



# DETECTING MOTION BY COMBINING THE STRUCTURE-TEXTURE IMAGE DECOMPOSITION AND SPACE-TIME INTEREST POINTS

<sup>1</sup>I.BELLAMINE, <sup>2</sup>H.TAIRI

<sup>1</sup> Sidi Mohamed Ben Abdellah University, LLIAN, Department of computer science, Morocco.

<sup>2</sup>Sidi Mohamed Ben Abdellah University, LLIAN, Department of computer science, Fes, Morocco.

E-mail: <sup>1</sup>[insaf\\_bellamine@hotmail.fr](mailto:insaf_bellamine@hotmail.fr), <sup>2</sup>[htairi@yahoo.fr](mailto:htairi@yahoo.fr)

## ABSTRACT

Among all the features which can be extracted from videos, we propose to use Space-Time Interest Points (STIP), these ones are particularly interesting because they are simple and robust. They allow a good characterization of a set of regions of interest corresponding to moving objects in a three-dimensional observed scene. In this paper, we show how the resulting features often reflect interesting events that can be used for a compact representation of video data as well as for tracking. For a good detection of moving objects, we propose to apply the algorithm of the detection of spatiotemporal interest points on both components of the decomposition which is based on a partial differential equation (PDE): a geometric structure component and a texture component. Proposed results are obtained from very different types of videos, namely sport videos and animation movies.

**Keywords:** *Space-Time Interest Points; Structure-Texture Image Decomposition; Motion Detection;*

## 1. INTRODUCTION

The motion analysis is a very active research area, which includes a number of issues: motion detection, optical flow, tracking and human action recognition.

To detect the moving objects in an image sequence is a very important low-level task for many computer vision applications, such as video surveillance, traffic monitoring, video indexing, recognition of gestures, analysis of sport-events, Sign language recognition, mobile robotics and the study of the objects' behavior (people, animals, vehicles, etc ...).

In the literature, there are many methods to detect moving objects, which are based on: optical flow [11], difference of consecutive images [2], Space-Time Interest Points [8] and modeling of the background (local, semi-local and global) [4].

Our method consists to use the notion of Space-Time Interest Points; these ones are especially interesting because they focus information initially contained in thousands of pixels on a few specific points which can be related to spatiotemporal events in an image. Laptev and Lindeberg were the first who proposed STIPs for action recognition [8], by introducing a space-time extension of the popular Harris detector [14]. They detect regions having high intensity variation in both space and time as spatio-temporal corners. The STIP detector of [8] usually suffers from sparse STIP detection.

Later, several other methods for detecting STIPs have been reported [10]. Dollar et al [10] improved the sparse STIP detector by applying temporal Gabor filters and selecting regions of high responses. Dense and scale-invariant spatio-temporal interest points were proposed by Willems et al. [9]. An evaluation of these approaches has been proposed in [12].

Our approach also uses Aujol Algorithm [29], this one decomposes the image  $f$  into a structure component  $u$  and a texture component  $v$ , ( $f = u + v$ ). The notion of Structure-Texture Image Decomposition is essential for understanding and analyzing images depending on their content.

In this paper, we propose to apply the algorithm of the detection of spatiotemporal interest points for a good detection of moving objects on both components of the decomposition: a geometric structure component and a texture component. Proposed results are obtained from different types of videos, namely sport videos and animation movies.

This paper is organized as follows: Section 2 presents the Spatio-Temporal Interest Points, Section 3 presents the Structure-Texture Image Decomposition, and finally, section 4 shows our experimental results.

## 2. SPACE-TIME INTEREST POINTS

The idea of interest points in the spatial domain can be extended into the spatio-temporal domain by

requiring the image values in space-time to have large variations in both the spatial and the temporal dimensions. Points with such properties will be spatial interest points with a distinct location in time corresponding to the moments with non-constant motion of the image in a local spatio-temporal neighborhood [15]. These points are especially interesting because they focus information initially contained in thousands of pixels on a few specific points which can be related to spatiotemporal events in an image.

Laptev et al. [15] proposed a spatio-temporal extension of the Harris detector to detect what they call "Space-Time Interest Points", denoted STIP in the following.

Detection of Space-Time Interest Points is performed by using the Hessian-Laplace matrix  $H$  [8], which is defined by:

$$H(x, y, t) = g(x, y, t; \sigma_s^2, \sigma_t^2) \otimes \begin{pmatrix} \frac{\partial^2 I(x, y, t)}{\partial x^2} & \frac{\partial^2 I(x, y, t)}{\partial x \partial y} & \frac{\partial^2 I(x, y, t)}{\partial x \partial t} \\ \frac{\partial^2 I(x, y, t)}{\partial x \partial y} & \frac{\partial^2 I(x, y, t)}{\partial y^2} & \frac{\partial^2 I(x, y, t)}{\partial y \partial t} \\ \frac{\partial^2 I(x, y, t)}{\partial x \partial t} & \frac{\partial^2 I(x, y, t)}{\partial y \partial t} & \frac{\partial^2 I(x, y, t)}{\partial t^2} \end{pmatrix} \quad (1)$$

$I(x, y, t)$  is the intensity of the pixel  $(x, y)$  at time  $t$ .

