilde{\mathbf{h}}$ in $\mathbf{D}_i(\mathbf{x})$, and S is the number of histograms used in the descriptor.

However, as discussed above, we should account for occlusions. The descriptor is computed over image patches and the formulation of (6) is not robust to partial occlusions. Even for a good match, if the pixel is close to an occlusion boundary, the histograms of the occluded parts have no reason to resemble each other.

We, therefore, introduce binary masks $\{\mathcal{M}_m(\mathbf{x})\}$ such as the ones depicted in Fig. 9, which allow DAISY to take into account only the visible parts when computing the distances between descriptors. The mask length is equal to the number of histograms used in the descriptor, the S of Table 1. We use these masks to rewrite the D' of (6) as

$$D = \frac{1}{\sum_{q=1}^S \mathcal{M}^{[q]}} \sum_{k=1}^S \mathcal{M}^{[k]} \|\mathbf{D}_i^{[k]}(\mathbf{x}) - \mathbf{D}_j^{[k]}(\mathbf{x})\|_2, \quad (7)$$

where $\mathcal{M}^{[k]}$ is the k th element of the binary mask \mathcal{M} . Following [8], we define the $p(\mathbf{D}_{1:N}(\mathbf{x}) | \mathbf{Z}, \mathbf{O})$ term of (5) as a Laplacian distribution $\text{Lap}(D(\mathbf{D}_{1:N}(\mathbf{x}) | \mathbf{Z}, \mathbf{O}); 0, \lambda_m)$ with our occlusion handling dissimilarity function.

To select the most likely mask at each pixel location, we rely on an EM algorithm and tried three different strategies:

- The simplest one, depicted by Fig. 9a, involves disabling the histograms that are marked as occluded in the current estimate of the occlusion map \mathbf{O} and obtaining a single binary mask $\mathcal{M}_m(\mathbf{x})$.
- A more sophisticated one is to use the predefined masks depicted by Fig. 9b which have a high spatial coherence. The probability of each mask is computed such that the masks that have large visible areas with similar depth values are favored. We write

$$p(\mathcal{M}_m(\mathbf{x}) | \mathbf{Z}, \mathbf{O}) = \frac{1}{Y} \left(\bar{v}_m + \frac{1}{\sigma_m^2(\mathbf{Z}) + 1} \right), \quad (8)$$

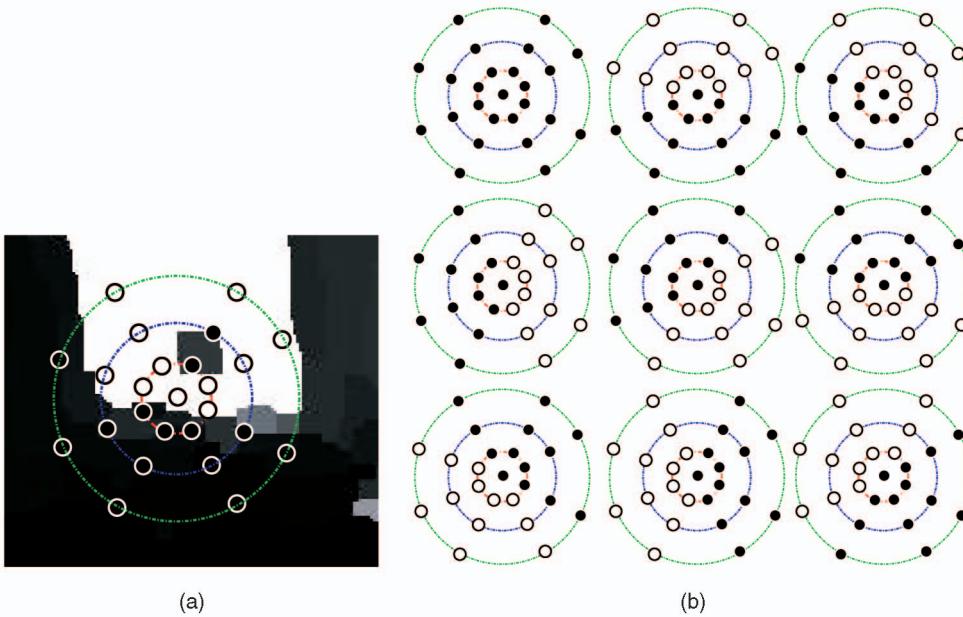


Fig. 9. Binary masks for occlusion handling: We use binary masks over the descriptors to estimate location similarities even near occlusion boundaries. In this figure, a black disk with a white circumference corresponds to “on” and a white disk to “off.” In (a), we use the occlusion map to define the masks and, in (b), predefined masks make it easy to enforce spatial coherence and to speed up the convergence of EM estimation.

where \bar{v}_m is the average visible pixel number, $\sigma_m(\mathbf{Z})$ is the depth variance within the mask region, and Y is the normalization term that is equal to the sum of all mask probabilities.

Then, we take the data posterior to be the weighted sum of individual mask responses as

$$\begin{aligned} p(\mathbf{D}_{1:N}(\mathbf{x})|\mathbf{Z}, \mathbf{O}) \\ = \sum_m p(\mathbf{D}_{1:N}(\mathbf{x})|\mathbf{Z}, \mathbf{O}, \mathcal{M}_m(\mathbf{x}))p(\mathcal{M}_m(\mathbf{x})|\mathbf{Z}, \mathbf{O}). \end{aligned} \quad (9)$$

- The third strategy is a simplified version of the second one, where we only use the result of the most probable mask instead of a mixture.

In the first and third strategy, we use only one mask. By contrast, the second strategy involves a mixture computed from several masks. Note that the mask probabilities are reestimated at each step of the EM algorithm. In our experiments, the second and third strategies always performed better than the first, mainly because they enforce spatial smoothness. The second strategy, however, is computationally more expensive than the third without any perceptible improvement in performance. Therefore, we use only the third strategy in the remainder of the paper. We generally run our EM algorithm for only two to three iterations, which results in better occlusion estimates around occlusion boundaries.

