

# On Spatio-Temporal Feature Point Detection for Animated Meshes

**Abstract** Although automatic feature detection has been a long-sought subject by researchers in computer graphics and computer vision, feature extraction on deforming models remains as relatively unexplored area. In this paper, we develop a new method for automatic detection of spatio-temporal feature points on animated meshes. Our algorithm consists of three main parts. We first define local deformation characteristics, based on strain and curvature values computed for each point at each frame. Next, we construct multi-resolution space-time Gaussians and Difference-of-Gaussian (DoG) pyramids on the deformation characteristics representing the input animated mesh, where each level contains 3D smoothed and subsampled representation of the previous level. Finally, we estimate locations and scales of spatio-temporal feature points by using a scale-normalized differential operator. A new, precise approximation of spatio-temporal scale-normalized Laplacian has been introduced, based on the space-time DoG. We have experimentally verified our algorithm on a number of examples, and conclude that our technique allows to detect spatio- and temporal feature points in a reliable manner.

**Keywords** Feature detection · Animated mesh · Multi-scale representation, Difference of Gaussian

---

## 1 Introduction

With the increasing advances in animation techniques and the motion capture devices, animation data has become more and more available today. Coupled with this, almost all geometry processing techniques (alignment, reconstruction, indexing, compression, segmentation, etc.) began to evolve around the new, time-varying data, which is an active research area in Computer Graphics. Many applications in medicine and engineering benefit from the increased availability and usability of animation data.

Since such data has considerably large sizes, it often becomes indispensable to be able to select distinctive features from it, so as to maintain efficiency in its representation and in the process applied to it. Consequently, the need for robust, repeatable, and consistent detection of meaningful features from animation data cannot be overemphasized. However, the feature detection in animated mesh remains as much less explored domain, despite the proliferation of feature detectors developed by many researchers in computer graphics and computer vision.

In this paper, we develop a spatio-temporal feature detection framework on animated meshes (an ordered sequence of static mesh frames with fixed number of vertices and connectivity), based on the scale space approaches. Our algorithm, which we call AniM-DoG, extends the spatial IP (interest point)

detectors on static meshes[PKG03][CCF\*08][ZBV\*09][DK12] to animated meshes, so as to detect spatio-temporal feature points on them. Based on a deformation characteristic computed at each vertex in each frame, we build the scale space by computing various smoothed versions of the given animation data. At the heart of our algorithm is a new space-time Difference of Gaussian (DoG) operator, which is an approximation of the spatio-temporal, scale-normalized Laplacian. By computing the local extrema of the new operator in space-time and scale, we obtain repeatable sets of spatio-temporal feature points over different deforming surfaces modeled as triangle mesh animations. We then validate the proposed AniM-DoG algorithm for its robustness and consistency. To the best of our knowledge, our work is the first that addresses the spatio-temporal feature detector in animated meshes.

The remainder of the paper is organized as follows. In Section 2, we survey related works on local feature extraction in videos and (static) meshes. After recapitulating some basic terminologies and notions in Section 3, we present an overview of the method's pipeline overview in Section 4. Next, we describe the scale space representation and our AniM-DoG algorithm in Section 5. In Section 6 we show results of proposed feature point extraction algorithm and evaluate robustness of the method. Finally, we present some useful applications of the spatio-temporal feature detection in Section 7 and conclude in Section 8.

---

## 2 Previous works

Feature extraction is essential in different domains of computer graphics and is frequently used for numerous tasks including registration, object query, object recognition etc. Scale-space representation has been widely used for feature extraction in image, video and triangle mesh data sets [Lin98]. However, almost no research has been done on the feature extraction of deforming surfaces, such as animated meshes.

**Interest point detection in images and videos.** Perhaps one of the most popular algorithms of feature extraction on images is Harris-Stephens detector [HS88], which uses second moment matrix and its eigen-values to choose points of interest. However, Harris method is not invariant to scale. Lindeberg [Lin98] tackled that problem and introduced automatic scale selection technique, which allows feature point detection at their characteristic scales. As Lindeberg has shown, local scale estimation using the normalised Laplace operator allows robustly detect interest point of different extents. Mikołajczyk and Schmid [MS01] further developed Lindeberg's idea. As an improvement to the work of Lindeberg, authors proposed to use simultaneously Harris and Laplacian operators to detect interest points in scale-space

representation of an image. First, feature point candidates are detected as local maxima of Harris function in the image plane. Further, in order to obtain a more compact representation, only those points are retained where Laplacian reaches maxima over scale space. This approach, however, requires dense sampling over the scale parameters and is therefore computationally expensive. As Lowe [Low04] proposed, Difference of Gaussians (DoG) is a good approximation of Laplacian and hence could be used to reduce computational complexity.

More recently, Laptev and co-authors [LL03] investigated how the notion of scale-space could be generalized to the detection of feature points in space-time data such as image sequences or videos. Interest points are identified as simultaneous maxima of the spatio-temporal Harris corner function as well as extrema of the normalized spatio-temporal Laplace operator. In order to avoid computational burden authors proposed to capture interest points in only sparse scale pyramid and then track these points in spatio-temporal scale-time-space towards the extrema of scale-normalized Laplacian. However, in their method there is no guarantee of convergence. In the work of [BET\*08] a novel detector-descriptor scheme SURF (Speeded up robust features) has been proposed. Authors extend existing Hessian-based approaches and introduce Fast-Hessian detector that employs integral images for fast Hessian approximation.

**Feature description and feature point (FP) extraction on static meshes.** There have been several approaches proposed for detecting feature points on 3D meshes. Most of them extend the detectors proposed for images. Pauly *et al.* [PKG03] has used ‘surface variation’ to measure the saliency of vertices on the mesh, from which they build multi-scale representation. After extracting points with high feature response values, they construct minimum spanning tree of the edge points to extract feature lines.

Lee *et al.* [LVJ05] proposed algorithm to compute the saliency of mesh points based on the center-surround operator of Gaussian-weighted mean curvatures. First mean surface curvatures are computed. Then for each vertex, they estimate saliency as an absolute value of the difference between mean curvatures filtered with Gaussians of smaller and larger variances. They repeat the procedure at different scales by increasing Gaussian variance. Non-linearly normalized aggregate of saliency at all scales is defined as the final vertex saliency.

Castellani *et al.* [CCF\*08] build scale-space over vertices in a mesh with successive decimations of the original shape. The displacements of a vertex throughout the decimation are used as a measure of saliency. Then vertices with high response in its DoG operator (inter-octave local maxima), and with high saliency in the neighborhood (intra-octave local maxima) are selected as feature points.

Zaharescu *et al.* [ZBV\*09] use photometric properties associated with each vertex as a scalar function defined on a 3D mesh. A discrete operator named as ‘MeshDoG’ is applied on this function, on which they apply Hessian operator to detect corner-like feature points. They extend MeshDOG to what they call MeshHOG, a feature descriptor, which essentially is a histogram of gradients in the neighborhood. The extracted features along with their descriptors were used for matching 3D model sequences they obtained from multi-view images.

Sipiran and Bustos [SB10] have used 3D Harris operator which is essentially an extension of the Harris corner detector for images. After fitting quadratic patch to the neighborhood, a vertex is treated as an image, on which the Harris corner detector can be applied.

Darom and Keller [DK12] propose a scale-invariant local feature descriptor for the repeatable feature point extraction on 3D mesh. Each point is characterized by its coordinates, and a scale-space is built by successive smoothing of each vertex with its 1-ring neighbors. Local maxima both in scale and location are chosen as features.

All these methods, however, have been concerned with mesh data defined on the spatial domain only. In this work, we propose a new feature detection technique in animated meshes which extends existing methods based on linear scale-space theory to spatio-temporal domain.

### 3 Preliminaries

At the heart of our algorithm is a scale-space representation. In this section we briefly recapitulate some basic notions that have been previously studied. Later, we develop its extensions to animated mesh data, which are described in Section 5.2 and Section 5.3.

Scale-space representations have been studied extensively in feature detection for images, and more recently, for videos. The basic idea is to represent an original image  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  at different scales as  $L : \mathbb{R}^d \times \mathbb{R}_+ \rightarrow \mathbb{R}$  by convolution of  $f$  with a Gaussian kernel with variance  $\sigma$ :

$$L(\mathbf{x}; \sigma) = G(\mathbf{x}; \sigma) * f(\mathbf{x}), \quad (1)$$

where

$$G(\mathbf{x}; \sigma) = \frac{1}{(\sqrt{2\pi}\sigma)^d} \exp\left(-\frac{x_1^2 + \dots + x_d^2}{2\sigma^2}\right). \quad (2)$$

One of the most successful feature detectors is based on DoG (Difference of Gaussians). To efficiently detect feature points in scale space, Lowe [Low04] proposed using convolution of the input image with the DoG functions. It is computed from the difference of two nearby scales:

$$\begin{aligned} D(\mathbf{x}; \sigma) &= (G(\mathbf{x}; k\sigma) - G(\mathbf{x}; \sigma)) * f(\mathbf{x}) \\ &= L(\mathbf{x}; k\sigma) - L(\mathbf{x}; \sigma), \end{aligned} \quad (3)$$

where  $k$  is a constant multiplicative factor separating the two nearby scales. Note that DoG is particularly efficient to compute, as the smoothed images  $L$  need to be computed in any case for the scale space feature description, and  $D$  can therefore be computed simply by image subtraction.

