# Learning to Assign Orientations to Feature Points

Kwang Moo Yi[1,*]     Yannick Verdie [1,*]     Pascal Fua[1]     Vincent Lepetit[2]

[1]Computer Vision Laboratory, École Polytechnique Fédérale de Lausanne (EPFL)
[2]Institute for Computer Graphics and Vision, Graz University of Technology

{kwang.yi, yannick.verdie, pascal.fua}@epfl.ch, lepetit@icg.tugraz.at

## Abstract

*We show how to train a Convolutional Neural Network to assign a canonical orientation to feature points given an image patch centered on the feature point. Our method improves feature point matching upon the state-of-the art and can be used in conjunction with any existing rotation sensitive descriptors. To avoid the tedious and almost impossible task of finding a target orientation to learn, we propose to use Siamese networks which implicitly find the optimal orientations during training. We also propose a new type of activation function for Neural Networks that generalizes the popular ReLU, maxout, and PReLU activation functions. This novel activation performs better for our task. We validate the effectiveness of our method extensively with four existing datasets, including two non-planar datasets, as well as our own dataset. We show that we outperform the state-of-the-art without the need of retraining for each dataset.*

## 1. Introduction

Feature points are an essential and ubiquitous tool in computer vision, and extensive research has been conducted on both detectors [3, 6, 22, 25, 27, 32, 48] and descriptors [2, 6, 22, 25, 32, 38, 43, 47], including using statistical approaches [33, 46]. However, the assignment of a canonical orientation, which is an important common step, has received almost no individual attention, probably since the dominant orientation of *SIFT* [25] is considered to give good results.

However, this is not necessarily true. In complex scenes, feature points lie on non-planar surfaces and their appearance can be drastically altered by viewpoint and illumination changes. This can easily produce errors in orientation estimates as shown in Fig. 1. In addition, rotation invariant descriptors [7, 14, 42] are not a definitive solution ei-
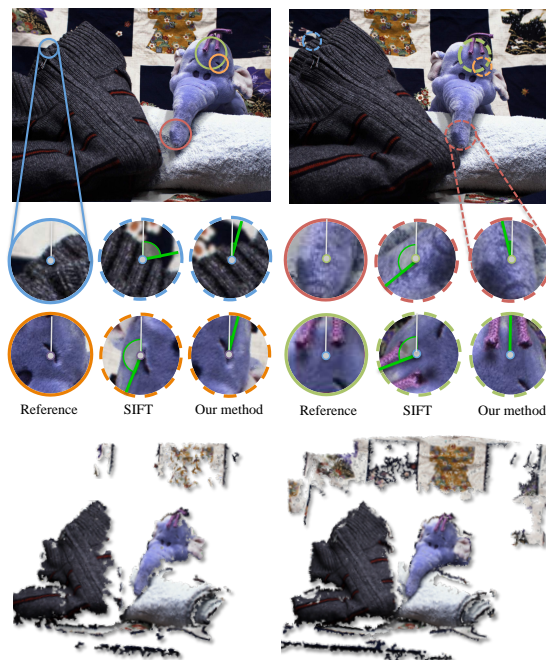


Figure 1. Multi-View Stereo (MVS) [44, 45] reconstruction with the orientation assignment of *SIFT* [25] and our orientation assignment. **Top:** two of the original images from [1] used for MVS. **Middle:** enlarged feature point regions in groups of three; *left* reference region from the left image, *center* region rotated back with *SIFT* orientations, *right* region rotated back with our learned orientations. Estimation errors are denoted by a green arc. **Bottom left:** MVS results with *SIFT* orientations, and **Bottom right:** MVS results with our orientations. As shown, due to viewpoints changes on non-planar surfaces, *SIFT* orientations are not stable. On the contrary, orientations provided by our method are stable, which leads to better reconstructions. 46272 vertices were obtained using *SIFT* orientations, and 84087 vertices with our orientations. *Edge Foci* feature points [48] were used in conjunction with *Daisy* [38] descriptors for both methods. [1]

ther as these descriptors discard rotation sensitive information which can be useful when ideal orientations are given. Thus, as we will show in our experiments, higher match-

---

[*]First two authors contributed equally.
[1]Figures are best viewed in color.

ing performances can be achieved with rotation sensitive descriptors and better orientation assignments.

In this paper, we show how to remedy this problem by training a regressor to estimate better orientations for matching, and to boost the performance of existing rotation sensitive descriptors. We train a Convolutional Neural Network to predict an orientation, given a patch around a feature point. To avoid the difficult task of finding the canonical orientation to learn, we treat the orientation to learn as an implicit variable, by training a Siamese network [9, 12] similar to descriptor learning methods [33, 46]. Also, to allow our method to work in conjunction with any existing rotation sensitive descriptors such as *SIFT* [25], *SURF* [6], and the learning-based *VGG* [34], we consider the descriptor component as a black box when learning.

We also propose a new activation function for the neurons based on Generalized Hinging Hyperplanes (GHH) [41], which plays a key role in our method. We will show that it generalizes the popular ReLU and maxout [16] activation functions, as well as the recent PReLU [17] activation function, with better performance for our task.

To evaluate the performance of descriptors with orientations from the proposed method, we use datasets with both planar or far away objects [26, 40, 48] and 3D objects [1, 36]. In addition, we created our own dataset as well, to further enrich the dataset with complex camera movements, such as in-plane rotations and viewpoint changes. We demonstrate that the proposed method gives significant improvement over the state-of-the-art for *all* the datasets, without the need of re-training for each dataset.

In the remainder of this paper, we first discuss related work, introduce our learning framework, detail our method as well as the proposed activation function. We then present our experimental results demonstrating the effectiveness of our orientation assignment compared to the state-of-the-art. We also investigate the influence of the proposed activation and of the datasets, and we conclude with several application results.

## 2. Related Work

As shown in the survey of [15], the importance of having a good orientation estimation has been overlooked, and thought to be a not very important step which either feature point detector or descriptor has to perform. The widely-used solution for assigning an orientation to a feature point is to use the dominant orientation of *SIFT* [25]. However, as pointed out by [24], dominant orientation-based methods do not work well for arbitrary positions, although it has critical impact on the descriptor performances [23]. Nevertheless, here we provide a brief review of existing methods related to orientation assignment and our method.