5 EXPERIMENTS AND RESULTS

In this section, we present various experiments we performed in order to measure the performance of DAISY. In Section 5.1, we present results of a parameter sweep experiment we performed to understand and optimize the DAISY parameters with respect to the baseline. We then compare DAISY against other descriptors for depth

estimation purposes. In Section 5.3, we pushed the baseline to very large values to explore the range within which DAISY yields acceptable depth accuracy and the quality of our occlusion estimates. Then, finally in Section 5.4, we test our approach on image pairs with various photometric and geometric transformations and show that it is robust to these and compare our reconstructions with that of a state-of-the-art *multiview* algorithm [26].

5.1 Parameter Selection

To understand the influence of the DAISY parameters of Table 1, we performed two parameter sweep experiments, one in the narrow-baseline case and the other in the wide-baseline case. We used the data set depicted by Fig. 11, which includes laser-scanned ground truth depth and occlusion maps as discussed in [28], [29]. For the narrow-baseline case, we used the image pairs $\{\{1, 2\}, \{2, 3\}, \{3, 4\}, \{4, 5\}, \{5, 6\}\}$. For the wide baseline, we used $\{\{1, 4\}, \{2, 5\}, \{3, 6\}\}$. This guarantees similar baselines within each one of the two groups.

Fig. 10 depicts the results. Correct depth estimates of 80+ percent can be achieved using a less complex descriptor for short baseline. However, as the baseline increases, Fig. 10 suggests that it becomes necessary to use a more complex descriptor at the expense of increased computation and matching time.

Most of the time devoted to descriptor computation is spent on convolutions. It can be reduced by using a smaller number of bins in the histogram (H) or by using a smaller number of layers (Q). We can use $H = 4$ with a little performance loss, but the layer number Q should be chosen carefully depending on the baseline. It appears that two or three layers give similar responses. As far as T , the discretization of the angular space, is concerned, four or eight levels perform similarly for both narrow and wide-baseline cases. However, when increasing the baseline, the

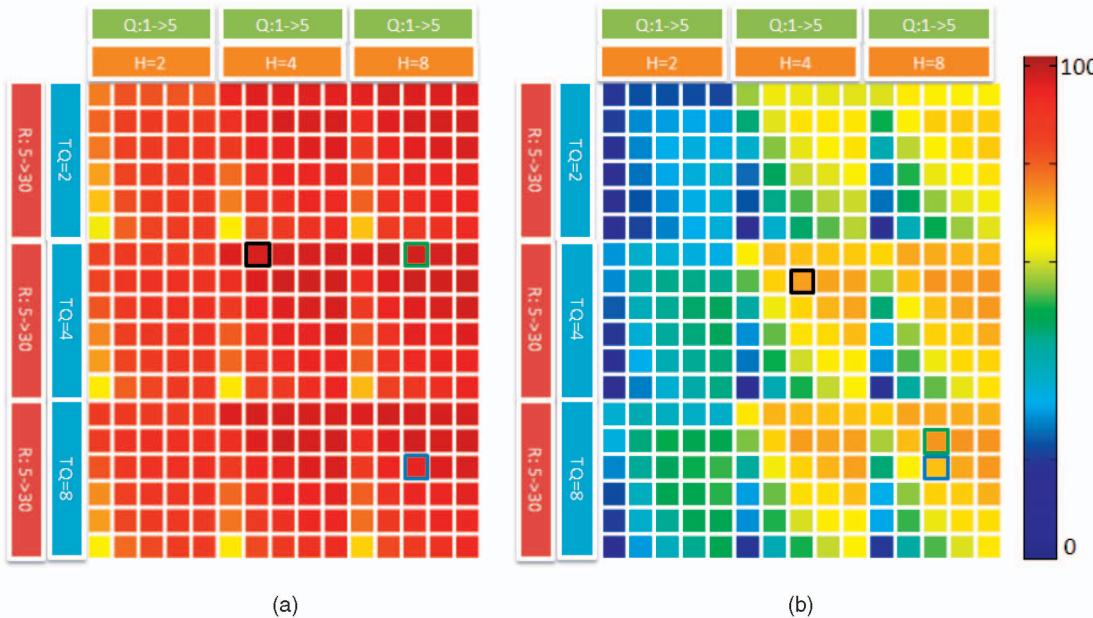


Fig. 10. Parameter sweep test for narrow- and wide-baseline cases: As described by Table 1, there are four parameters that specify the shape and size of DAISY: radius (R), radius quantization (Q), angular quantization (T), and number of bins of the histogram (H). The above figures depict the results of a 4D sweep of these parameters and the color of each square represents the percentage of depths that are correctly estimated. The value associated to a color is given by the color scale on the right. To assess the correctness of an estimated depth, we used laser-scanned depth maps and assumed an estimate as correct if the estimate error is within 1 percent of the scene's depth range. (a) The averaged result for five narrow-baseline image pairs of the Fountain sequence of Fig. 11. The green rectangle denotes the best parameter set for this configuration which is $R = 5, Q = 3, T = 4, H = 8$, resulting in a descriptor of size 104 with a 81.2 percent correct depth estimates. However, upon closer inspection, we see that many other configurations produce similar results (80+ percent). Among these, the configuration $R = 5, Q = 2, T = 4, H = 4$ produces the shortest descriptor length of 36. It is denoted by the black rectangle. (b) The averaged result for three wide-baseline image pairs of the Fountain sequence. The best result, again denoted by a green rectangle, 73 percent correct depth estimate is achieved with $R = 10, Q = 3, T = 8, H = 8$ which yields a 200 length descriptor. However, as in the narrow-baseline case, there are many other configurations that produce a very similar performance (71+ percent) with shorter descriptor sizes. The shortest one (black rectangle) is $R = 10, Q = 3, T = 4, H = 4$ with 52 length. Having multiple configurations that result in similarly high performance shows that we don't really need to change our descriptor parameters depending on the baseline to improve performance, but we can change to meet speed or memory requirements. The configuration we used ($R = 15, Q = 3, T = 8, H = 8$) in all of the other experiments presented in this paper is outlined in blue.

radius R should also be increased, but only up to a point. Going beyond this point causes a loss of discriminative power and a performance drop, especially in the wide-baseline case.

One might argue that there is no need to use more than four bins in the histograms as one could generate the in-between responses of a gradient from the horizontal and vertical directions only. However, this is not the case when summing over a group of pixels because aggregating the low-resolution responses will lose the gradient distribution information of individual pixels and computing a higher resolution version of the histogram from the low-resolution one will not be equal to summing individual high-resolution responses. This is why increasing the histogram resolution makes the descriptor more distinctive at the cost of some computational overhead.