The DoG provides a close approximation to the scale-normalized Laplacian of Gaussian [Lin94],  $\sigma^2 \nabla^2 G$ , which has been proven to produce the most stable, scale-invariant image features [MS01]. The DoG and scale-normalized LoG are related through the heat-diffusion equation:

$$\frac{\partial G(\mathbf{x})}{\partial \sigma} = \sigma \nabla^2 G(\mathbf{x}), \quad (4)$$

where the Laplacian on the right side is taken only with respect to the  $\mathbf{x}$  variables. From this, we see that  $\nabla^2 G(\mathbf{x})$  can be computed from the finite difference approximation to  $\partial G(\mathbf{x})/\partial \sigma$ , using the difference of nearby scales at  $k\sigma$  and  $\sigma$ :

$$\frac{\partial G(\mathbf{x})}{\partial \sigma} = \lim_{k \rightarrow 1} \frac{G(\mathbf{x}; k\sigma) - G(\mathbf{x}; \sigma)}{k\sigma - \sigma} = \sigma \cdot \nabla^2 G(\mathbf{x}), \quad (5)$$

and therefore,

$$G(\mathbf{x}; k\sigma) - G(\mathbf{x}; \sigma) \approx (k - 1) \cdot \sigma^2 \cdot \nabla^2 G. \quad (6)$$

## 4 Overview

Our goal is to develop a feature detector on animated mesh based on space-time DoG, which has been reported to be efficient approximation of robust Laplacian blob detector in space domain. Note that animated meshes that we are dealing with are assumed to have no clutters or holes, and maintain fixed topology over time, without tearing or changing genus. The spatial samplings can vary from one mesh to another, but it is desirable to have uniform sampling across one surface. The temporal sampling rate can also vary ( $\sim 30\text{Hz}$  in our experiments), depending on how the animation has been obtained. In any case, the temporal sampling is considered uniform.

The features we want to extract are the corners/blob-like structures, which are located in regions that exhibit a high variation of deformation spatially and temporally. We first define local deformation attributes on the animated mesh, from which we build a multi-scale representation of it. One of the main motivations to base our method on local surface deformation rather than vertex trajectories can be explained by the fact that (1) local deformation on a surface can be effectively measured by some well-defined principles, and that (2) the domain has intrinsic dimension of  $2D+time$  (rather than  $3D+time$ ) with some reasonable assumption on the data, i.e. differentiable 2-manifold with time-varying embedding.

We then compute the deformation characteristics at different scales, by defining an appropriate spatio-temporal Gaussian-like smoothing method. However, real Gaussian smoothing on mesh animation is problematic and expensive. Therefore we follow the other alternative and approximate Gaussian low-pass filter by a sequence of spatio-temporal box average filters of fixed width. We obtain different space and time scales of deformation field over animation by varying the number box filtering is applied in space and in time.

To estimate positions and scales of mesh animation feature points, we define a scale-normalized differential operator that assumes simultaneous extrema over space-time and scale neighborhood. Theoretically, it is possible to compute spatio-temporal, scale-normalized Laplacian on every vertex of the animated mesh. For example, one could extend the work by Zaharescu *et al.* and compute the 3D gradient and Laplacian on the animated mesh. However, it would be too much costly as it requires computing the normal plane, on which principal direction should be determined. Therefore, we introduce a new precise approximation of spatio-temporal scale-normalized Laplacian based on space-time Difference of Gaussian. Then local extrema of the space-time DoG operator are captured as feature points. Space-time DoG operator is cheap to compute and allows to robustly detect feature points over mesh animation in a repetitive and consistent manner.

## 5 Dynamic feature detector (AniM-DoG)

### 5.1 Deformation characteristics definition

We are interested in quantities that are related to local deformation characteristics associated to each point of the mesh, at each frame. Thus, we base our algorithm on locally computed strain and curvature values computed as follows.

**Strain computation.** We first consider the degree of deformation associated to each triangle on the mesh at each frame. Our method requires to specify the reference pose, the rest shape of the mesh before deformation (Fig. 1). In certain cases the reference pose can be found in one of the frames of the animation. If none of the given frames is appropriate as the rest pose, some prior works [LGX13] could be adopted to compute a canonical mesh frame by taking the average of all frames.

Let  $\mathbf{v}_i$  and  $\tilde{\mathbf{v}}_i$  be the vertices of a triangle before and after the deformation, respectively. A  $3$  by  $3$  affine matrix  $\mathbf{F}$  and displacement vector  $\mathbf{d}$  transforms  $\mathbf{v}_i$  into  $\tilde{\mathbf{v}}_i$  as follows.

$$\mathbf{F} \cdot \mathbf{v}_i + \mathbf{d} = \tilde{\mathbf{v}}_i, i = 1, \dots, 3$$

Similarly to Sumner *et al.* [SZG\*05], we add a fourth vertex in the direction of the normal vector of the triangle and subtract the first equation from the others to eliminate  $\mathbf{d}$ . Then, we get  $\mathbf{F} = \tilde{\mathbf{V}} \cdot \mathbf{V}^{-1}$  where

$$\mathbf{V} = [\mathbf{v}_2 - \mathbf{v}_1 \quad \mathbf{v}_3 - \mathbf{v}_1 \quad \mathbf{v}_4 - \mathbf{v}_1],$$

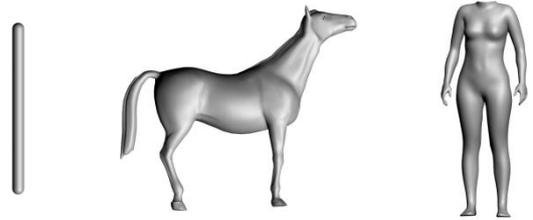
and

$$\tilde{\mathbf{V}} = [\tilde{\mathbf{v}}_2 - \tilde{\mathbf{v}}_1 \quad \tilde{\mathbf{v}}_3 - \tilde{\mathbf{v}}_1 \quad \tilde{\mathbf{v}}_4 - \tilde{\mathbf{v}}_1].$$

Non-translational component of  $\mathbf{F}$  encodes the change in orientation, scale, and skew induced by the deformation. Note that this representation specifies the deformation in per-triangle basis, so that it will be independent of the specific position and orientation of the mesh in world coordinates. Without loss of generality, we assume that the triangle is stretched first and then rotated. Then we have  $\mathbf{F} = \mathbf{R}\mathbf{U}$ , where  $\mathbf{R}$  denotes the rotation tensor and  $\mathbf{U}$  the stretch tensor. Since we want to describe the triangle only with its degree of stretch, we eliminate the rotation component of  $\mathbf{F}$  by computing the right Cauchy deformation tensor  $\mathbf{C}$  as defined by:

$$\mathbf{C} = \mathbf{F}^T \mathbf{F} = (\mathbf{R}\mathbf{U})^T (\mathbf{R}\mathbf{U}) = \mathbf{U}^T \mathbf{U}.$$

It can be shown that  $\mathbf{C}$  is equal to the square of the right stretch tensor. We obtain principal stretches by the Eigenanalysis on  $\mathbf{C}$ , and use the largest eigenvalue  $\lambda_1$  (maximum principal strain) as the in-plane deformation of the triangle.



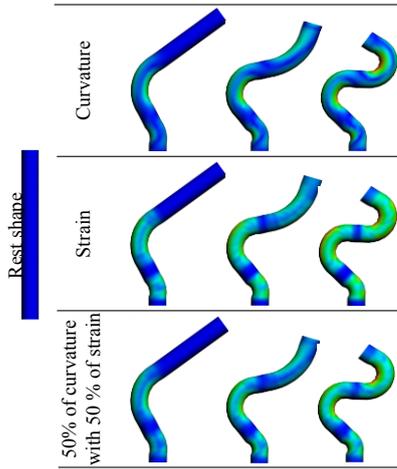
**Fig. 1** Rest shapes are chosen as the reference frame for defining the deformation characteristics.

**Curvature computation.** Computing the curvature at the vertices of a mesh is known to be non-trivial because of the piecewise-linear nature of meshes. One simple way of computing the curvature would be to compute the angle between two neighboring triangles along an edge. However, such curvature measurement is too sensitive to the noise on the surface of the mesh because its computation relies on two triangles only. Instead, we compute the curvature over a set of edges as described in [ACD\*03]. Given a vertex  $v_i$ , we first compute the set of edges  $E_i$  whose two vertices are within a user-defined geodesic distance to  $v_i$ . Next we compute the curvature at each of the edges of  $E_i$ . The curvature at  $v_i$  is then calculated as the average of the curvatures at the edges of  $E_i$ .

**Deformation measure.** Let  $\mathcal{M}$  with  $M$  frames and  $N$  vertices be a given deforming mesh. For each vertex  $v_i^f \in \mathcal{M}$  ( $f = 1, \dots, M, i = 1, \dots, N$ ) on which we have computed strain  $s(v_i^f)$  and curvature  $c(v_i^f)$ , we define the deformation characteristics  $d(v_i^f)$  as follows:

$$d(v_i^f) = s(v_i^f) + \alpha \cdot |c(v_i^f) - c(v_i^1)|.$$

The first term is obtained by transferring the above described per-triangle strain values to per-vertex ones, computed at each frame. At each vertex, we take the average strain values of its adjacent triangles as its strain.