**Orientation assignment of feature point detectors.** In *SIFT* [25], histograms of gradient orientations are used to determine the dominant orientation. It remains the most popular method and has also been extended to 3D [4]. *SURF* [6] uses Haar-wavelet responses of sample points to extract the dominant orientation. *MOPs* [11] simply uses the gradient at the center of a patch after some smoothing for robustness to noise. *ORB* [32] uses image moments to compute the center of mass as well as the main orientation. *HIP* [37] considers intensity differences over a circle centered around the feature point to estimate the orientation. Although this is rather fast, it is also very sensitive to noise.

In summary, despite the variation, the main idea of these methods remain the same: finding a reliable dominant orientation in their respective ways. Thus, when computation time constraints are not too drastic, using the SIFT orientation remains to be the first solution to try [15].

**Rotation invariant descriptors.** As existing orientation assignment methods are not always robust enough to guarantee good matching performances, interest has been drawn to descriptors which are inherently rotation invariant [14, 42]. *MROGH* [14] uses local intensity order pooling with rotation invariant gradients, and *LIOP* [42] constructs the descriptor in a similar way but with a different strategy for aggregating the gradient information. *BRISK* [22] and *FREAK* [2] claims rotation invariance as well, but they still depend on the orientation estimation which is included in the descriptor extraction process.

Besides descriptors that are rotation invariant by construction, [21] uses concentric rings for generating orientation histogram bins with spin images, and a specific distance function for rotation invariant matching. *sGLOH* [7] also proposes to use a rotation invariant distance function which computes distances for all possible rotation combinations and takes the minimum. The authors further extend their method by proposing a general method for histogram-based feature descriptors taking into account the main orientation of the scene [8].

Although these methods may be better than the original *SIFT* descriptor [25], *SIFT* descriptor combined with our learning-based orientation estimation outperforms them, as we will show in the experiments. This is probably due to the fact that rotation sensitive information is discarded when computing these descriptors. Furthermore, [8] is only applicable when the entire scene is the object of interest, and is impractical as the main orientation is obtained by computing all the possible matching pairs of features to keep the configuration with best matches.

**Learning-based methods.** Learning-based methods have been already used in the context of feature point matching, but only for problems other than orientation assignment of general feature points. For example, [18] learns to predict

(a) Matches with *SIFT* orientations     (b) Matches with our orientations

Figure 2. Image matching example from the duckhunt sequence in the *Viewpoints* dataset with *SIFT* descriptors using the orientation assigned by (a) *SIFT* and (b) our method. The yellow lines denote the inlier matches after RANSAC. The homography is correctly estimated only when using our orientations.

the pose of patches, but uses one regressor per patch, which is not a viable solution for general feature points. [33, 46] use Siamese networks—as we do—to directly compare image patches [46], or to learn to compute descriptors [33]. *VGG* [34] as well as [13, 39] also learn descriptors, through convex optimization, boosting, and greedy optimization, respectively.

One caveat in these learning-based descriptors is that they still rely on the orientation estimation of local feature detectors which are traditionally handcrafted. Moreover, they typically use the Brown dataset [10] for learning, with patches extracted using ground truth orientations from Structure from Motion (SfM) techniques. This ground truth orientation assignment is not something one can expect to have in practical use, and may lead to performance degradation when tested on other data with inaccurate orientation assignments [34]. These methods will also benefit from better orientation assignments on test time, as we will show in our experiments with *VGG*.

## 3. Method

In this section we first introduce our learning strategy, then formalize it. We also describe our activation function based on GHH.

### 3.1. Canonical Orientation as an Implicit Variable

As illustrated in Fig. 2, orientation assignment plays a critical role in the descriptor matching performances. However, a major problem we face in our approach is that it is not clear which orientation should be learned. For example, one can try to learn to predict the dominant orientation of *SIFT* [25], or maybe the median of the dominant orientations for the same feature point extracted from multiple images. However, there is no guarantee that the orientations retrieved from such approach is the ideal canonical orientation we want to learn. Our early experiments, based on such heuristics to decide which orientation should be learned, remained unfruitful.

Since it is hard to define a canonical orientation to learn, we instead take into account that it is actually the descriptor distances of the feature points that are important, not the orientation values themselves. We formulate the problem by learning to assign orientations which minimize the descriptor distances of pairs of local features corresponding to the same physical point. In this way, we do not have to decide which orientations should be learned. We let the learning optimization find which orientations are both reliably predictable and improve the matching performance. We formalize this approach in the next subsection.

### 3.2. Formalization

Our approach is related to Siamese networks used for descriptor learning [33, 46], but the loss function and its computation are different since we learn to estimate the orientation and not the descriptor itself. In fact, we treat the descriptor as a black box so that various rotation variant descriptors can be used. However, this is not necessarily a restriction, and can be easily adapted to include learning of the descriptors as well.

Our training data is made of pairs of image patches centered on feature points in two images but corresponding to the same physical 3D points. We minimize a loss function $\sum_i \mathcal{L}_i$ over the parameters $\mathbf{W}$ of a CNN, with

$$\mathcal{L}(\mathbf{p}_i) = \left\| g(\mathbf{p}_i^1, f_{\mathbf{W}}(\mathbf{p}_i^1)) - g(\mathbf{p}_i^2, f_{\mathbf{W}}(\mathbf{p}_i^2)) \right\|_2^2 \quad , \quad (1)$$

where $\mathcal{L}_i = \mathcal{L}(\mathbf{p}_i)$, the pairs $\mathbf{p}_i = \{\mathbf{p}_i^1, \mathbf{p}_i^2\}$ are pairs of image patches from the training set, $f_{\mathbf{W}}(\mathbf{p}_i^*)$ denotes the orientation computed for image patch $\mathbf{p}_i^*$ using a CNN with parameters $\mathbf{W}$, and $g(\mathbf{p}_i^*, \theta_i^*)$ is the descriptor for patch $\mathbf{p}_i^*$ and orientation $\theta_i^*$. As discussed in the previous subsection, there is no target orientation in the loss function of Eq. (1): the predicted orientations will be optimized implicitly during training.

**Predicting an angle.** Learning angles requires a special care. Directly predicting an angle with a CNN did not work well in our early experiments, probably because the periodicity of $f_{\mathbf{W}}(\mathbf{p}_i^*)$ in Eq. (1) generates many local minima. An alternative way would be to learn to provide histogram-like outputs, which then can be used with $\arg\max$ to give angular outputs, in a way reminiscent of SIFT. However, this approach also did not work well, as the estimated orientations have to be discretized and the network becomes too large when we want fine resolutions.