Although descriptor parameters could be adapted depending on the baseline, scene complexity, and texturedness, it is difficult to automate this process. The purpose of this experiment was to see whether there exists a set of parameters that clearly outperform other parameter sets. However, the experimental results suggest that the descriptor is relatively insensitive to parameter choice for an extended range for both narrow and wide-baseline image pairs: different parameter sets produce similar results. Looking at this experiment, we can conclude about three

of the four parameters of the descriptor, namely the radius quantization (Q), angular quantization (T), and histogram quantization (H). However, the effect of the size of the descriptor radius (R) is not so clear since a relatively textured scene is used in the experiment, and we believe the effect of R will be more apparent for less textured scenes. Although we do not have a data set with a ground truth depth map of such a scene and therefore cannot accurately quantify this effect, we have observed that using a larger R is beneficial for less textured scenes from other data sets used in this paper. Hence, in practice, we use the most generic parameter set $R = 15, Q = 3, T = 8, H = 8$ for all of the experiments presented in this paper. Admittedly, this produces a longer descriptor than strictly necessary, but it performs well for both narrow and wide baselines. However, depending on the application DAISY is used for, the parameters can be set accordingly. For example, if it is known that input images have a short baseline and scene is more or less textured, $R = 10, Q = 3, T = 4, H = 4$ will produce good results with a shorter footprint of length 52.

5.2 Comparison with Other Descriptors

To compare DAISY's performance with that of the other descriptors we again used the data set of Fig. 11 and also the data set of Fig. 12, for which laser-scanner data are available as well. Arguably our results on the data of Fig. 11 should

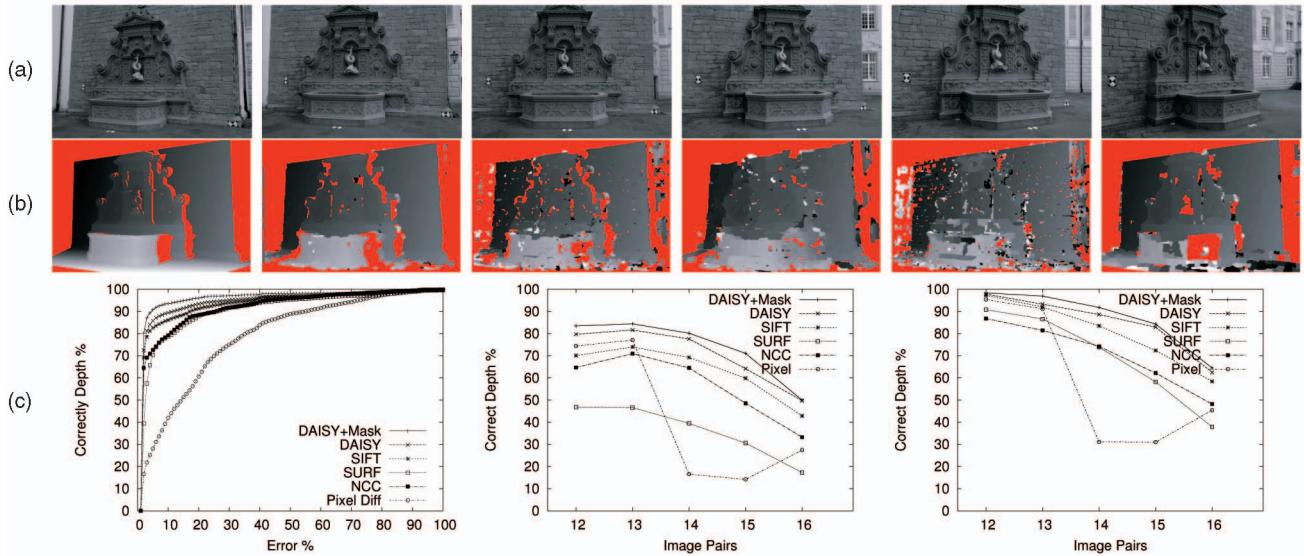


Fig. 11. Comparing different descriptors: Fountain Sequence [28]. (a) In our tests, we match the leftmost image against each one of the other five. (b) The laser-scan depth map we use as a reference and five depth maps computed from the first and third images. From left to right, we used DAISY, SIFT, SURF, NCC, and Pixel Difference. (c) The leftmost plot shows the corresponding distributions of deviations from the laser-scan data, expressed as a fraction of the scene's depth range. The other plots summarize these distributions for the five stereo pairs of increasing baseline with discrete error thresholds set to be 1 and 5 percent of the scene's depth range, respectively. Each data point represents a pair where the baseline increases gradually from left to right and individual curves correspond to DAISY with masks, DAISY without masks, SIFT, SURF, NCC, and Pixel Difference. In all cases, DAISY does better than the others and using masks further improves the results.