**Fig. 2** Local deformation characteristics are shown on a bending cylinder mesh.

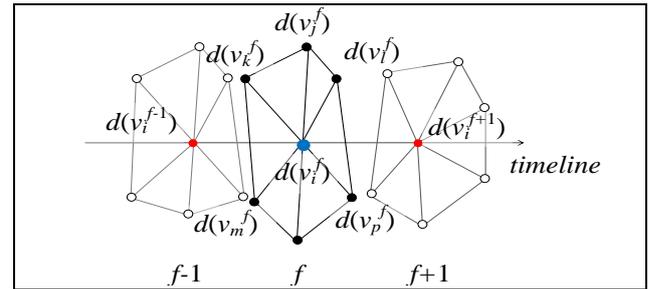
The second term encodes the curvature change with respect to the initial, reference frame. Note that  $d(v_i^f) \geq 0$  for  $\forall v_i^f$ , which we use later for the feature detection (Section 5.3). We set  $\alpha$  typically to 7 in our experiments. Color coded deformation characteristics on a bending cylinder data is shown in Fig. 2.

## 5.2 Scale space construction (Multi-scale representation)

Given the deformation measures  $d$  for all vertices of the input animated mesh  $\mathcal{M}$ , we re-compute  $d$  at K-L different scale representations, obtaining octaves  $O_{kl}$  ( $k \in \Sigma = 0, \dots, K, l \in T = 0, \dots, L$ ) of deformation characteristics at different spatio-temporal resolutions. Theoretically, the octaves are obtained by applying an approximated Gaussian filter for meshes. In practice, the approximation consists of subsequent convolutions of the given mesh with a box (average) filter [DK12]. In our work, we define a spatio-temporal average filter on the deformation characteristics of the animated mesh and compute a set of filtered deformation scalar fields, which we call as anim-octaves. As shown in Fig. 3, we define spatio-temporal neighborhood  $\mathcal{N}_{st}$  of a vertex in animation as a union of its spacial and temporal neighborhoods. A spatio-temporal average smoothing over  $\mathcal{N}_{st}$  is obtained by applying a local spatial filter followed by a local temporal one.

More specifically, for each vertex  $v_i^f$  at an anim-octave of scale  $(\sigma_k, \tau_l)$ , we compute deformation measures at next spatial octave  $(\sigma_{k+1}, \tau_l)$ , by averaging deformation measurements in current vertex of current octave  $d(v_i^f, \sigma_k, \tau_l)$  and its 1-ring's spatial neighborhood  $d(\mathcal{N}_s^1(v_i^f), \sigma_k, \tau_l)$  i.e. at adjacent vertices. For the next

temporal octave  $(\sigma_k, \tau_{l+1})$  we repeat similar procedure but this time averaging deformation values in 1-ring temporal neighborhood  $\mathcal{N}_t^1(v_i^f)$  as in Fig. 3. And for the next spatio-temporal octave, we start from deformations in octave  $(\sigma_{k+1}, \tau_l)$  and apply temporal average filter again in the way described above, which yields  $d(v_i^f, \sigma_{k+1}, \tau_{l+1})$ . We continue this procedure until we build the desired number of spatio-temporal octaves. Fig. 4 illustrates our anim-octaves structure. We denote an anim-octave as  $O_{kl} = d(\mathcal{M}, \sigma_k, \tau_l)$ , where  $O_{00} = d(\mathcal{M})$ . We note that although the term octave is widely used to refer to a discrete interval in the scale space, it may be misleading since in a strict sense, our octaves do not represent the interval of half or double the frequency. In Fig. 5, we illustrate multi-scale deformation characteristics we computed on an animated mesh. The horizontal axis represent the spatial scale  $\sigma_k$ , and the vertical axis the temporal scale  $\tau_l$ .



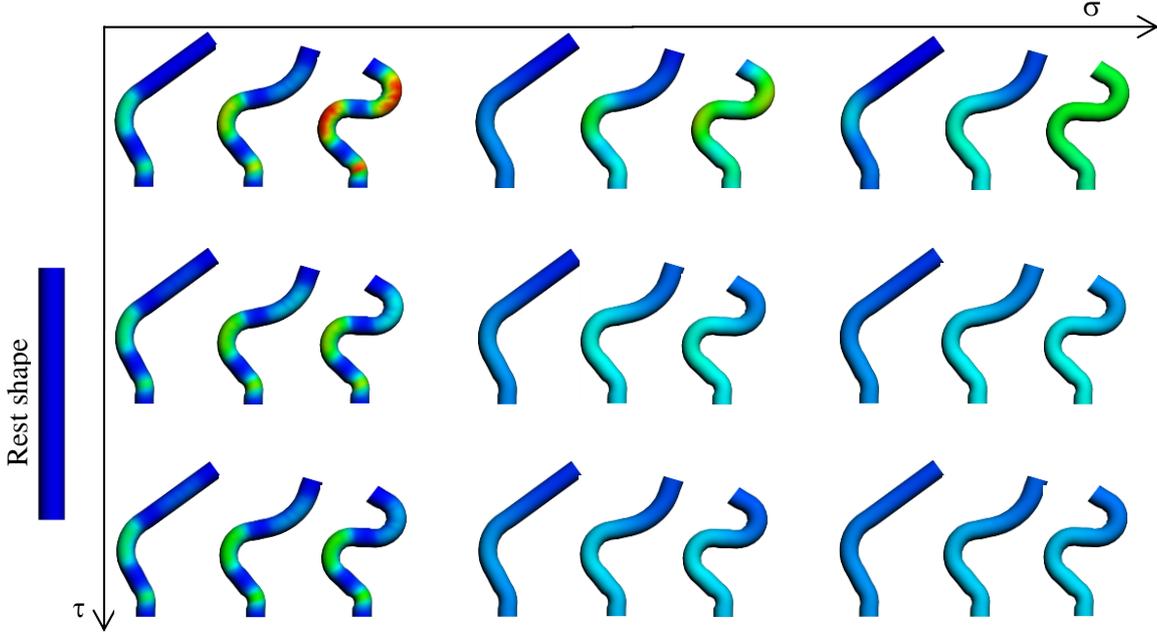
**Fig. 3** The smallest possible spatio-temporal neighborhood  $\mathcal{N}_{st}$  of a vertex  $v_i^f$  (blue dot) is composed of 1-ring spatial neighbors in frame  $f$  (black vertices) and 1-ring temporal neighbors (red vertices). Note that considering the temporal neighbors implies considering their spatial neighbors (white vertices) as well.

octave scale	$\sigma_1$	$\sigma_2$	...	$\sigma_k$
$\tau_1$	$O_{11}$	$O_{12}$	...	$O_{1k}$
$\tau_2$	$O_{21}$	$O_{22}$	...	$O_{2k}$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
$\tau_l$	$O_{l1}$	$O_{l2}$	...	$O_{lk}$

**Fig. 4** Scale space is built by computing a set of octaves of an input animated mesh.

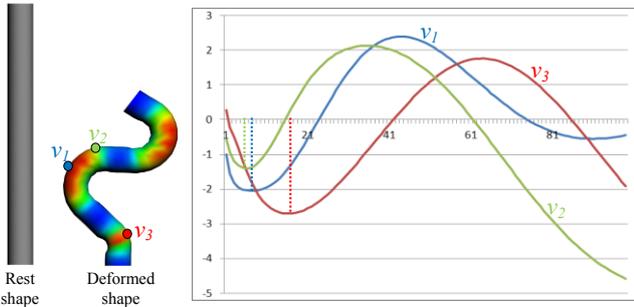
**Widths of the average filters.** We set the width of the spatial filter as the average edge length of the mesh taken at the initial frame, assuming that spatial sampling of the mesh is moderately regular, and that the edge lengths in the initial frame represent well those in other frames of animation. Note that it can be done in a per-vertex manner, by computing for each vertex the average distance to its 1-ring neighbors, as it has been proposed by Darom and Keller [DK12]. However, since this will drastically increase the computation time for the octave construction stage, we have chosen to use the same filter width for all vertices.

Determining the width of the temporal filter is simpler than the spatial one, as almost all data have regular temporal sampling rate (fps) throughout the duration of animation. Similarly to the spatial case, the inter-frame time is used to set the width of the temporal filter. Instead of averaging over immediate neighbors, however, we consider larger number of frame neighbors, in most cases. This is especially true when the animated mesh is densely sampled in time. The filter widths we used for each dataset are summarized in Table 1



**Fig. 5** Multi-scale deformation characteristics on an animated mesh. From left to right, spatial scale  $\sigma_j$  increases, and from top to bottom, temporal scale  $\tau_i$  increases.

**Maximum number of smoothings.** Since an animated mesh can be highly redundant and heavy in size, the memory space occupied by the anim-octaves can be large as the number of scales increases. This becomes problematic in practice. With an insufficient number of smoothings, on the other hand, features of large characteristic scale will not be detected. Indeed, when the variance of the Gaussian filter is not sufficiently large, only boundary features will be extracted. Fig. 6 illustrates the principle behind the characteristic scale and the maximum required scale level. Given a spatio-temporal location on the mesh, we can evaluate the DoG response function and plot the resulting value as a function of the scale (number of smoothings). Here, the spatial scale has been chosen as a parameter for the simplicity. The characteristic scales of the chosen vertices are shown as vertical lines, which can be determined by searching for scale-space extrema of the response function. To detect the feature points on the bending region in the middle ( $v_1$ ), for instance, the octaves should be built up to 12 level. This means the maximum number of smoothings must be carefully set in order to be able to extract feature points of all scales while maintaining moderate number of maximum smoothing.



**Fig. 6** The DoG response function has been evaluated as a function of the spatial scale (number of smoothings). The characteristic scales of the chosen vertices are shown as vertical lines.