To alleviate the problem of periodicity, similarly to how manifolds are embed in [30, 31], we train a CNN $\hat{f}_{\mathbf{W}}(.)$ to predict two values, which can be seen as a scaled cosine and sine, and compute an angle by taking:

$$f_{\mathbf{W}}(\mathbf{p}_i^*) = \arctan2(\hat{f}_{\mathbf{W}}^{(1)}(\mathbf{p}_i^*), \hat{f}_{\mathbf{W}}^{(2)}(\mathbf{p}_i^*)) \quad , \quad (2)$$

where $\hat{f}_{\mathbf{W}}^{(1)}(\mathbf{p}_i^*)$ and $\hat{f}_{\mathbf{W}}^{(2)}(\mathbf{p}_i^*)$ are the two values returned by the CNN for patch $\mathbf{p}_i^*$, and $\arctan2(y,x)$ is the four-quadrant inverse tangent function[2].

This function is not defined at the origin, which turned out to be a problem only happening in rare occasions at the first iteration, after random initialization of the CNN parameters $\mathbf{W}$. To prevent this, we use the following approximation for its gradient:

$$\nabla\arctan2(y,x) = \left( \frac{-y}{x^2+y^2+\epsilon}, \frac{x}{x^2+y^2+\epsilon} \right) \ , \quad (3)$$

where $\epsilon$ is a very small value.

**Computing the derivatives.** The derivatives of the loss function for a given training pair $\mathbf{p}_i$ can be computed using the chain rule:

$$\frac{\partial \mathcal{L}_i}{\partial \mathbf{W}}(\mathbf{p}_i) = \frac{\partial \mathcal{L}_i}{\partial \mathbf{g}_{1,i}} \frac{\partial \mathbf{g}_{1,i}}{\partial \theta_{1,i}} \frac{\partial \theta_{1,i}}{\partial \mathbf{W}}(\mathbf{p}_i^1) + \frac{\partial \mathcal{L}_i}{\partial \mathbf{g}_{2,i}} \frac{\partial \mathbf{g}_{2,i}}{\partial \theta_{2,i}} \frac{\partial \theta_{2,i}}{\partial \mathbf{W}}(\mathbf{p}_i^2) \ , \quad (4)$$

with $\theta_{*,i} = f_{\mathbf{W}}(\mathbf{p}_i^*)$ and $\mathbf{g}_{*,i} = g(\mathbf{p}_i^*, \theta_{*,i})$.

Jacobians $\frac{\partial \mathcal{L}_i}{\partial \mathbf{g}_{*,i}}$ and $\frac{\partial \theta_{1,i}}{\partial \mathbf{W}}$ are straightforward to compute. $\frac{\partial \mathbf{g}_{*,i}}{\partial \theta_{*,i}}$ is not as easy, since $\mathbf{g}_{*,i}$ is the descriptor for patch $\mathbf{p}_i^*$ after rotation by an amount given by $\theta_{*,i}$. For example, in case of *SIFT*, the descriptor extraction process involves building histograms, which cannot be expressed as a differentiable function. Moreover, depending on the descriptor, pooling region for extracting the descriptor changes as a different orientation is provided.

We therefore use a numerical approximation of the gradients: when we form the training data we also compute the descriptors for many possible orientations, every 5 degrees in our current implementation. We can then efficiently compute the derivatives in $\frac{\partial \mathbf{g}_{*,i}}{\partial \theta_{*,i}}$ by numerical differentiation. Note that in case of descriptors that can be expressed in analytic form, for example learning based descriptors [33], we can also easily compute $\frac{\partial \mathbf{g}_{*,i}}{\partial \theta_{*,i}}$, instead of using numerical approximations.

To implement the CNN $\hat{f}_{\mathbf{W}}(.)$, we use three convolution layers with the ReLU activation function, each followed by a max-pooling layer, followed by two fully connected layers with GHH activation. We detail the GHH activation below. We also use dropout regularization [35] for better generalization. Implementation details are provided in Section 4.1.

### 3.3. Generalized Hinging Hyperplane Activation

To achieve state-of-the-art results with CNNs, we propose to use a new activation function in our network layers that works better for our problem than standard ones. This activation function is a generalization of the popular ReLU,

maxout [16], and the recent PReLU [17] activation functions based on Generalized Hinging Hyperplanes (GHH), which is a general form for continuous piece-wise linear functions [41]. As GHH activation function is more general, it has less restrictions in shape, and allows for more flexibility in what a single layer can learn. This activation function plays one of the key roles in our method for obtaining good orientations, as we will show in Section 4.3.

Mathematically, for a given layer output $\mathbf{y} = [\mathbf{y}_{1,1}, \mathbf{y}_{1,2}, \ldots, \mathbf{y}_{2,1}, \ldots, \mathbf{y}_{S,M}]$ before activation, we consider the following activation function:

$$o(\mathbf{y}) = \sum_{s \in \{1,2,\ldots,S\}} \delta_s \max_{m \in \{1,2,\ldots,M\}} \mathbf{y}_{s,m} \ , \quad (5)$$

where

$$\delta_s = \begin{cases} 1, & \text{if } s \text{ is odd} \\ -1, & \text{otherwise} \end{cases} , \quad (6)$$

and $S$ and $M$ are meta-parameters controlling the number of planar segments and thus the complexity of the function. When $S = 1$, Eq. (5) reduces to maxout activation [16], and when additionally $M = 2$ with $\mathbf{y}_{s,1} = 0$, the equation reduces to the ReLU activation function. Finally, when $S = 2$, $M = 2$, $\mathbf{y}_{s,1} = 0$ and $\mathbf{y}_{1,m} = -\alpha_m \mathbf{y}_{2,m}$, where $\alpha_m$ is a scalar variable, Eq. (5) is equivalent to the PReLU activation function proposed in [17].

Therefore, instead of having to choose a non-linear activation function, we also learn it under the constraint that it is piece-wise linear.