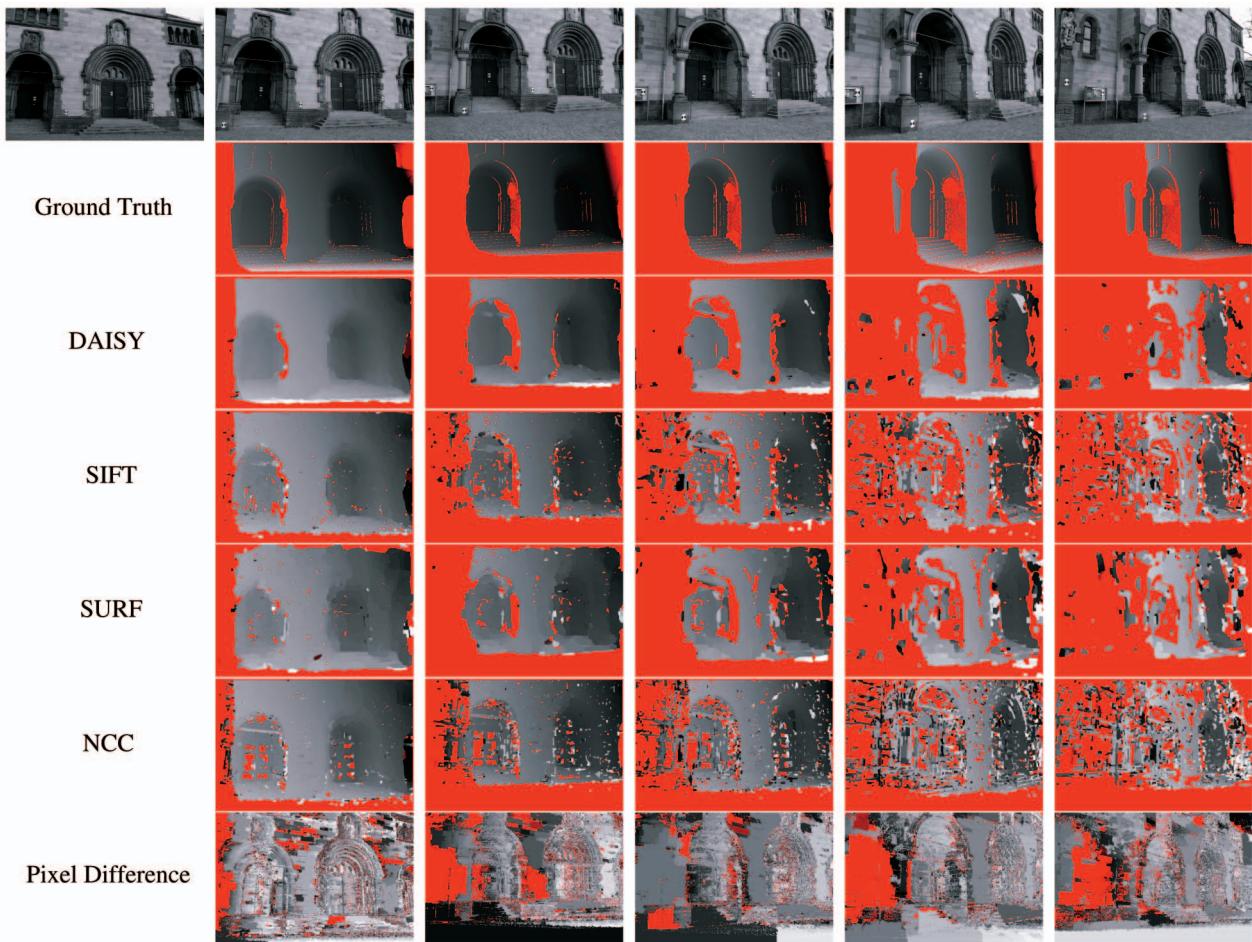


Fig. 12. Comparing different descriptors: HerzJesu Sequence [28]. As in the Fountain sequence of Fig. 11, we use different descriptors to match the leftmost images with each one of the other images. We compute the depth map in the reference frame of these. Second row: Ground truth depth maps with overlaid occlusion masks. Remaining rows: Depth maps computed using DAISY, SIFT, SURF, NCC, and Pixel differencing, in that order.

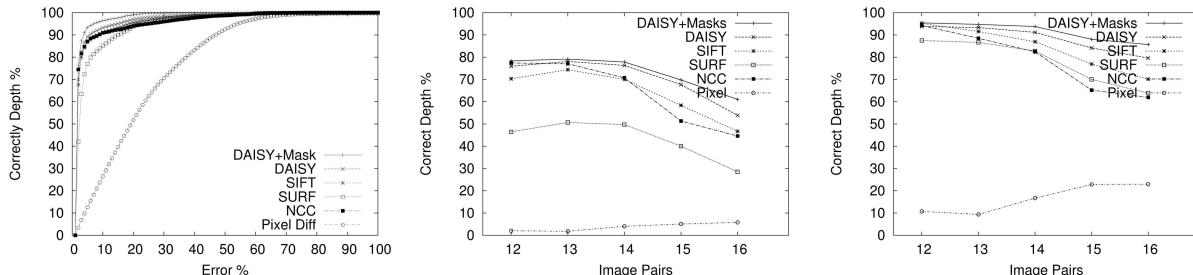


Fig. 13. Quantitative comparison results for Fig. 12: The leftmost plot shows the corresponding distributions of deviations from the laser-scan data, expressed as a fraction of the scene's depth range. The other plots summarize these distributions for the five stereo pairs of increasing baseline with discrete error thresholds set to be 1 and 5 percent of the scene's depth range, respectively. Each data point represents a pair where the baseline increases gradually from left to right and individual curves correspond to DAISY with masks, DAISY without masks, SIFT, SURF, NCC, and Pixel Difference. In all cases, DAISY does better than the others and using masks further improves the result.

be treated with caution since we have used these images to set our parameters. However, we have done no such thing with the data of Fig. 12 and obtain very similar results.

We used DAISY with occlusion masks, DAISY, SIFT, SURF, NCC, and Pixel differencing to densely compute matching scores. They are then all handled similarly, as described in Section 4, to produce depth maps. The only difference is that we do not use binary masks to modify matching scores for descriptors other than DAISY. All of the region-based descriptors are computed perpendicular to the epipolar lines; SURF and SIFT descriptors are 128-length vectors and NCC is 11×11 .

For the data set of Fig. 11, the leftmost image in the first row is matched against each one of the other five, which implies a wider and wider baseline. The second row depicts the laser-scanner data on the left and the depth maps computed from the first and third images using DAISY, SIFT, SURF, NCC, and pixel differencing. The third row summarizes the comparison of different descriptors against DAISY. The leftmost graph shows the result for the first and third image pairs by plotting the percentage of correctly estimated depths against the amount of allowed error which is represented as a fraction of scene's depth range and remaining graphs summarize these curves for all image pairs at discrete error levels. In one case, the depths are considered to be correct if they are within 1 percent of the scene's depth range and in the other within 5 percent. We present results using DAISY with and without using the occlusion masks. DAISY by itself outperforms the other descriptors and the masks provide a further boost.

For the data set of Fig. 12, we match the leftmost image with each one of the other first row images in turn and compute the depth map with respect to the latter image. We display the ground truth maps in the second row and show the estimated depth maps for different descriptors in the remaining rows. Fig. 13 depicts the quantitative results for this data set as in the previous data set.

Both of the data sets show that DAISY performs better than all of the other descriptors. Note that, although the results of SIFT and DAISY are close for the 5 percent threshold, DAISY finds substantially more correct depths for the 1 percent threshold. This indicates that the depths found using DAISY are more accurate than those found using SIFT.