As with the Harris detector, a Gaussian smoothing is applied both in spatial domain (2D filter) and temporal domain (1D filter).

$$g(x, y, t; \sigma_s^2, \sigma_t^2) = \frac{\exp\left(-\frac{x^2 + y^2}{2\sigma_s^2} - \frac{t^2}{2\sigma_t^2}\right)}{\sqrt{(2\pi)^3 \sigma_s^4 \sigma_t^2}} \quad (2)$$

The two parameters ( $\sigma_s$  and  $\sigma_t$ ) control the spatial and temporal scale. As in [8], the spatio-temporal extension of the Harris corner function, entitled "saliency function", is defined by:

$$R(x, y, t) = \det(H(x, y, t)) - k \times \text{trace}(H(x, y, t))^3 \quad (3)$$

Where  $k$  is a parameter empirically adjusted at 0.04,  $\det$  is the determinant of the matrix  $H$  and  $\text{trace}$  is the trace of the same matrix.

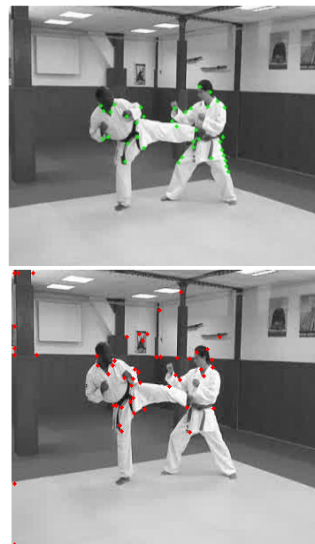
STIP correspond to high values of the saliency function  $R$  and they are obtained by using a thresholding step.

A. Tests.

In what follows, we represent some examples of clouds of space-time interest points detected in these sequences (Fig1):

- Sport video: (karate's fight) lasts for 2 minutes and 49 seconds with 200 images and the size of each image frame is 400 by 300 pixels.
- Animation movie: lasts for 3 minutes and 32 seconds with 230 images and the size of each image frame is 352 by 288 pixels.
- KTH dataset [31]: It was provided by Schuldt et al. [31] in 2004 and is one of the largest public human activity video dataset. It contains six types of actions (boxing, hand clapping, hand waving, jogging, running and walking) performed by 25 subjects in four different scenarios including in door, outdoor, changes in clothing and variations in scale. Each video clip contains one subject performing a single action. Each subject is captured in a total of 23 or 24 clips, giving a total of 599 video clips. Each clip has a frame rate of 25Hz and lasts between 10 and 15 s. The size of each image frame is 160 by 120 pixels. Two examples of the KTH dataset are shown in Fig1

We chose the value of 1.5 for the two standards deviation  $\sigma_s$  and  $\sigma_t$ , according to a study that was done by Alain Simac-Lejeune et al [13].



a) karate's fight (image t=54)



b) Animation movie (Trois petits points) (Images t=25 et t=30)



c) KTH dataset: hand waving, walking (Images t=20 and t=60)

Fig.1. Examples of clouds of space-time interest points. We have  $\sigma = 1.5$  and  $\sigma_t = 1.5$  with  $k=0, 04$ . In each frame, the red points represent the extracted Space Interest Points (Harris [14] detector), and the green points are the extracted Space-Time Interest Points (Laptev [15] detector).

### 3. STRUCTURE-TEXTURE IMAGE DECOMPOSITION

Let  $f$  be an observed image which contains texture and/or noise. Texture is characterized as repeated and meaningful structure of small patterns.

Noise is characterized as uncorrelated random patterns. The rest of an image, which is called *cartoon*, contains object hues and sharp edges (boundaries). Thus an image  $f$  can be decomposed as

$f = u + v$ , where  $u$  represents image cartoon and  $v$  is texture and/or noise.

Decomposing an image  $f$  into a geometric structure component  $u$  and a texture component  $v$  is an inverse estimation problem, essential for understanding and analyzing images depending on their content.

Many image decomposition models have been proposed, those based on the total variation as the

model of the ROF minimization was proposed by Rudin, Osher and Fatimi in [24], this model has demonstrated its effectiveness and has eliminated oscillations while preserving discontinuities in the image, it has given satisfactory results in image restoration [25, 26] because the minimization of the total variation smooths images without destroying the structure edges.