In order to make sure that the features representing blobs of large scale are detected, we start by an average filter. Multiple

applications of a box (average) filter approximates a Gaussian filter [AC92]. More precisely,  $n$  averagings with a box filter of width  $w$  produce overall filtering effect equivalent to the Gaussian filter with a standard deviation of:

$$\sigma = \sqrt{\frac{n(w^2-1)}{12}}. \quad (8)$$

When the Laplacian of Gaussian is used for detecting blob centers (rather than boundaries), the Laplacian achieves a maximum response with

$$\sigma = \frac{r}{\sqrt{2}}, \quad (9)$$

where  $r$  is the radius of the blob we want to detect.

Now assuming that the maximum radius  $r_{\max}$  of the blob we want to detect is known, we can compute the required number of average smoothing that is sufficient to detect blob centers from Eq. (8) and Eq. (9):

$$r_{\max}^2 = \frac{n(w^2-1)}{6} \quad (10)$$

$$\Leftrightarrow n = \frac{6r^2}{w^2-1}. \quad (11)$$

The maximum number of application of box filter for each dataset is listed in Table 1.

**Maximum radius of all possible blobs.** Along the spatial scale space, we consider the average edge length of the initial shape as the width of the average filter  $w$ , as described above. For the maximum possible radius of a blob, we compare the axis-length change of the tight bounding box of the mesh during animation, with respect to its initial shape. The half of the largest change in axis-length is taken as  $r_{\max}$ .

Along the temporal scale space, we assume that the maximum radius  $r_{\max}$  of all possible blobs is the half of the total time of duration of animation. By fixing the maximum number of

smoothing to some moderate value, we obtain the desirable box filter width from Eq. (10) or Eq. (11).

### 5.3 FP detection by DoG

In this section, we extend the idea of scale representation in spatial domain to spatio-temporal domain and adopt it to the case of animated mesh. Next, we propose our feature point detector and discuss some of its implementation related issues.

**Spatio-temporal scale space principles.** Given time varying input signal  $f(\mathbf{x}, t)$ ,  $f: \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}$ , one could build its scale-space representation  $L(\mathbf{x}, t; \sigma, \tau)$  by convoluting  $f$  with anisotropic Gaussian

$$L(\mathbf{x}, t; \sigma, \tau): \mathbb{R}^d \times \mathbb{R} \times \mathbb{R}_+^2 \rightarrow \mathbb{R}.$$

The motivation behind the introduction of separate scale parameters in space  $\sigma$  and time  $\tau$  is that the space- and the time- extents of feature points are independent in general [LL03].

Alternatively, another useful formulation of spatio-temporal scale space was reported in the work of Salden *et al.* [SRV98]. The spatio-temporal scale space  $L(\mathbf{x}, t; \sigma, \tau)$  for a signal  $f(\mathbf{x}, t)$  could be defined as a solution of two diffusion equations

$$\frac{\partial L}{\partial \sigma} = \sum_i \frac{\partial^2 L}{\partial x^i \partial x^i}, \quad (12a)$$

$$\frac{\partial L}{\partial \tau} = \frac{\partial^2 L}{\partial t^2}, \quad (12b)$$

with an initial condition

$$\lim_{\sigma \rightarrow 0^+ \tau \rightarrow 0^+} L(\mathbf{x}, t; \sigma, \tau) = f(\mathbf{x}, t).$$

In our case, the input animated mesh  $\mathcal{M}$  can be considered as 2-manifold with time-varying embedding, i.e.  $m(u, v, t)$  in 3D Euclidean space. Measuring deformation scalar field  $d(\mathcal{M})$  in 2-manifold over space and time yields a 3D input signal of the form  $d(u, v, t)$ ,  $d: \mathbb{R}^2 \times \mathbb{R} \rightarrow \mathbb{R}$ , and its scale space of the form  $L_d(u, v, t; \sigma, \tau): \mathbb{R}^2 \times \mathbb{R} \times \mathbb{R}_+^2 \rightarrow \mathbb{R}$ . Given the scale space representation  $L_d(\mathbf{x}, t; \sigma, \tau)$  of the input animated mesh, we proceed with the construction of the DoG feature response pyramid, which we as describe below.

**Computing DoG pyramid.** To achieve the invariance in both space and time, we introduce a spatio-temporal DoG operator, which is a new contribution. Our idea is to combine the spatial and the temporal parts of Laplacian and Difference-of-Gaussians. Given the property of DoG (Eq.6) and Eq. 12a, Eq. 12b, we obtain the following:

$$D(\mathbf{x}; \sigma) = L(\mathbf{x}; k\sigma) - L(\mathbf{x}; \sigma) = \sum_i \frac{\partial^2 L}{\partial x^i \partial x^i}, \quad (13a)$$

$$D(t; \tau) = L(t; k\tau) - L(t; \tau) = \frac{\partial^2 L}{\partial t^2}. \quad (13b)$$

Then we propose to define the spatio-temporal Laplacian by adding (13a) and (13b):

$$D(\mathbf{x}; \sigma) + D(t; \tau) = \sum_i \frac{\partial^2 L}{\partial x^i \partial x^i} + \frac{\partial^2 L}{\partial t^2} = \nabla^2 L. \quad (14)$$

The new spatio-temporal Laplacian operator is just a sum of DoG in space scale and DoG in time scale, which is computationally very efficient.

In order to be able to extract features of all scales correctly, we need to scale-normalize the DoG response function. Choosing the exponent coefficients for the spatial Eq.13a (rightmost-hand term) and the temporal Eq.13b (rightmost-hand term) parts of Laplacian [LL03]:

$$\nabla_{norm}^2 L = \sigma^2 \tau^{1/2} \sum_i \frac{\partial^2 L}{\partial x^i \partial x^i} + \sigma \tau^{3/2} \frac{\partial^2 L}{\partial t^2} \quad (15)$$

Therefore, to achieve scale-normalized approximation of Laplacian through DoG, we multiply both sides of (13a) with  $\sigma^2 \tau^{1/2}$  and both sides of (13b) with  $\sigma \tau^{3/2}$  obtaining

$$\sigma^2 \tau^{1/2} D(\mathbf{x}; \sigma) = \sigma^2 \tau^{1/2} \sum_i \frac{\partial^2 L}{\partial x^i \partial x^i}, \quad (16a)$$

$$\sigma \tau^{3/2} D(t; \tau) = \sigma \tau^{3/2} \frac{\partial^2 L}{\partial t^2}. \quad (16b)$$

And from (16a-16b) we see that

$$\nabla_{norm}^2 L = \sigma^2 \tau^{1/2} D(\mathbf{x}; \sigma) + \sigma \tau^{3/2} D(t; \tau).$$

On the other hand, we can get a formulation of spatio-temporal DoG that approximates scale-normalized Laplacian

$$D_{st}(\mathbf{x}, t; \sigma, \tau) = \sigma^2 \tau^{1/2} D(\mathbf{x}; \sigma) + \sigma \tau^{3/2} D(t; \tau).$$

Thus, given the definition of spatio-temporal Difference of Gaussians we can compute feature response pyramid in the following way. For each vertex  $(u, v, t)$  in the animated mesh  $\mathcal{M}$ , and for every scale  $(\sigma_k, \tau_l) \in \Sigma \times T$  of the surface deformation pyramid we compute  $D_{st}(u, v, t; \sigma_k, \tau_l)$ .

**FP detection.** Once the spatio-temporal DoG pyramid  $\{D_{st}(u, v, t; \sigma_k, \tau_l) \mid (\sigma_k, \tau_l) \in \Sigma \times T\}$  is constructed, we extract feature points by identifying local extrema of the adjacent regions in the space, time, and scales. In contrast to Mikolajczyk and Schmid [MS01] who computes Harris and Laplacian operators, our method requires only DoG, which makes itself computationally efficient. This is particularly interesting for the animated mesh data which is generally much heavier than the image. Considering that our surface deformation function is always non-negative (and consequently its scale-space representation), it is worth mentioning that Laplacian of Gaussian and its DoG approximation reach local minima in centers of blobs. Such specific LoG behavior is illustrated in Fig. 7.

For each scale  $(\sigma_k, \tau_l) \in \Sigma \times T$  of 2D DoG pyramid, we first detect vertices in animation that are local minima in DoG response over their spatio-temporal neighborhood  $\mathcal{N}_{st}$ :

$$P_{kl} = \{p \in \mathcal{M} \mid \forall p_i \in \mathcal{N}_{st}(p), D_{kl}(p) < D_{kl}(p_k) \text{ and } D_{kl}(p) < \varepsilon_{st}\},$$

where  $\mathcal{N}_{st}(p)$  is a spatio-temporal neighbourhood of vertex  $p$  in the animation  $\mathcal{M}$  (Fig. 3).