5.3 Occlusion Handling

We tested the performance of our occlusion detection scheme with the extended version of the HerzJesu sequence

of Fig. 12, depicted by Fig. 14. The matching is done using two images, one from the first row and one from the first column, and the depth map is shown in the referential of the second image. The resulting depth map is displayed on the intersection of the respective column and row together with the ground truth on the diagonal. We use different colors to highlight correctly estimated occlusions, missed occlusions, and falsely labeled occlusions. An example pair of images from this sequence is shown in Fig. 15. Table 3 gives the percentage of the correctly estimated depths in visible areas where the correctness threshold is set to 5 percent of the scene's depth range.

As discussed in Section 4, the value of the occlusion cost in the graph structure has a direct influence on how many pixels are labeled as occluded. In Fig. 16, we plot ROC curves obtained by using this for all of the image pairs of Fig. 14. Here, each data point represents the result with a different occlusion cost, true positive rate shows the percentage of the visible areas we detect as visible, and false positive rate shows the amount of missed occlusions. By using these plots, we picked a single value, 25 percent of the maximum cost, for the occlusion cost for all the results shown in this paper.

To show the effect of the baseline on the performance, we plot the area under the curve (AUC) of these ROC curves with respect to the angle between the cameras for all the image pairs of Fig. 14. This curve shows that our approach works for a wide variety of camera configurations and is robust up to 30-40 degree changes. We also show two example ROC curves for narrow- and wide-baseline cases in the same figure.

The progress of our EM-based occlusion detection algorithm can be seen in Fig. 17. In this figure, we give an example for the evolution of the depth map with occlusion estimates at each iteration. The initial estimate is quickly improved in the next iteration with occlusions receding and new depths being estimated for these previously occluded-marked regions. We see that further iterations do not improve the result significantly, and in practice we stop the process after one or two iterations.

5.4 Robustness to Image Transformations

Although we designed DAISY with only wide-baseline conditions in mind, it exhibits the same robustness as histogram-based descriptors to changes in contrast (Fig. 1), scale (Fig. 2), image quality (Fig. 3), viewpoint (Figs. 4 and 18), and brightness (Fig. 19). In these figures, we also present the

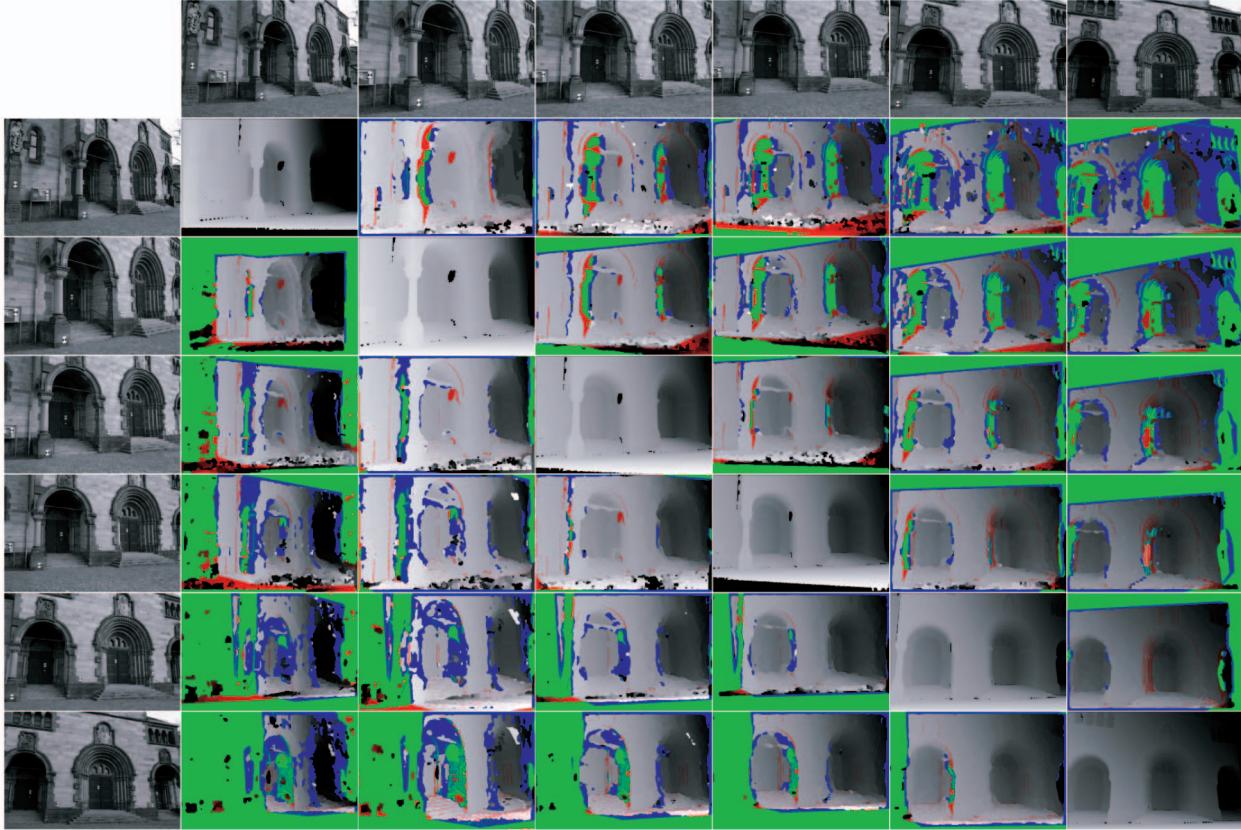


Fig. 14. HerzJesu Grid: By using two images, one from the leftmost column and one from the upper row, we compute depth and occlusion maps from the viewpoint of the row image. In the diagonal, we display the ground truth depth maps. We marked the correctly detected occlusions with green, incorrectly detected ones with blue, and the missed ones with red. From this figure, it is apparent that DAISY can handle quite large baselines without losing too much from its accuracy, as can be seen from Table 3.