## 4. Results

In this section, we first introduce the datasets used for evaluation and the setup for training our regressor. We then demonstrate the effectiveness of our method by comparing the descriptor performances using the original and our learned orientations. We show that the best matching performance can be achieved with our learned orientations, outperforming state-of-the-art. We also demonstrate the performance gain obtained by using the GHH activation compared to other activation functions, and investigate influence of datasets on the descriptor performances. We finally show a Multi-View Stereo (MVS) application.[3]

### 4.1. Dataset, Training, and Evaluation Setup

**Dataset.** Fig. 3 shows example images from the datasets we use for evaluation and training. Note that our collection of data is not only composed of planar objects but also of 3D objects with self occlusions. We also have various imaging changes including changes in the camera pose. We use the *Oxford* dataset [26] for training, and the Edge Foci (*EF*) dataset [48], the *Webcam* dataset [40], the *Strecha*

---

[2]We follow the standard implementation for the C language for this function.

[3]Datasets and source code available at http://cvlab.epfl.ch/.

Figure 3. Selected images from all the datasets. **First row:** images from the *Oxford* dataset, **Second row:** images from the *Strecha* dataset, **Third row:** images from the *DTU* dataset, **Fourth row:** images from the *EF* dataset, **Fifth row:** images from the *Webcam* dataset, and **Last row:** images from our *Viewpoints* dataset. We use the *Oxford* dataset for training and the other datasets for testing.

dataset [36], the *DTU* dataset [1], and our own *Viewpoints* dataset for testing. Details on the datasets are as follows:

- *Oxford* dataset [26]: 8 sequences with 48 images in total. The dataset contains various imaging changes including viewpoint, rotation, blur, illumination, scale, JPEG compression changes. We thus use this dataset for training.

- *EF* dataset [48]: 5 sequences with 38 images in total. The dataset exhibits drastic lighting changes as well as daytime changes and viewpoint changes.

- *Webcam* dataset [40]: 6 sequences with 120 images in total. The dataset exhibits seasonal changes as well as daytime changes of scenes taken from far away.

- *Strecha* dataset [36]: composed of two sequences, fountain-P11 (11 images) and Herz-Jesu-P8 (8 images). The scene is non-planar and 3D. The dataset exhibits large viewpoint changes with self occlusions.

- *DTU* dataset [1]: 60 sequences with 600 images in total. This dataset also has multiple lighting settings for selected viewpoints, but we consider here only one lighting setting as we are mostly interested in the changes that occur on non-planar scenes undergoing camera movements. We also sample the viewpoints from the original dataset in regular intervals to make the dataset a manageable size.

- *Viewpoints* dataset: 5 sequences with 30 images in total. We created our own dataset to further enrich the dataset. The dataset exhibits large viewpoint changes and in-plane rotations up to 45 degrees from the reference image, which is when commercial cameras compensate the image orientation as landscape or portrait.

**Implementation details and training.** We use a patch size of $28 \times 28$ as input to the CNN. For the convolution layers, the first convolution layer uses a filter size of $5 \times 5$ and 10 output channels, the second convolution layer a filter size of $5 \times 5$ and 20 output channels, and the third convolution layer a filter size of $3 \times 3$ and 50 output channels. All max-pooling layers perform $2 \times 2$ max pooling. The size of the output of the first fully connected layer is 100, with the second fully connected layer having two outputs with the arctan2 mapping into orientations as described on Section 3.2.

For optimization, We use the ADAM [20] method with default parameters and exponentially decaying learning rate. We run 100 epochs with batch size of 10. The learning rate decay is set to half the learning rate every ten epochs.

Our method is also computationally efficient as we are only estimating the orientations. On an Intel Xeon E5-2680 2.5GHz Processor, our current implementation in Python with Theano [5] takes *0.47* milliseconds per feature point to compute orientations without any multi-threading. When used with *SIFT* descriptors, it overall takes *1.39* milliseconds per feature point to obtain the final descriptor. Note that the C++ implementation of the *MROGH* descriptor, which is the best performing rotation invariant descriptor in our experiments, takes *1.94* milliseconds per feature point.

**Evaluation methodology.** To demonstrate the effectiveness of our method, we compare the descriptor matching performances with our orientation estimation against other
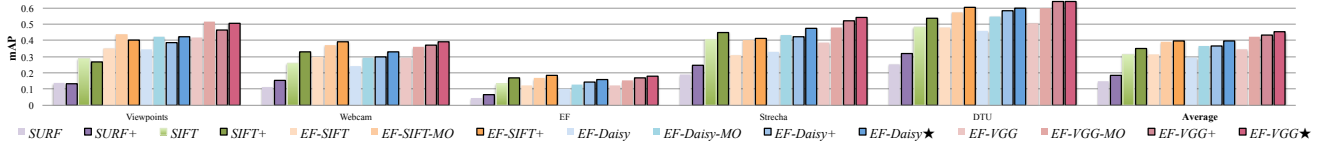
Figure 4. Descriptor performances with and without our orientation assignment. Methods with the multiple orientation strategy are denoted with MO, our orientation assignments learned with respective descriptors are denoted with a **+** at the end, and those learned with *SIFT* are denoted with a ⋆. Note that on average, methods with our orientation assignments perform better for all descriptors, than using orientation assignments from respective detectors and using multiple orientations.

state-of-the-art descriptors[4]. We use the standard precision-recall measure of [26] with nearest neighbor matching, and with a maximum of 1000 feature points per image. In case of the *DTU* and *Strecha* datasets, the scenes are non-planar and we rely on the 3D models and camera projection matrices to map a point from one viewpoint to another. Such mapping is used instead of the homography, followed by the overlap test in [26]. Results are summarized with the mean Average Precision (mAP) as in [34], where mAP is effectively the Area Under Curve of the precision-recall graph.

We compare against both descriptors that require orientation estimations (*ORB* [32], *BRISK* [22], *FREAK* [2], *SURF* [6], *SIFT* [25], *KAZE* [3], *BiCE* [47], *Daisy* [38], and the learning-based *VGG* [34]), as well as rotation invariant descriptors (*LIOP* [42], *MROGH* [14], and *sGLOH* [7]). Note that descriptors are generally designed for a specific detector (for example they typically have different range of scale of operation) and for fair comparisons, we do not interchange the detector and descriptors. We use the feature point detectors presented when the descriptors were introduced. In case of the *VGG* descriptor, we use the descriptor pre-learned with the *liberty* dataset [10], as other sequences are partially included in our test set. We employ the Edge Foci (*EF*) [48] detector for *VGG* as it showed better performance than the Difference of Gaussians (DoG) detector, which was used in the original work of *VGG*. We will denote this method as *EF-VGG*.