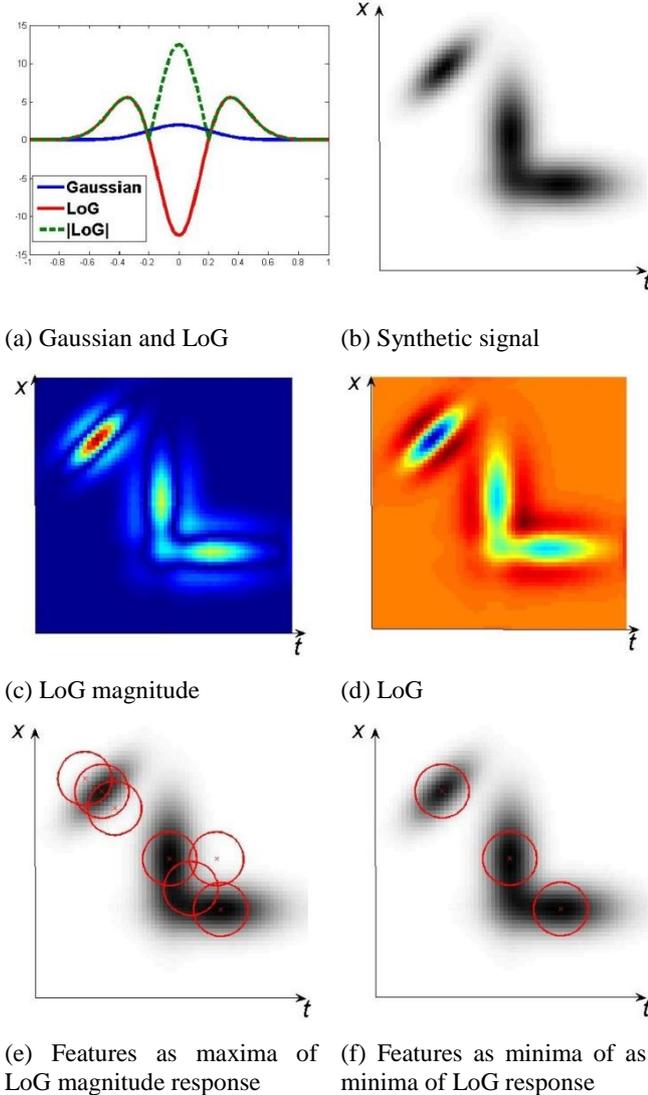
Next, out of preselected feature candidates  $P_{kl}$ , we retain only those vertices which are simultaneous minima over neighboring spatio-temporal scales of DoG pyramid :

$$P = \{p \in P_{kl} \mid \forall (i, j) \in \mathcal{N}_{\sigma\tau}(k, l), D_{ij}(p) > D_{kl}(p) \text{ and } D_{kl}(p) < \varepsilon_{\sigma\tau}\},$$

where  $\mathcal{N}_{\sigma\tau}(k, l)$  is a set of eight neighboring scales  $\mathcal{N}_{\sigma\tau}(k, l) = \{D_{(k+1)l}, D_{(k+1)(l-1)}, D_{(k+1)(l+1)}, D_{k(l-1)},$

$$D_{k(l+1)}, D_{(k-1)l}, D_{(k-1)(l-1)}, D_{(k-1)(l+1)},$$

and  $\varepsilon_{st}$ ,  $\varepsilon_{\sigma\tau}$  are user-controlled thresholds. The spatial scale of a feature point corresponds to the size of the neighborhood where some distinctive deformation is exhibited. Similarly, the temporal scale corresponds to the duration (or speed) of the deformation.



**Fig. 7** A 2D illustration of our feature detection method. (a) LoG yields the valley in blob's center and peaks around the boundary, while the magnitude of LoG has peaks in both cases. (b) Synthetic input signal consisting of 3 Gaussian blobs in 2d. (c) Response of synthetic 2d signal as the absolute value of LoG. (d) Response of the 2d signal computed as LoG. (e) Working with LoG magnitude response we observe several false secondary blobs. (f) Features captured as the local minima of LoG response are reliable.

**Dealing with secondary (border) blobs.** However, in case we consider local maxima of DoG (LoG) magnitude, we may detect artifacts. Undesirable secondary blobs are caused by the shape of Laplacian of Gaussian which yields peaks around the border of real blob (Fig. 7a). Consider a perfect Gaussian blob as an input signal. If we assume magnitude (i.e. absolute value) of LoG to be feature response, we get strong peak in the center of the blob and two other secondary peaks around edges, and that is troublesome. In contrast, dealing with signed LoG (not absolute) we observe valleys (local minima) in blob centers and peaks (local maxima) on borders. Hence searching for

local minima of LoG, rather than local maxima of LoG magnitude, prevents from the detection of false secondary features (Fig. 7b-f). The other way around could be to use LoG magnitude but discard local maxima which are not strong enough in initial signal, and therefore are false findings. Though, previous works in feature extraction on images/video/static meshes [MS01, LL03, ZBV\*09] often adopt Hessian detector, which does not detect secondary blobs. However, in contrast to DoG detector, estimation of Hessian on a mesh surface is significantly heavier. And even more problematic and challenging to estimate Hessian matrix on animated mesh.

**Implementation notes.** Often, animated meshes are rather heavy data. As we increase the number of anim-octaves in the pyramid, we can easily run out of memory, since each octave is essentially a full animation itself but at different scale. Consequently, we have to address that issue in the implementation stage. In order to minimize memory footprint, we compute pyramids and detect feature points progressively. We fully load into main memory only space scale dimension of Gaussian and DoG pyramids. As for time scale, we keep only two neighboring time octaves simultaneously, which are required for DoG computation. Then we construct the pyramid from bottom to top by iteratively increasing time scale. On each iteration of Gaussian/DoG pyramid movement along time scale, we apply our feature detection method to capture interest points (if any on current layer). We repeat the procedure until all pyramid scales have been processed.

## 6 Experiments

Deforming meshes used in our experiments include both synthetic animations and motion capture sequences, which are summarized in Table 1. We synthesized a simple deforming Cylinder animation by rigging and used it for initial tests. Also we captured two person's facial expressions using Vicon system [Vicon] and then transferred the expressions to the scanned faces of the two persons Face1 and Face2 (Table 1).

**Table 1** The animated meshes used in our experiments.

Name	No.vertices/ No. triangles	No. frames	Filter widths (space/time)	Max. no. smoothings (space/time)
Cylinder	587/1170	40	10.0/0.83	50/100
Face1(happy)	608/1174	139	8.96/8.45	118/113
Face1(surprise)	608/1174	169	9.39/13.2	96/107
Face2(happy)	662/1272	159	9.31/13.2	112/94
Face2(surprise)	662/1272	99	8.95/8.45	109/57
Horse	5000/9984	48	3.48/5.33	77/54
Camel	4999/10000	48	2.62/5.33	102/54
Woman1	4250/8476	100	5.12/5.2	72/150
Woman2	4250/8476	100	4.44/5.2	82/150
Woman3	4250/8476	100	4.54/5.2	99/150
Face3	5192/9999	71	5.18/3.65	90/114
Head	7966/15809	71	9.06/3.65	28/114

Fig. 10 shows selected frames of several animated meshes we used in our experiments. Spatio-temporal feature points we have extracted using our algorithm are illustrated as spheres. For the complete sequences along with the extracted feature points, please take a look at our accompanying demo video.

Face1 (happy), Face1 (surprise), Face2 (happy), Face2 (surprise) contain facial expressions of happiness and surprise of those scanned subjects. The horse and the camel were obtained from results of Sumner and Jovan Popović’s work [SP04] that are available online [Mesh]. Furthermore, we produced two mesh animations of seventy frames, Face3 and Head, by using nine facial expressions of the face and the head from [SP04]. More specifically, given an ordered sequence of nine facial expressions, we smoothly morphed each mesh to the next one through a linear interpolation of their vertex coordinates. We also used “walk and whirl” skeletal animations of three women models. Those models share the same mesh topology and were obtained by deforming a template mesh onto three body scans of different subjects. Note that there is high semantic similarity between animation pairs of Face1/Face2, horse/camel, Face3/Head. It is also the case for three women models.

The color of a sphere represents the temporal scale (red color corresponds to more fast deformations) of the feature point, and radius of sphere indicates the spatial scale. Vertex color on surfaces corresponds to amount of deformation (strain and curvature change) observed in each of animation frame. During experiments we have discovered that our method captures spatio-temporal scales in a robust manner. For example, surface patches around joints of cylinder (Fig. 11 1a-1e) exhibit different amount of deformation that occurs at different speed. The top joint is moving fast and consequently corresponding feature was detected at low temporal scale (red color). However, the mid-joint is deforming for a long time and we identify it at high temporal scale (blue color). Moreover large radii of deforming spheres for both joints make sense and indicate large deforming regions around the features, rather than very local deformation (Fig. 11 1c). Second row in (Fig. 11 2a-2e) depicts some of the feature points in Horse mesh animation, and third row (Fig. 8 3a-3e) corresponds to Camel animation. Those two meshes deform in a coherent manner [SP04], and eventually we detect their spatio-temporal features quite consistently. In the last two rows (Fig. 11 4a-4e, 5a-5e) we present feature points in mocap-driven face animations of two different subjects. Our subjects were asked to mimic of slightly exaggerated emotions during the mocap session. Notice that people normally use different set of muscles when they show up facial expressions, and therefore naturally we observe some variations in the way their skin deforms.

Our algorithm is implemented in C++. All our tests have been conducted on an Intel Core i7–2600 3.4 GHz machine, with 8GB of RAM. The computation time devoted to full pipeline of the algorithm is approximately 2 minutes for most of our example data.

**Invariance to rotation and scale.** Invariance of our detector to rotation as well as scale is evident from the definition of our deformation characteristics. Both the strain and the curvature measure we use are invariant to rotation and scale of the animated mesh.

**Robustness to changes in spatial and temporal sampling.** Robustness of our feature detector to changes in spatial sampling is obtained by the adaptive setting of the widths of the box filters. As described in Section 5.2, we set the width of the spatial filter as the average edge length of the mesh taken at the initial frame. In order to demonstrate the invariance to spatial density of the input mesh, we have conducted comparative experiments on two bending cylinders. These two cylinders have identical shape and deformation; they only differ by the number of vertices and the inter-vertex distance.

As shown in the 1st and 3rd rows of Fig. 10, the features are extracted at the same spatio-temporal locations.