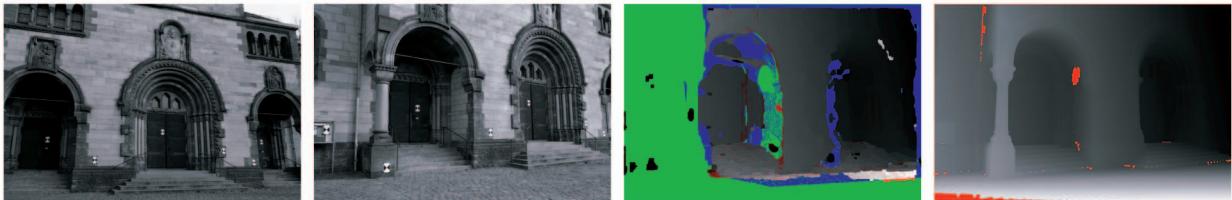


Fig. 15. HerzJesu Close-up: Depth map is computed from the second image's point of view using DAISY. Correctly detected occlusions are shown in green, incorrectly detected ones with blue, and the missed ones with red. The last image is the laser-scanned ground truth depth map.

depth maps obtained using NCC and SIFT which include more artifacts than ours. This result is noteworthy because, although it is a well-known fact that histogram-based

descriptors are robust against these transforms at *feature point locations*, these experiments show that such robustness can also be found at many other point locations.

To compare our method to one of the best current techniques [26], we ran our algorithm on two sets of image pairs that were used in that paper, the Rathaus sequence of Fig. 20 and the Brussels sequence of Fig. 21. But instead of using the original $3,072 \times 2,048$ images, whose resolution is high enough for apparently blank areas to exhibit usable texture, we used 768×512 images in which this is not true. DAISY nevertheless achieved visually similar results.

Fig. 21 also highlights the effectiveness of our occlusion handling. When using only two images, the parts of the church that are hidden by people in one image and not in the other are correctly detected as occluded. When using three images, the algorithm returns an almost full depth map that lets us erase the people in the synthetic images we produce.

TABLE 3
Correctly Estimated Depth Percentage for Fig. 14

Frame	1	2	3	4	5	6
Number	-	90.8	88.7	86.9	87.5	90.1
1	86.9	-	94.3	92.4	91.4	92.9
2	85.5	90.3	-	93.2	93.9	96.4
3	83.6	87.1	93.0	-	95.4	96.8
4	83.9	86.3	91.7	93.3	-	97.7
5	86.4	89.2	94.5	95.6	96.7	-

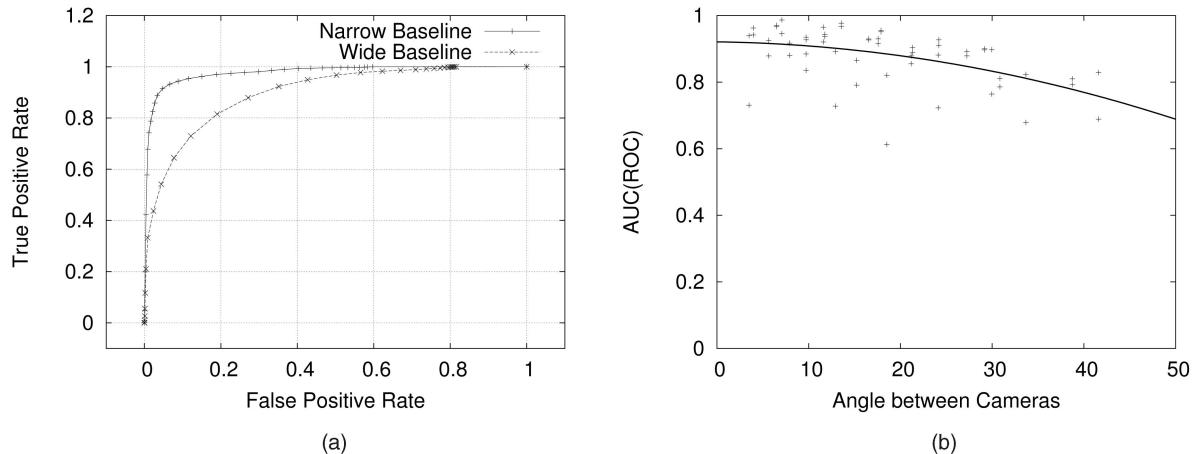


Fig. 16. ROC for occlusion threshold: We plot ROC curves for the selection of the occlusion threshold for all of the image pairs of Fig. 14. (a) The ROC curves of the narrow-baseline image pair {4, 5} and the wide-baseline image pair {1, 5}. (b) The area under the curve (AUC) of the ROC graphs of the image pairs with respect to the angle between the principal axes of the cameras. The result of each such pair is represented with a data point and the curve shows the fitted line to these.

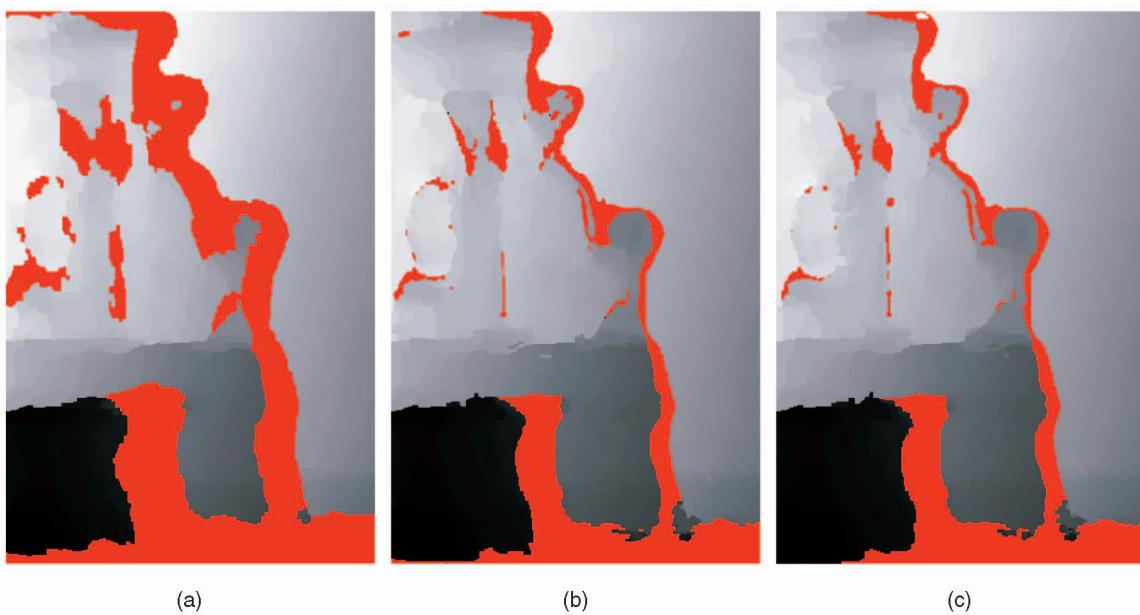


Fig. 17. Evolution of the occlusions during EM: The three images show the evolution of the occlusion estimate during the iterations of the EM for Fig. 11 images. The initial solution (a) is quickly improved even after a single iteration (b) and does not change much thereafter (c).