We also use *EF* detector with the *Daisy* descriptor and the *SIFT* descriptor, as this particular detector was designed with these two descriptors in mind. We will refer them as *EF-Daisy* and *EF-SIFT*, respectively.

### 4.2. Descriptor Matching Performances

To demonstrate the effectiveness of our method, we evaluate the descriptor matching performances with and without our orientation assignment. We first show the performance gain we obtain for *SIFT*, *SURF*, *Daisy*, and *VGG* descriptors and then compare our performance against other state-of-the-art methods. Note that for each descriptor, we only train our method *once* using the *Oxford* dataset and test on all the other datasets.

---

[4]Details on the implementations of these methods are provided as appendix in the supplementary material.

**Performance gain with our orientations.** To demonstrate the performance gain we obtain by using our orientation assignments, we learned orientations for *SIFT*, *SURF*, *Daisy*, and *VGG* descriptors. We denote descriptors computed using our learned orientation assignments with a **+** and a ⋆ at the end; we use **+** when orientations are learned with respective descriptors, and ⋆ when learned with *SIFT* descriptors. We also compare against using multiple dominant orientations. Note that using multiple orientations effectively amounts to creating duplicate feature points, which resulted in 34% increase in descriptor extraction time and 79% increase in matching time within our evaluation framework.

As shown in Fig. 4, we gain a consistent boost in descriptor matching performance with our orientation estimation. This includes the learning-based *VGG* descriptor, showing that learning-based methods also can benefit from a better orientation assignment. We also obtain a larger gain on average compared to using multiple orientations. The best performance is achieved with *EF-VGG⋆*.

Interestingly, for *EF-Daisy* and *EF-VGG*, learning with the *SIFT* descriptor gave larger boost in performances than learning with the respective descriptors. We suspect that this is due to the characteristics of the two descriptors being less sensitive to orientations than *SIFT* descriptors, resulting in the Jacobians with respect to orientations to vanish.

Based on the comparison results in Fig. 4, in the remainder of the results section we will report the performance of the best performing handcrafted descriptor with our orientations, *EF-Daisy⋆*, and the best performing learning-based descriptor with our orientations, *EF-VGG⋆*. In Section 4.3, as the best performance was achieved by learning with the *SIFT* descriptor (*EF-VGG⋆*), we will use *EF-SIFT***+** to evaluate the influence of different activation functions.

**Comparison with the state-of-the-art.** As shown in Fig. 5, both *EF-Daisy⋆* and *EF-VGG⋆* outperform all compared methods, with *EF-VGG⋆* outperforming all others by a large margin. Specifically, *EF-VGG⋆* performs *27.4%* better in terms of mAP compared to *EF-VGG*, which is the best performing competitor. Note that without our orientation estimation, although the best among the competitors, the gap is small. Also, as pointed out in descriptor performance surveys [1, 19, 28, 29], *SIFT* or *EF-SIFT* generally
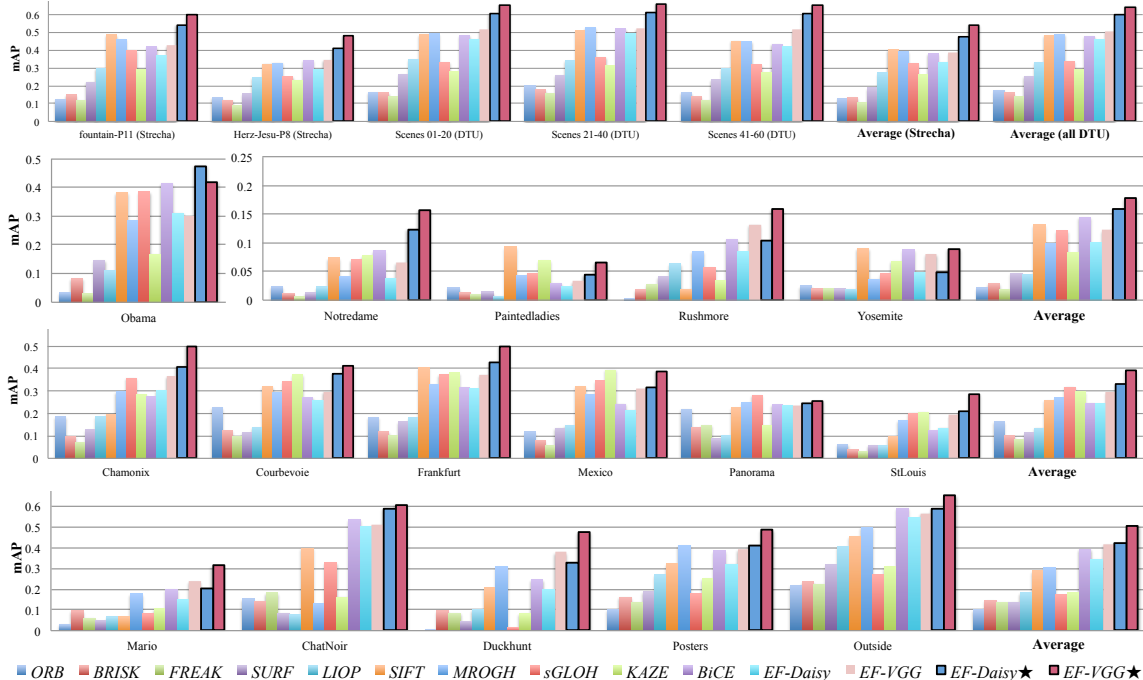
Figure 5. Mean Average Precision (mAP) for all datasets. **First row:** results for the *Strecha* dataset and the *DTU* dataset, **Second row:** results for the *EF* dataset, **Third row:** results for the *Webcam* dataset and **Last row:** results for the *Viewpoints* dataset. Best results are achieved with our orientation assignments, *EF-VGG⋆*.