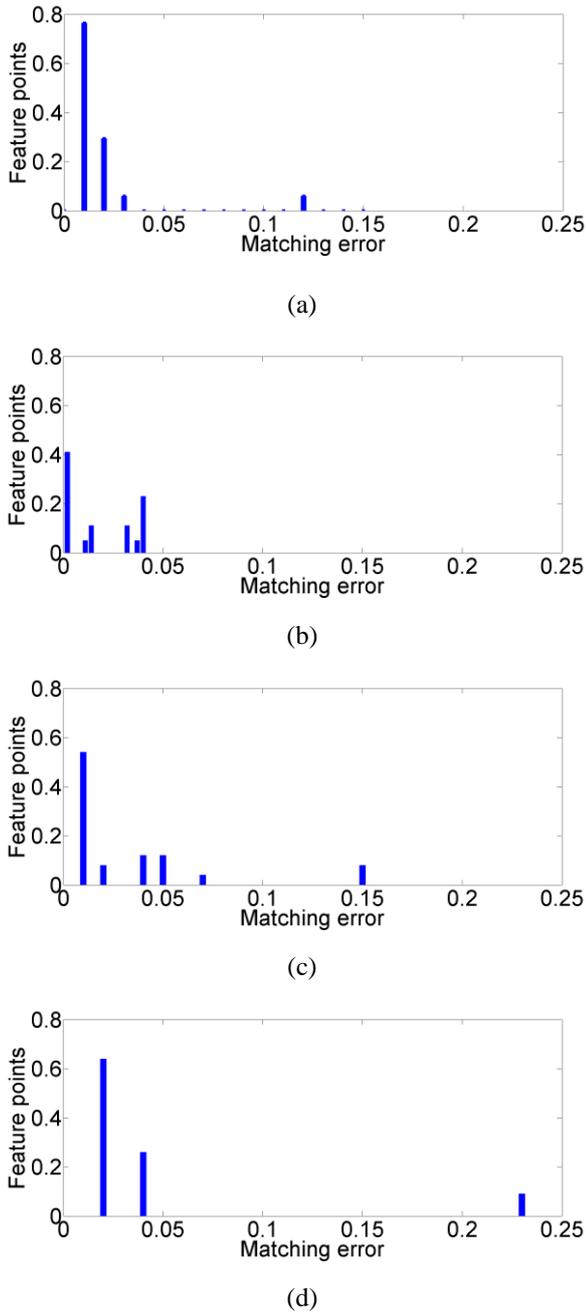
Robustness to changes in temporal sampling is obtained similarly to the above, i.e. by the adaptive setting of the widths of the box filters. Similar experiments have been conducted by using the two bending cylinders as shown in the 1st and 2nd rows of Fig. 10. They are perfectly identical except that the temporal sampling of the first one is twice higher than that of the second one. Once again, the extracted feature points are identical in their locations in space and time.

We have further experimented with datasets of similar animations, but with different shape, spatial- and temporal-samplings (4th row of Fig. 11, galloping animals and two face models in Fig. 11). Although the extracted features show a good level consistency, they are not always identical. For example, feature points for the galloping horse and camel do not have the same properties (location, time, tau and sigma). Similar results have been observed for the “face” models. This can be explained by the following facts: Firstly, although the two meshes have deformations that are semantically identical, the level of deformation (curvature and strain) might differ greatly. Secondly, most of these models have irregular vertex sampling whereas in our computation of the spatial filter width, we assume that the vertex sampling is regular.

## 6.1 Consistency

Since our method is based on deformation characteristics, it has an advantage of consistent feature point extraction across mesh animations with similar motions. To demonstrate mutual consistency among feature points in different animations, we used animation data that exhibit semantically similar motions. Our technique captures similarity of surface deformations and therefore ensures feature point detection consistency (Fig. 12). In most cases our method demonstrates high coherency not only in space and time locations of extracted features but also in their space-time scales  $\sigma$  and  $\tau$ . The only data sets for which we observed relatively lower consistency of feature detection are the two face mocap sequences. The reason for this lies in inherent difference of people’s facial expressions and underlying muscle anatomy.

Additionally, we have performed the quantitative evaluation of the feature extraction consistency as follows. For all feature points we consider only their spatial locations disregarding the time coordinates. Then a pair of similarly deforming meshes  $\mathcal{M}_S$  and  $\mathcal{M}_T$  whose full correspondence  $f: \mathcal{M}_S \rightarrow \mathcal{M}_T$  is known, we find the matching between their feature points  $F_S$  and  $F_T$  based on the spatial proximity. More precisely, for each feature point  $p_S^i \in F_S$ , the feature point  $p_T^j \in F_T$  that minimizes  $d_T^i = \|f(p_S^i) - p_T^j\|$  is considered to be the matching one. The distance  $d_T^i$  is what we call feature matching error. Histogram plots of feature matching errors are depicted in Fig. 8. Obtaining the full correspondence for walking women models was straightforward because they share the same mesh topology. For horse/camel, we obtained a full per-triangle correspondence from [SP04], which we converted to a per-vertex correspondence.



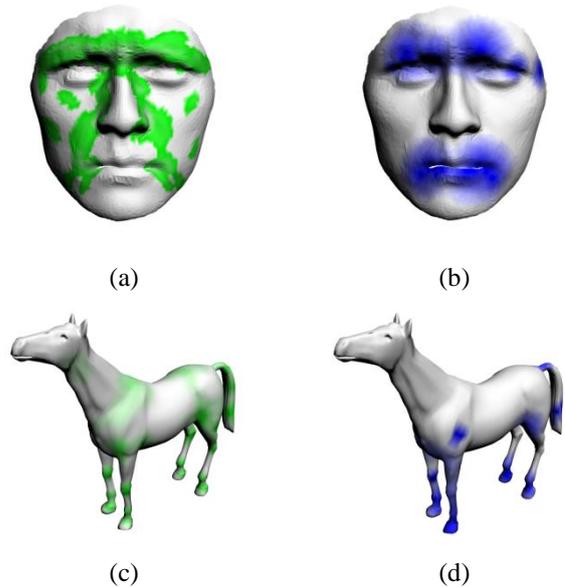
**Fig. 8** The error plots of feature points for pairs (a) Woman1-Woman2, (b) Woman1 – Woman3, (c) Woman2 - Woman3, Camel-Horse (d). We depict the feature matching error on the x-axis as the error (percentage of the error with respect to the diagonal length of the mesh bounding box). The percentage of features with prescribed matching error is depicted on the y-axis. For all four pairs of animations, more than 90% of features have a matching error less than 0.05.

## 6.2 Comparison to the ground truth

We have validated our method by comparing the feature points to the manually defined ground truth. We asked six volunteers to paint feature regions on the animated meshes using an interactive tool. The task was to mark locations at which salient surface deformation behavior can be observed from the user point of view. Each of them could play and pause the animation at any moment and mark feature regions by a color. To simplify the task, the time duration of each

feature region was not considered. Since the per-vertex selection can be error-prone, we deliberately allow users to select a region on the surface instead of a single vertex. By aggregating the feature regions from all volunteers, we generated a color map of feature regions. More specifically, for each vertex we summed up and averaged the number of times it has been included in the user-selected regions. The aggregated ground truth was then converted into a surface color map, as depicted in Fig. 9a, c. Note, that eyes do not belong to feature regions of face animations, since the user’s task was to define features based on the surface deformation rather than geometric saliency or naturally eye-catching regions.

To compare our results with respect to the ground truth, we compute for every feature point  $p$  its feature region of neighboring vertices  $q$  such that  $\mathcal{L} = \{q : d_{\mathcal{M}}(q, p) < \sigma_k\}$ , where  $d_{\mathcal{M}}(\cdot, \cdot)$  is a within-surface geodesic distance and  $\sigma_k$  is the corresponding scale value at which the feature was detected. Similarly to the ground truth, for each vertex of the mesh we count the number of occurrences in feature regions during the animation, and convert the numbers to the surface color map as shown in Fig. 9b, d. We observe a good level of positive correlation between the computed feature regions and the ground truth.



**Fig. 9** Comparison of the ground truth (a, c) to the feature regions computed by our method (b, d). For each vertex, color intensity depicts the accumulative number of its appearances during the animation. Green and blue colors were used for the ground truth and our computed feature regions, respectively.

## 7 Discussion

Our method and results could be extended and applied to a number of useful applications. We describe some of the ideas below while leaving their developments as future works.

**Animated mesh simplification.** As it has been noted in earlier works on simplification of dynamic (deforming) meshes [SPB00][KG05], it is preferable to allocate bigger triangle budget for regions of high surface deformation while simplifying mostly rigid regions. Our algorithm could be adopted in these as it detects feature points that are exactly in deforming regions. Their spatial scales  $\sigma_k$  can be used to

define regions around features where the mesh must keep denser sampling during simplification. For instance, the spatial scale of the feature points can be used to define regions where the mesh must be densely sampled during simplification. The temporal scale can also be used to dynamically determine the triangle budget around the feature point, when designing a time-varying simplification technique. A very small temporal scale implies either a short duration or a high speed of the animation, thus one may assign low priority to the feature point. In the same way, the region around a feature point with large temporal scale will be prioritized when allocating the triangle budget. Another use of the temporal scale is in the maintenance of the hierarchy. When transferring the previous frame's hierarchy to one better suited for the current frame in a time-critical fashion, the algorithm can use the temporal scale of a FP as a “counter” to determine whether to update or reuse the node corresponding to the region around the FP. By processing the nodes corresponding to the spatio-temporal feature points in an order of decreasing temporal scale, one can economize the time for the per-frame maintenance of the hierarchy while keeping the animation quality as much as possible.