Fig. 18. Valencia Cathedral: The reconstruction results of the exterior of the Valencia Cathedral from two very different viewpoints. The depth map is computed from the second image's point of view.

6 CONCLUSION

In this paper, we introduced DAISY, a new local descriptor, which is inspired from earlier ones such as SIFT and GLOH

but can be computed much more efficiently for dense-matching purposes. Speed increase comes from replacing weighted sums used by the earlier descriptors by sums of



Fig. 19. Brightness Change: We compute the depth map from the second image's point of view using DAISY. There is a brightness change between the two images.

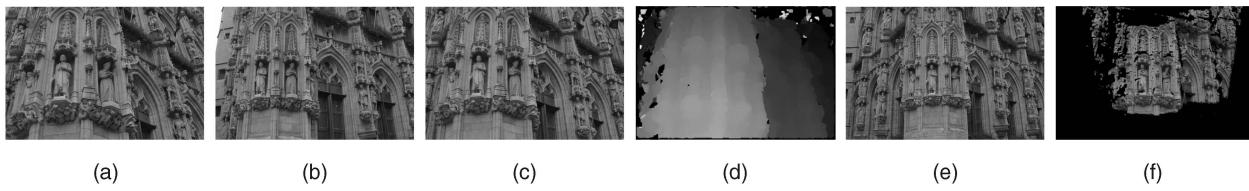


Fig. 20. Results on low-resolution versions of the Rathaus images [27]: (a)-(c) Three input images of size 768×512 instead of the $3,072 \times 2,048$ versions that were used in [26]. (d) Depth map computed using all three images. (e) A fourth image not used for reconstruction. (f) Image synthesized using the depth map and the image texture in (a) with respect to the view point of (e). Note how similar it is to (e). The holes are caused by the fact that a lot of the texture in (e) is not visible in (a).

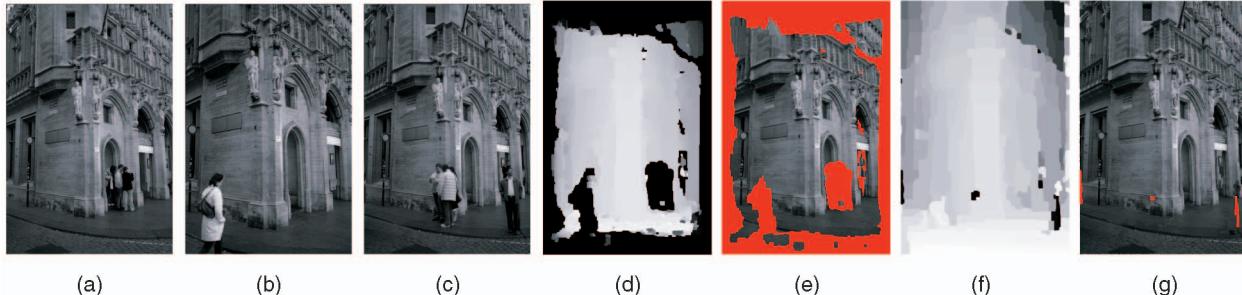


Fig. 21. Low-resolution versions of the Brussels images [26]: (a)-(c) Three 768×510 versions of the original $2,048 \times 1,360$ images. (d) and (e) The depth map computed using images (a) and (b) seen in the perspective of image (c) and the corresponding resynthesized image. Note that the locations where there are people in one image and not in the other are correctly marked as occlusions. (f) and (g) The depth map and synthetic image generated using all three images. Note that the previously occluded areas are now filled and that the people have been erased from the synthetic image.

convolutions, which can be computed very quickly and from using a circularly symmetrical weighting kernel. The experiments suggest that, although pixel differencing or correlation is good for short-baseline stereo, wide baseline requires a more advanced measure for comparison. We showed DAISY to be very effective for this purpose.

Our method gives good results, even when using small images for stereo reconstruction. This means that we could use our algorithm to process video streams whose resolution is often lower than that of still images. When dealing with slanted surfaces and foreshortening, these results could be further improved by explicitly taking into account 3D surface orientation and warping the DAISY grid accordingly, which would not involve any significant computational overhead. This would fit naturally in a warp

stereo approach [23] in which we would begin with unwarped detectors to compute a first surface estimate, use the corresponding orientations to warp the detectors, and iterate.

Computing our descriptor primarily involves performing Gaussian convolutions, which are amenable to hardware exportation or GPU implementation. This could lead to real-time, or even faster, computation of the descriptor for all image pixels. This could have implications beyond stereo reconstruction because dense computation of image descriptors is fast becoming an important technique in other fields, such as object recognition [7], [18]. To encourage such developments, a C++ and MATLAB implementation of DAISY is available for download from our webpage [30].