|  | Viewpoints | Webcam | EF | Strecha | DTU | Rank | Avg. mAP |
|---|---|---|---|---|---|---|---|
| *ORB* | 13.20 | 10.33 | 11.40 | 12.50 | 12.33 | 12 | 0.12 |
| *BRISK* | 10.40 | 12.67 | 12.60 | 12.50 | 12.43 | 13 | 0.11 |
| *FREAK* | 11.20 | 13.33 | 12.60 | 14.00 | 13.45 | 14 | 0.10 |
| *SURF* | 11.40 | 12.17 | 10.80 | 11.00 | 10.57 | 11 | 0.15 |
| *LIOP* | 9.80 | 11.17 | 11.40 | 9.00 | 9.12 | 10 | 0.19 |
| *SIFT* | 7.00 | 6.33 | 4.80 | 4.50 | 5.33 | 5 | 0.31 |
| *MROGH* | 5.80 | 6.00 | 7.20 | 4.50 | 4.85 | 6 | 0.31 |
| *sGLOH* | 10.20 | 3.50 | 6.00 | 7.50 | 8.93 | 8 | 0.25 |
| *KAZE* | 9.20 | 4.83 | 6.00 | 10.00 | 9.43 | 9 | 0.22 |
| *BiCE* | 3.80 | 7.50 | 4.00 | 5.00 | 5.28 | 4 | 0.33 |
| *EF-Daisy* | 6.00 | 7.50 | 7.20 | 7.50 | 6.17 | 7 | 0.30 |
| *EF-VGG* | 3.20 | 5.50 | 5.40 | 4.00 | 3.87 | 3 | 0.35 |
| *EF-Daisy⋆* | 2.80 | 2.83 | 3.80 | 2.00 | 2.12 | 2 | 0.40 |
| *EF-VGG⋆* | **1.00** | **1.33** | **1.80** | **1.00** | **1.12** | **1** | **0.45** |

Table 1. Average rank of each method which summarizes the results in Fig. 5. Average rank for each dataset is given on the left, and the rank by averaging this value is provided on the right, as well as the mAP for all datasets. Bold denotes best performance. Our method *EF-VGG⋆* ranks first, followed by *EF-Daisy⋆*.

give comparable results to the state-of-the-art.

As average results can be influenced by certain sequences being too easy or hard, we also investigate the average rank of each method on the entire dataset similarly to [29]. In Table 1, we show the rank of each method on the datasets according to the average ranks of their mAP on each sequence. We also show the average mAP with all datasets for each method. Again, best results are obtained with our methods, *EF-Daisy⋆* and *EF-VGG⋆*.

### 4.3. Performance With Different Activations

To evaluate the influence of the proposed GHH activation function, we compared the matching performance of *EF-SIFT+* with different activation functions. All parameters were set to be identical except for the activation type and the number of outputs in the fully connected layers. Specifically, we used 1600 hidden nodes for ReLU, Tanh, and PReLU [17], 400 for maxout [16] with four outputs inside the max. Note that PReLU has slightly more parameters than other activations, as an additional parameter is introduced for each output of the layer.

As shown in Fig. 6, we have a consistent gain in performance when using the proposed GHH activation function instead of ReLU, Tanh, maxout, and PReLU. This shows that indeed using the GHH activation, which is a generalization of several common activation functions, is suitable for learning orientations.

### 4.4. Results on Upright and Non-upright Datasets

We observed in the existing datasets a general tendency for the images to be carefully taken with an upright posture. As a result, it is possible to assign a ground truth orientation to them by using a constant orientation. We will denote this upright assignment of orientations with the suffix "Up", and compare their assignments with our orientation assignments. We group the datasets into "upright"
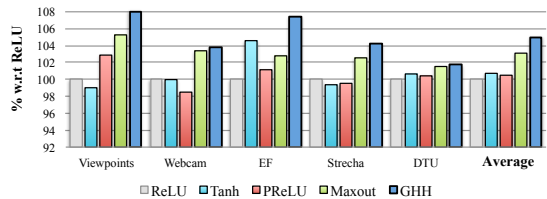
Figure 6. Relative descriptor performance obtained with the proposed GHH activation compared to ReLU, Tanh, PReLU and max-out activation functions. We use the mAP of ReLU activation as the reference (100%). Best performance is achieved with GHH activation.
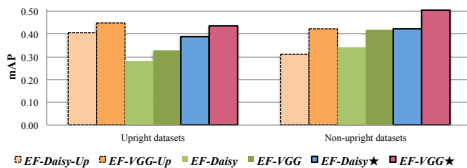


Figure 7. Performance of *EF-Daisy-Up*, *EF-VGG-Up*, *EF-Daisy*, *EF-VGG*, *EF-Daisy⋆*, and *EF-VGG⋆* on upright datasets and non-upright datasets. The "Up" suffix in the method names denotes that the feature points are assigned zero orientations. *EF-Daisy⋆* and *EF-VGG⋆* perform well on both cases, and are close to *EF-Daisy-Up*, and *EF-VGG-Up* on the upright datasets.

and "non-upright" ones, depending on whether the systematic assignment to an upright orientation performs better than using the original orientation assignments, and compare the performance of *EF-Daisy*, *EF-VGG*, *EF-Daisy⋆*, *EF-VGG⋆*, *EF-Daisy-Up*, and *EF-VGG-Up*.

Fig. 7 shows the results of these experiments. As expected, in case of upright datasets, using a systematic upright orientation performs the best, which can be seen as a upper bound for the descriptor performances. The performances however degrade when tested on non-upright datasets. However, our methods *EF-Daisy⋆* and *EF-VGG⋆* perform comparably to the upper bounds for upright datasets and are significantly better on non-upright datasets, achieving state-of-the-art. Note that *EF-VGG* performs similar to *EF-VGG-Up*, showing that inaccurate orientation assignments are not helpful.

### 4.5. Application to Multi-View Stereo

We also apply our orientation estimations for a MVS application [44, 45]. Fig. 8 shows MVS results using *EF-Daisy*, *EF-VGG* and *EF-VGG⋆*. Due to better matching performances, our method *EF-VGG⋆* gives best MVS results, followed by *EF-VGG*. Specifically, for the fountain sequence of the *Strecha* dataset, we obtain 260747 vertices with *EF-Daisy*, 323979 with *EF-VGG*, and 365261 with *EF-VGG⋆*. For Scene 55 of the *DTU* dataset, we obtain 72972, 86826, and 95014 vertices for *EF-Daisy*, *EF-VGG*, and *EF-VGG⋆*, respectively.