**Viewpoint selection.** With increasing advances in scanning and motion capture technologies animated mesh data becomes more and more available today. Thus it is very practical to have a tool for automatic viewpoint selection for the preview of the motion in animation repositories. The idea behind that is to let a user to quickly browse the animation data from the point that maximizes the visibility of mesh deformations. With such viewpoint selection, the user benefits from a better perception of the animation. One equally handy and straight forward way to automatically select optimal viewpoint is to compute the one which maximizes the number of visible feature points through the optimization. We note that our spatio-temporal feature points can simplify the selection of good viewpoint(s). For instance, the quality of a viewpoint could be defined as a function of the visibility of the spatio-temporal feature points in terms of the total number, temporal variability and the concavity of the projected feature region (as defined by the spatial and temporal scales), etc. Interested reader may refer to an optimization technique proposed in [LVJ05] on saliency-based viewpoint selection for static meshes.

**Animation alignment.** Another interesting application could be animated mesh alignment. Considering the consistency of the extracted feature points, their scales values can be employed for the temporal alignment. Given sets of features  $P$  and  $P'$  extracted from pair of similar animations, we consider corresponding sequences  $\{(\sigma_1, \tau_1), \dots, (\sigma_n, \tau_n)\}$ ,  $\{(\sigma'_1, \tau'_1), \dots, (\sigma'_m, \tau'_m)\}$  of spatio-temporal feature scales

aligned along the time they were detected. Existing algorithms of sequence alignment such as [Got90] can then be used to compute the temporal alignment between them. In addition to the spatial- and temporal- scales, more sophisticated feature descriptors can also be used to compose the sequences.

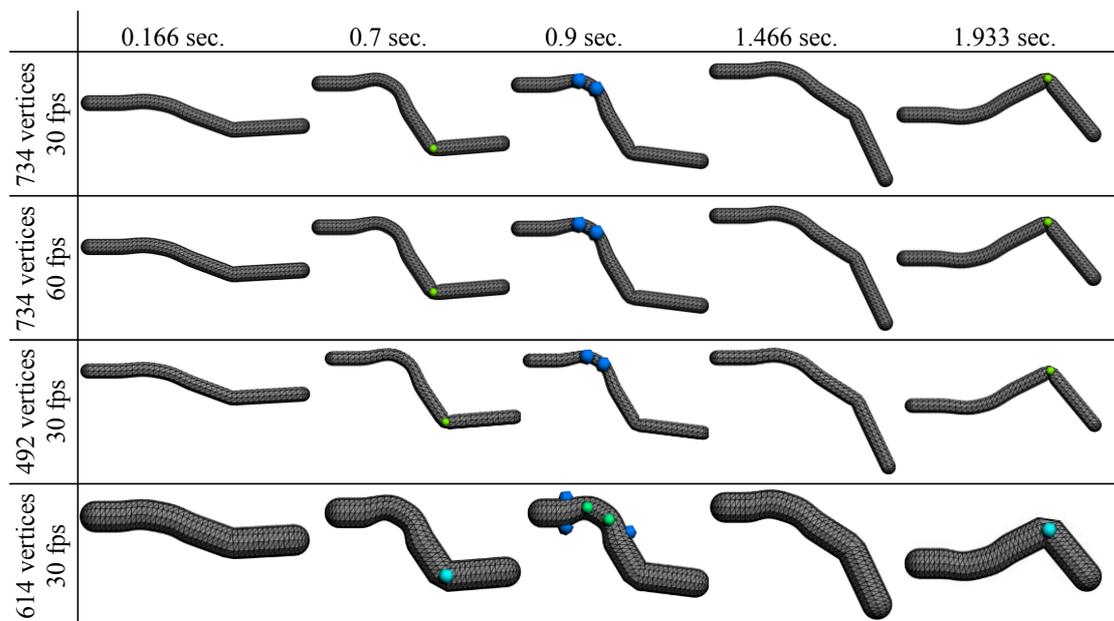
**Animation similarity.** We can also think of extending the above mentioned animation alignment algorithm towards a measurement of animation similarity. From the feature sequence alignment map, we can sum up all penalty gaps i.e. some predefined costs for all features for which no match can be found. That cost function could serve as a distance metric between the animations and hence be a measure of dissimilarity/similarity. Note that an important byproduct of the animation similarity is the animated mesh retrieval, which is particularly beneficial in emerging dynamic data repositories.

---

## 8 Conclusion

We have presented a new feature detection technique on triangle mesh animations based on linear scale-space theory. We introduced a new spatio-temporal scale representation of surface deformation in mesh animations. Furthermore, we developed extension of classical DoG filter to spatio-temporal case. The later allows our method to robustly extract repeatable sets of feature points over different deforming surfaces modeled as triangle mesh animations. We carried out experimental validation of detected features on various types of data sets and observed consistent results. Our approach has shown robustness to spatial and temporal sampling of mesh animation. In our future research, we intend to focus on feature point descriptor that could be useful for applications such as matching between animations.

**Descriptors.** Our feature detector could be extended to detector-descriptor. One straightforward idea of the feature point descriptor could be the following. Given space-time neighborhood around feature point  $v_i^f$  which consists of its  $k$ -ring space neighborhood over range of  $[f-l, \dots, f+l]$  frames. Note, that values of  $k$  and  $l$  could be adjusted to reflect characteristic scale of the feature. Next, flattening of  $k$ -ring regions for each frame of the range produces a volumetric stack of planar mesh patches. Then, proceeding in a spirit of robust 3D SIFT descriptor [SAS07], we can estimate histograms of DoG gradients computed inside the spatio-temporal volume around the feature point. Further, Euclidean or Earth Mover's distance between the histograms could be used for measuring similarities of features within mesh animation or between different animations.



**Fig. 10** Results we obtained on varying datasets of bending cylinder animations demonstrate the consistent behavior of our feature detector.

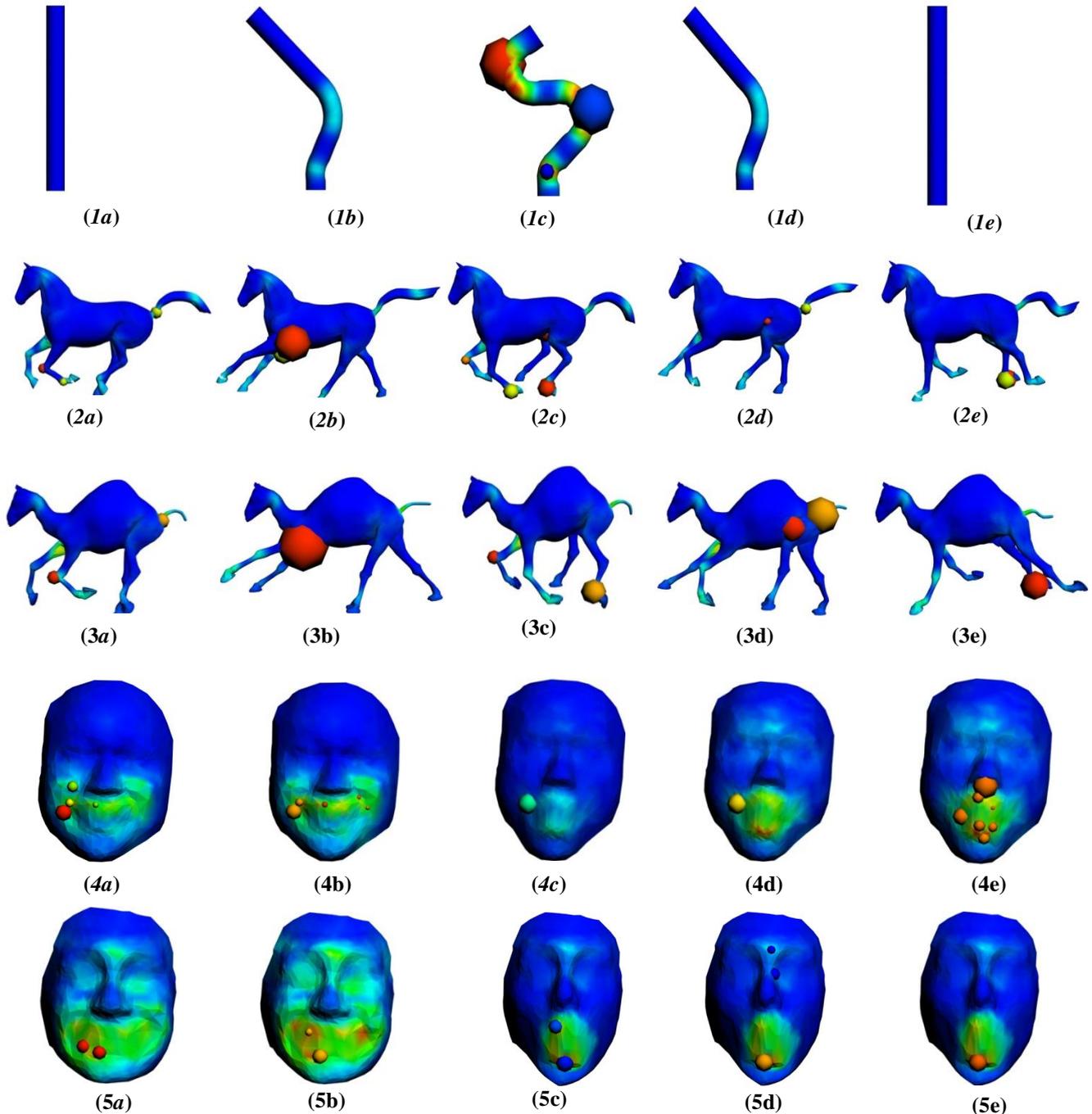
## Acknowledgements

We acknowledge Robert W. Sumner for providing triangle correspondences of the horse and camel models. We also thank Frederic Larue and Olivier G enevaux for their assistance with the facial motion capture. This work has been supported by the French national project SHARED (Shape Analysis and Registration of People Using Dynamic Data, No.10-CHEX-014-01).