REFERENCES

- [1] L. Alvarez, R. Deriche, J. Weickert, and J. Sanchez, "Dense Disparity Map Estimation Respecting Image Discontinuities: A PDE and Scale-Space Based Approach," *J. Visual Comm. and Image Representation*, vol. 13, nos. 1/2, pp. 3-21, Mar. 2002.
- [2] N. Ayache and F. Lustman, "Fast and Reliable Passive Trinocular Stereovision," *Proc. Int'l Conf. Computer Vision*, June 1987.
- [3] H.H. Baker and T.O. Binford, "Depth from Edge and Intensity Based Stereo," *Proc. Int'l Joint Conf. Artificial Intelligence*, vol. 2, pp. 631-636, Aug. 1981.
- [4] H. Bay, T. Tuytelaars, and L. Van Gool, "SURF: Speeded Up Robust Features," *Proc. European Conf. Computer Vision*, 2006.
- [5] A.C. Berg and J. Malik, "Geometric Blur for Template Matching," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 607-614, 2001.
- [6] S. Birchfield and C. Tomasi, "A Pixel Dissimilarity Measure that is Insensitive to Image Sampling," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 20, no. 4, pp. 401-406, Apr. 1998.
- [7] A. Bosch, A. Zisserman, and X. Munoz, "Scene Classification via pLSA," *Proc. European Conf. Computer Vision*, 2006.
- [8] Y. Boykov, O. Veksler, and R. Zabih, "Fast Approximate Energy Minimization via Graph Cuts," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 23, no. 11, pp. 1222-1239, Nov. 2001.
- [9] M.Z. Brown, D. Burschka, and G.D. Hager, "Advances in Computational Stereo," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 25, no. 8, pp. 993-1008, Aug. 2003.
- [10] F. Devernay and O.D. Faugeras, "Computing Differential Properties of 3D Shapes from Stereoscopic Images without 3D Models," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 208-213, June 1994.
- [11] O.D. Faugeras and R. Keriven, "Complete Dense Stereovision Using Level Set Methods," *Proc. European Conf. Computer Vision*, June 1998.
- [12] L. Fei-Fei and P. Perona, "A Bayesian Hierarchical Model for Learning Natural Scene Categories," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2005.
- [13] D. Geiger, B. Ladendorf, and A. Yuille, "Occlusions and Binocular Stereo," *Int'l J. Computer Vision*, vol. 14, pp. 211-226, 1995.
- [14] S.S. Intille and A.F. Bobick, "Disparity-Space Images and Large Occlusion Stereo," *Proc. European Conf. Computer Vision*, pp. 179-186, May 1994.
- [15] T. Kanade and M. Okutomi, "A Stereo Matching Algorithm with an Adaptive Window: Theory and Experiment," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 16, no. 9, pp. 920-932, Sept. 1994.
- [16] V. Kolmogorov and R. Zabih, "Multi-Camera Scene Reconstruction via Graph Cuts," *Proc. European Conf. Computer Vision*, May 2002.
- [17] K.N. Kutulakos and S.M. Seitz, "A Theory of Shape by Space Carving," *Int'l J. Computer Vision*, vol. 38, no. 3, pp. 197-216, July 2000.
- [18] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond Bags of Features: Spatial Pyramid Matching for Recognizing Natural Scene Categories," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2006.
- [19] D.G. Lowe, "Distinctive Image Features from Scale Invariant Keypoints," *Int'l J. Computer Vision*, vol. 20, no. 2, pp. 91-110, 2004.
- [20] G. Medioni, C.K. Tang, and M.S. Lee, "Tensor Voting: Theory and Applications," *Proc. Reconnaissance des Formes et Intelligence en Artificielle*, 2000.
- [21] K. Mikolajczyk and C. Schmid, "A Performance Evaluation of Local Descriptors," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 27, no. 10, pp. 1615-1630, Oct. 2005.
- [22] F. Porikli, "Integral Histogram: A Fast Way to Extract Histograms in Cartesian Spaces," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, vol. 1, pp. 829-836, 2005.
- [23] L.H. Quam, "Hierarchical Warp Stereo," *Readings in Computer Vision: Issues, Problems, Principles, and Paradigms*, pp. 80-86, Morgan Kaufmann, 1987.
- [24] S. Roy and I.J. Cox, "A Maximum-Flow Formulation of the N-Camera Stereo Correspondence Problem," *Proc. Int'l Conf. Computer Vision*, pp. 492-499, 1998.
- [25] D. Scharstein and R. Szeliski, "A Taxonomy and Evaluation of Dense Two-Frame Stereo Correspondence Algorithms," *Int'l J. Computer Vision*, vol. 47, nos. 1-3, pp. 7-42, Apr.-June 2002.
- [26] C. Strecha, R. Fransens, and L. Van Gool, "Combined Depth and Outlier Estimation in Multi-View Stereo," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2006.
- [27] C. Strecha, T. Tuytelaars, and L. Van Gool, "Dense Matching of Multiple Wide Baseline Views," *Proc. Int'l Conf. Computer Vision*, 2003.
- [28] C. Strecha, W. von Hansen, L. Van Gool, P. Fua, and U. Thoennessen, "On Benchmarking Camera Calibration and Multi-View Stereo for High Resolution Imagery," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, June 2008.
- [29] C. Strecha, Multi-View Evaluation, <http://cvlab.epfl.ch/data/>, 2008.
- [30] E. Tola, Daisy Code, <http://cvlab.epfl.ch/software/>, 2008.
- [31] T. Tuytelaars and L. Van Gool, "Wide Baseline Stereo Matching Based on Local, Affinely Invariant Regions," *Proc. British Machine Vision Conf.*, pp. 412-422, 2000.
- [32] S.A. Winder and M. Brown, "Learning Local Image Descriptors," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, June 2007.
- [33] J. Yao and W.-K. Cham, "3D Modeling and Rendering from Multiple Wide Baseline Images," *Signal Processing: Image Comm.*, vol. 21, pp. 506-518, 2006.



virtual view synthesis.

Engin Tola received the undergraduate degree in electrical and electronics engineering from the Middle East Technical University (METU), Turkey, in 2003. Afterward, he received the MSc degree in signal processing in 2005 from METU, and he is presently a PhD student at the Ecole Polytechnique Fédérale de Lausanne (EPFL), Switzerland, in the Computer Vision Laboratory. His research interests focus on point descriptors, scene reconstruction, and



3D camera tracking, and object recognition.



Pascal Fua received the engineering degree from the Ecole Polytechnique, Paris, in 1984 and the PhD degree in computer science from the University of Orsay in 1989. He joined EPFL (Swiss Federal Institute of Technology) in 1996, where he is now a professor in the School of Computer and Communication Science. Before that, he worked at SRI International and at INRIA Sophia-Antipolis as a computer scientist. His research interests include shape modeling and motion recovery from images, human body modeling, and optimization-based techniques for image analysis and synthesis. He has (co)authored more than 150 publications in refereed journals and conferences. He has been an associate editor of the *IEEE Transactions for Pattern Analysis and Machine Intelligence* and has been a program committee member and an area chair of several major vision conferences. He is a senior member of the IEEE.