(a) *EF-Daisy* results (b) *EF-VGG* results (c) *EF-VGG⋆* results



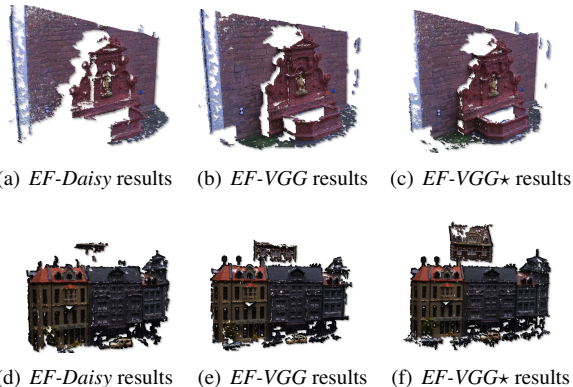(d) *EF-Daisy* results (e) *EF-VGG* results (f) *EF-VGG⋆* results

Figure 8. Multi-View Stereo (MVS) application [44, 45] example with *EF-Daisy*, *EF-VGG* and our method *EF-VGG⋆*. (a) – (c) results for the fountain sequence of the *Strecha* dataset, (d) – (f) results for the Scene 55 of the *DTU* dataset. Original images are shown as the first image for each dataset in Fig. 3. As better matches are provided with *EF-VGG⋆*, more detailed MVS reconstructions are obtained. Same *EF* feature points were used, differing only in orientation assignments.

## 5. Conclusion

We have introduced a learning scheme using a Convolutional Neural Network for the estimation of a canonical orientation for feature points, which improves the performance of existing descriptors. We proposed to train Siamese network to predict an orientation, which avoided the need of explicitly defining a "good" orientation to learn. We also proposed a new GHH activation function, which generalizes existing piece-wise linear activation functions and performs better for our task. We evaluated the effectiveness of our learned orientations by comparing the descriptor performances with and without our orientation assignment. Descriptors using our orientations gained consistent performance increase and outperformed state-of-the-art descriptors on all datasets. We finally investigated the influence of the GHH activation function showing its effectiveness.

Although we were able to enhance the performance of the learning-based *VGG* descriptor as well, an interesting future research direction is to fully integrate our method with learning-based descriptors, such as the recent descriptor presented in [33]. In which case, we can have a fully differentiable Siamese network which learns both the orientation assignment and the descriptor at the same time.

# References

[1] H. Aanæs, A. L. Dahl, and K. S. Pedersen. Interesting Interest Points. *IJCV*, 97:18–35, 2012. 1, 2, 4, 5, 6

[2] A. Alahi, R. Ortiz, and P. Vandergheynst. FREAK: Fast Retina Keypoint. In *CVPR*, 2012. 1, 2, 6, 11

[3] P. Alcantarilla, P. Fernández, A. Bartoli, and A. J. Davidson. KAZE Features. In *ECCV*, 2012. 1, 6, 11

[4] S. Allaire, J. J. Kim, S. L. Breen, D. A. Jaffray, and V. Pekar. Full Orientation Invariance and Improved Feature Selectivity of 3D SIFT with Application to Medical Image Analysis. In *CVPR*, 2008. 2

[5] F. Bastien, P. Lamblin, R. Pascanu, J. Bergstra, I. Goodfellow., A. Bergeron, N. Bouchard, and Y. Bengio. Theano: New Features and Speed Improvements. In *NIPS*, 2012. 5, 11

[6] H. Bay, T. Tuytelaars, and L. Van Gool. SURF: Speeded Up Robust Features. In *ECCV*, 2006. 1, 2, 6, 11

[7] F. Bellavia and D. Tegolo. Improving Sift-Based Descriptors Stability to Rotations. In *ICPR*, 2010. 1, 2, 6, 11

[8] F. Bellavia, D. Tegolo, and C. Valenti. Keypoint Descriptor Matching with Context-Based Orientation Estimation. *IVC*, 32(9):559–567, 2014. 2

[9] J. Bromley, I. Guyon, Y. LeCun, E. Sackinger, and R. Shah. Signature Verification Using a Siamese Time Delay Neural Network. In *NIPS*, 1993. 2

[10] M. Brown, G. Hua, and S. Winder. Discriminative Learning of Local Image Descriptors. *PAMI*, 2011. 3, 6

[11] M. Brown, R. Szeliski, and S. Winder. Multi-Image Matching Using Multi-Scale Oriented Patches. In *CVPR*, 2005. 2

[12] S. Chopra, R. Hadsell, and Y. LeCun. Learning a Similarity Metric Discriminatively, with Application to Face Verification. In *CVPR*, 2005. 2

[13] B. Fan, Q. Kong, T. Trzcinski, Z. Wang, C. Pan, and P. Fua. Receptive Fields Selection for Binary Feature Description. *TIP*, 23(6):2583–2595, 2014. 3

[14] B. Fan, F. Wu, and Z. Hu. Aggregating Gradient Distributions into Intensity Orders: A Novel Local Image Descriptor. In *CVPR*, 2011. 1, 2, 6, 11

[15] S. Gauglitz, M. Turk, and T. Höllerer. Improving Keypoint Orientation Assignment. In *BMVC*, 2011. 2

[16] I. Goodfellow, D. Warde-farley, M. Mirza, A. Courville, and Y. Bengio. Maxout Networks. In *ICML*, 2013. 2, 4, 7

[17] K. He, X. Zhang, R. Ren, and J. Sun. Delving Deep into Rectifiers: Surpassing Human-Level Performance on Imagenet Classification. In *ICCV*, 2015. 2, 4, 7

[18] S. Hinterstoisser, S. Benhimane, N. Navab, P. Fua, and V. Lepetit. Online Learning of Patch Perspective Rectification for Efficient Object Detection. In *CVPR*, 2008. 2

[19] N. Khan, B. Mccane, and S. Mills. Better Than SIFT? *MVA*, 26(6):819–836, 2015. 6

[20] D. Kingma and J. Ba. Adam: A Method for Stochastic Optimisation. In *ICLR*, 2015. 5

[21] S. Lazebnik, C. Schmid, and J. Ponce. Semi-Local Affine Parts for Object Recognition. In *BMVC*, 2004. 2

[22] S. Leutenegger, M. Chli, and R. Siegwart. BRISK: Binary Robust Invariant Scalable Keypoints. In *ICCV*, 2011. 1, 2, 6, 11

[23] W.-Y. Lin, L. Liu, Y. Matsushita, K.-L. Low, and S. Liu. Aligning Images in the Wild. In *CVPR*, 2012. 2

[24] K. Liu, H. Skibbe, T. Schmidt, T. Blein, K. Palme, T. Brox, and O. Ronneberger. Rotation-Invariant HOG Descriptors Using Fourier Analysis in Polar and Spherical Coordinates. *IJCV*, 106(3):342–364, 2014. 2