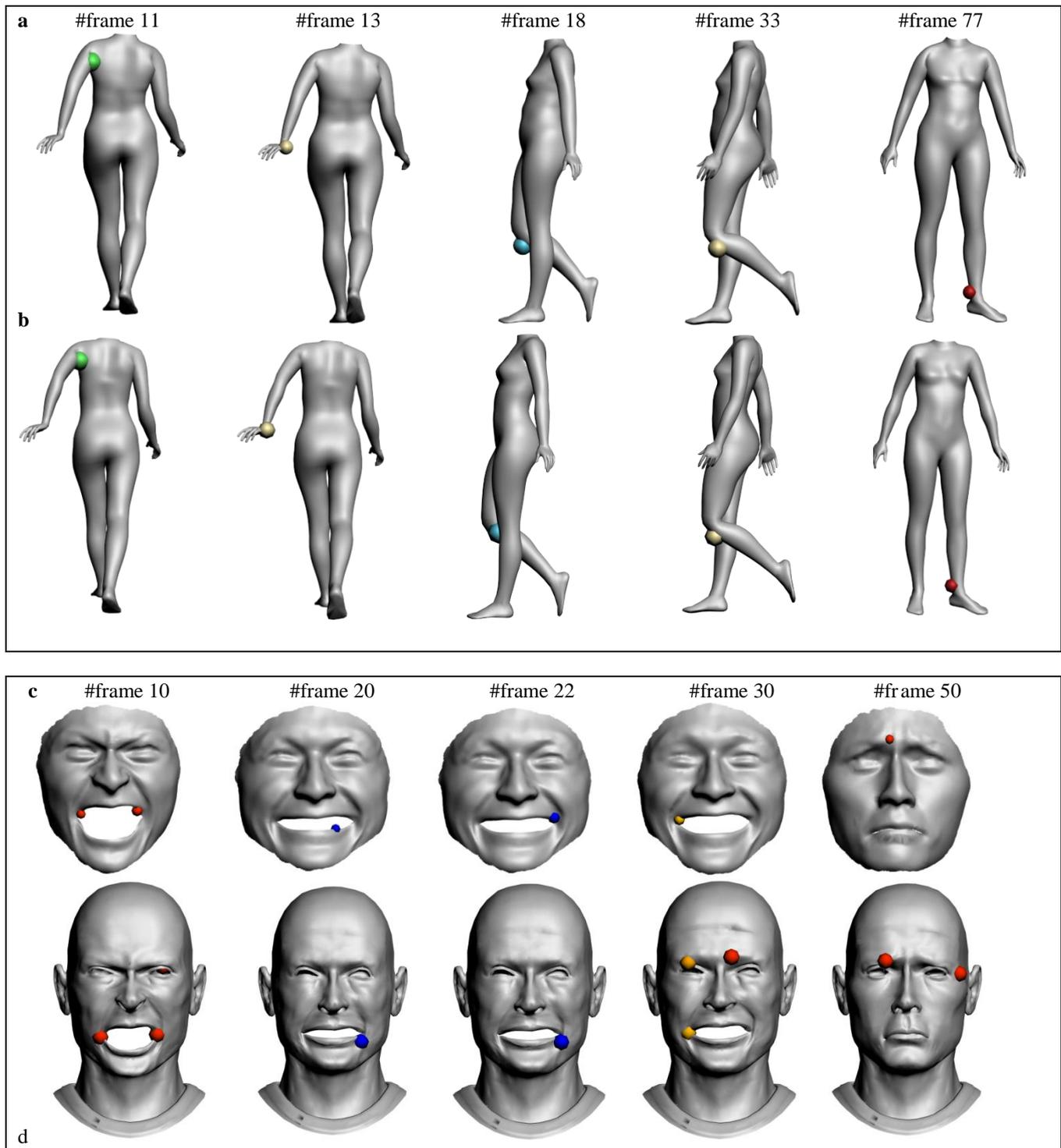
## References

- [AC92] R. Andonie, E. Carai: Gaussian Smoothing by Optimal Iterated Uniform Convolutions, *Computers and Artificial Intelligence*, 11(4): pp. 363-373, 1992.
- [ACD\*03] P. Alliez, D. Cohen-Steiner, O. Devillers, B. Levy, and M. Desbrun: Anisotropic Polygonal Remeshing, *ACM Transactions on Graphics*, 2003.
- [BET\*08] H. Bay, A. Ess, T. Tuytelaars, L. van Gool: Speeded-up Robust Features (SURF), *Computer Vision and Image Understanding (CVIU)*, Vol. 110, No. 3, pp. 346-359, 2008.
- [CCF\*08] U. Castellani, M. Cristani, S. Fantoni, V. Murino: Sparse points matching by combining 3D mesh saliency with statistical descriptors. *Comput. Graph. Forum* 27(2): 643-652 (2008).
- [DK12] T. Darrow, Y. Keller: Scale-Invariant Features for 3-D Mesh Models. *IEEE Transactions on Image Processing* 21(5): 2758-2769 (2012).
- [Got90] O. Gotoh: Optimal sequence alignment allowing for long gaps. *Bull. Math. Biol.* 52, 359-373, 1990.
- [HS88] C. Harris, M. Stephens: A combined corner and edge detector, In *Proc. 4th Alvey Vision Conference*, 147-151, 1988.
- [KG05] S. Kircher, M. Garland: Progressive multiresolution meshes for deforming surfaces, In *ACM SIGGRAPH Symposium on Computer Animation*, 191-200, 2005.
- [LGX13] Z. Lian, A. Godil, J. Xiao: Feature-Preserved 3D Canonical Form, 102: pp.221-238, *Int'l J. Computer Vision*, 2013.
- [LL03] I. Laptev, T. Lindeberg: Interest Point Detection and Scale Selection in Space-Time. *Scale-Space 2003*: pp. 372-387.
- [Lin94] T. Lindeberg: Scale-space theory: A basic tool for analyzing structures at different scales. *Journal of Applied Statistics*, 21(2): 224-270, 1994.
- [Lin98] T. Lindeberg: Feature detection with automatic scale selection. *International Journal of Computer Vision*, 30, pp.79-116,1998.
- [Low04] D. Lowe: Distinctive Image Features from Scale-Invariant Keypoints, 60(2): pp. 91-110, *Int'l Journal of Computer Vision*, 2004.
- [LVJ05] C.-H. Lee, A. Varshney, D. W. Jacobs: Mesh Saliency, *ACM Transactions on Graphics (Proc. SIGGRAPH)*, 2005.
- [Mesh] Mesh data, [http://people.csail.mit.edu/sumner/research/deft\\_ransfer/data.html](http://people.csail.mit.edu/sumner/research/deft_ransfer/data.html).
- [MS01] K. Mikolajczyk and C. Schmid: Indexing based on scale invariant interest points, *IEEE International Conference on Computer Vision (ICCV)*, pp. 525-531, 2001.
- [PKG03] M. Pauly, R. Keiser, M. Gross: Multi-scale Feature Extraction on Point-Sampled Surfaces, *Computer Graphics Forum*, 22(3), pp. 281-289, 2003.
- [SAS07] P. Scovanner, S. Ali, and M. Shah: A 3-Dimensional SIFT Descriptor and Its Application to Action Recognition. *Proceedings of the 15th International Conference on Multimedia*, pp. 357-360, (2007).
- [SB10] I. Sipiran and B. Bustos: A Robust 3D Interest Points Detector Based on Harris Operator, *Eurographics workshop on 3D Object Retrieval (3DOR)*, pp. 7-14 2010.
- [SP04] R. Sumner, J. Popovic: Deformation Transfer for Triangle Meshes. *ACM Transactions on Graphics*. 23, 3. August 2004.
- [SPB00] A. Shamir, V. Pascucci and C. Bajaj: Multi-resolution dynamic meshes with arbitrary deformations, pp. 423-430, *Proc. IEEE Visualization*, 2000.
- [SRV98] A.H. Salden, B.M. ter Haar Romeny and M.A. Viergever: Linear Scale-Space Theory from Physical

- Principles, In Journal of Mathematical Imaging and Vision, Vol. 9, pp. 103-139, 1998.
- [SZG\*05] R. W. Sumner, M. Zwicker, C. Gotsman, J. Popović: Mesh-based inverse kinematics. ACM Transactions on Graphics (TOG) 24(3), pp.488–495, 2005.
- [Vicon] Vicon motion capture system, <http://vicon.com>.
- [ZBV\*09] A. Zaharescu, E. Boyer, K. Varanasi, R. Horaud: Surface feature detection and description with applications to mesh matching, IEEE Computer Vision and Pattern Recognition (CVPR), pp. 373-380, 2009.



**Fig. 11** Dynamic feature points detected by our AniM-DoG framework are illustrated on a number of selected frames of animated meshes. The color of a sphere represents the temporal scale (from blue to red) of the feature point, and radius of sphere indicates the spacial scale.



**Fig. 12** Inter-subject consistency of feature points extracted from semantically similar mesh animations. Rows (a-b) depict subset of feature points extracted from walking subject sequences and (c-d) from face animations. Note that each column corresponds to identical frame of the animations.