[25] D. Lowe. Distinctive Image Features from Scale-Invariant Keypoints. *IJCV*, 20(2), 2004. 1, 2, 3, 6, 11

[26] K. Mikolajczyk and C. Schmid. A Performance Evaluation of Local Descriptors. *PAMI*, 27(10):1615–1630, 2004. 2, 4, 5, 6

[27] K. Mikolajczyk and C. Schmid. Scale and Affine Invariant Interest Point Detectors. *IJCV*, 60:63–86, 2004. 1

[28] P. Moreels and P. Perona. Evaluation of Features Detectors and Descriptors Base on 3D Objects. In *IJCV*, 2006. 6

[29] D. Mukherjee, Q. M. J. Wu, and G. Wang. A Comparative Experimental Study of Image Feature Detectors and Descriptors. *MVA*, 26(4):443–466, 2015. 6, 7

[30] M. Osadchy, Y. LeCun, and M. Miller. Synergistic Face Detection and Pose Estimation with Energy-Based Models. *JMLR*, 8:1197–1215, 2007. 3

[31] H. Penedones, R. Collobert, F. Fleuret, and D. Grangier. Improving Object Classification using Pose Information . Technical report, Idiap Research Institute, 2011. 3

[32] E. Rublee, V. Rabaud, K. Konolidge, and G. Bradski. ORB: An Efficient Alternative to SIFT or SURF. In *ICCV*, 2011. 1, 2, 6, 11

[33] E. Simo-Serra, E. Trulls, L. Ferraz, I. Kokkinos, P. Fua, and F. Moreno-Noguer. Discriminative Learning of Deep Convolutional Feature Point Descriptors. In *ICCV*, 2015. 1, 2, 3, 4, 8

[34] K. Simonyan, A. Vedaldi, and A. Zisserman. Learning Local Feature Descriptors Using Convex Optimisation. *PAMI*, 2014. 2, 3, 6, 11

[35] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *JMLR*, 15:1929–1958, 2014. 4

[36] C. Strecha, W. Hansen, L. Van Gool, P. Fua, and U. Thoennessen. On Benchmarking Camera Calibration and Multi-View Stereo for High Resolution Imagery. In *CVPR*, 2008. 2, 4, 5

[37] S. Taylor and T. Drummond. Binary Histogrammed Intensity Patches for Efficient and Robust Matching. *IJCV*, 94(2):241–265, 2011. 2

[38] E. Tola, V. Lepetit, and P. Fua. Daisy: An Efficient Dense Descriptor Applied to Wide Baseline Stereo. *PAMI*, 32(5):815–830, 2010. 1, 6

[39] T. Trzcinski, M. Christoudias, and V. Lepetit. Learning Image Descriptors with Boosting. *PAMI*, 2015. 3

[40] Y. Verdie, K. M. Yi, P. Fua, and V. Lepetit. TILDE: A Temporally Invariant Learned DEtector. In *CVPR*, 2015. 2, 4, 5

[41] S. Wang and X. Sun. Generalization of Hinging Hyperplanes. *TIT*, 51(12):4425–4431, 2005. 2, 4

[42] Z. Wang, B. Fan, and F. Wu. Local Intensity Order Pattern for Feature Description. In *ICCV*, 2011. 1, 2, 6, 11

[43] S. Winder, G. Hua, and M. Brown. Picking the Best Daisy. In *CVPR*, June 2009. 1

[44] C. Wu. Towards Linear-Time Incremental Structure from Motion. In *3DV*, 2013. 1, 8

[45] C. Wu, S. Agarwal, B. Curless, and S. M. Seitz. Multicore Bundle Adjustment. In *CVPR*, 2011. 1, 8

[46] S. Zagoruyko and N. Komodakis. Learning to Compare Image Patches via Convolutional Neural Networks. In *CVPR*, 2015. 1, 2, 3

[47] C. Zitnick. Binary Coherent Edge Descriptors. In *ECCV*, 2010. 1, 6, 11

[48] C. Zitnick and K. Ramnath. Edge Foci Interest Points. In *ICCV*, 2011. 1, 2, 4, 5, 6, 11

# A. Supplementary Appendix

In this appendix, we provide details on the implementations used in the experiments.

## A.1. Implementations of the Methods

**Compared Methods** To keep the maximum number of features points to 1000, we sort the detected feature points according to their respective response scores and keep the best 1000. Details for the implementations of the compared methods are as follows:

- *ORB* [32]: OpenCV library – http://opencv.org/downloads.html
  We used nFeatures=1000, nLevels=3, and default values for other parameters.

- *BRISK* [22]: Provided by the authors – http://www.asl.ethz.ch/people/lestefan/personal/BRISK
  We used threshold of 20, with default values for other parameters.

- *FREAK* [2]: Provided by the authors – https://github.com/kikohs/freak
  Default parameters were used.

- *SURF* [6]: OpenCV library – http://opencv.org/downloads.html
  Default parameters were used.

- *LIOP* [42]: VLFeat library – http://www.vlfeat.org/
  Default parameters were used.

- *SIFT* [25]: OpenCV library – http://opencv.org/downloads.html
  Default parameters were used.

- *MROGH* [14]: Provided by the authors – https://github.com/bfan/MROGH-feature-descriptor
  Default parameters were used.

- *sGLOH* [7]: Provided by the authors – http://www.math.unipa.it/fbellavia/htm/research.html
  Default parameters were used.

- *KAZE* [3]: Provided by the authors – https://github.com/pablofdezalc/kaze
  Default parameters were used.

- *EF* [48] and *BiCE* [47]: Provided by the authors – http://research.microsoft.com/en-us/um/people/larryz/edgefoci/edge_foci.htm
  Default parameters were used.

- *Daisy* [47]: Provided by the authors – https://github.com/etola/libdaisy
  Patches were extracted to be four times the scale, which was the value authors used in [48]. Other parameters are set to default values.

- *VGG* [34]: Provided by the authors – http://www.robots.ox.ac.uk/~vgg/software/learn_desc/
  Patches were extracted with the VLFeat library, with a relativeExtent of 7.5, which is the same as what *SIFT* uses. We use the pre-learned model learned with the *liberty* dataset from [34], as the other two datasets are partially included in our test set.

**Our Methods**

We used the Python Theano library [5] for implementation – http://deeplearning.net/software/theano/