In recent years, several models based on total variation, which are inspired by the ROF model, were created [21, 23]. In the literature there is also another model called Mayer [27] that is more efficient than the ROF model. Many algorithms have been proposed to solve numerically this model. In the following, we represent the most popular algorithm, Aujol algorithm.

#### The Aujol Algorithm

In [29], J. F. Aujol et al propose a new algorithm for computing the numerical solution of the Meyer's model. An image  $f$  can be decomposed as  $f = u + v$ , where  $u$  represents a geometric structure component and  $v$  is a texture component.

The algorithm of Aujol is represented as follows:

#### Step 1: Initialisation

$$u_0 = v_0 = 0$$

#### Step 2: Iterations

$$v_{n+1} = P_{G_\mu}(f - u_n)$$

$$u_{n+1} = f - v_{n+1} - P_{G_\lambda}(f - v_{n+1})$$

where  $P_{G_\mu}$  and  $P_{G_\lambda}$  are the operators of projections(see [30])

#### Step 3: Stop condition

We stop the algorithm if

$$\max(|u_{n+1} - u_n|, |v_{n+1} - v_n|) \leq \epsilon$$

Or if it reaches a maximum of iterations required.

The regularization of parameters ( $\lambda$  and  $\mu$ ) play an important role in the decomposition of the image:  $\lambda$  controls the amount of the residue  $f - u - v$  and  $\mu$  influences on the texture component  $v$ . The choice of  $\lambda$  does not pose a problem. It gives it just a small value, but the  $\mu$  parameter is not easy to adjust.

Well the Aujol algorithm can extract the textures in the same way as the Osher-Vese

algorithm. Moreover, this algorithm has some advantages if it is compared to Osher-Vese [32]:

- No problem of stability and convergence
- Easy to implement (requiring only a few lines of code).

*B. Decomposition results*

Let  $f$  the image to decompose, then  $f$  can be written as follows:  $f = u + v$

The Structure-Texture Image Decomposition has been applied on the Barbara image of size  $512 \times 512$ . The result of the decomposition using the parameters ( $\mu = 1000, \lambda = 0.1$ ) is shown in Fig 2.

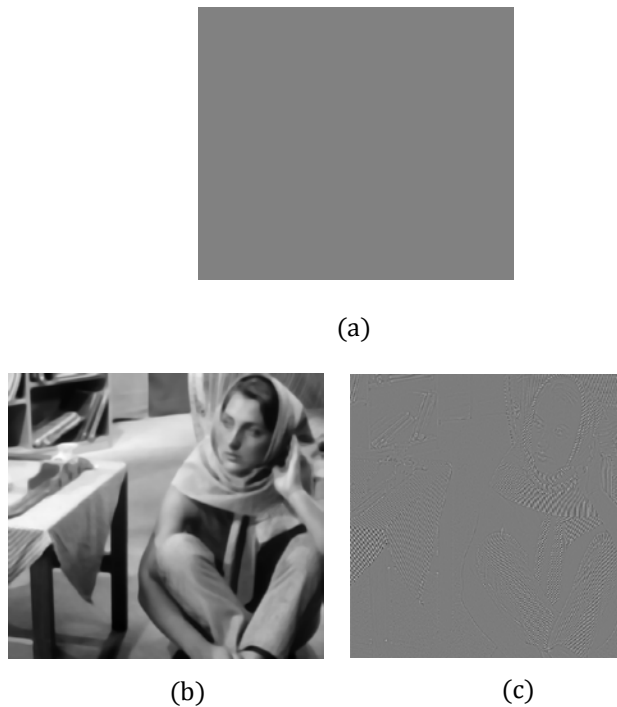


Fig. 2. Decomposition with Aujol Algorithm [29]: (a) Barbara image (b) the  $u$  component (c) the  $v$  component.

The program was run in a PC with a 2.13 GHz Intel core (TM) i3 CPU with 3 GB RAM.

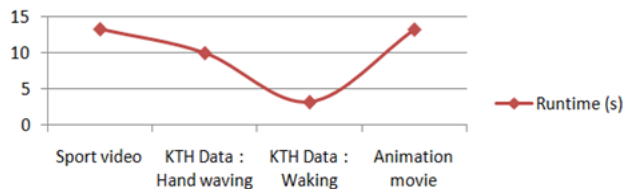


Fig 3:Extraction time of Structure-Texture Image Decomposition

Decomposing an image  $f$  into a geometric structure component and a texture component requires relatively low computation time, which gives us the opportunity to use this decomposition in motion detection in real time.

4. PROPOSED APPROACH

The most famous algorithm to detect Space-Time Interest Points is that of Laptev; however we can reveal three major problems when a local method is used:

- Texture, Background and Objects that may influence the results.
- Noisy datasets such as the KTH dataset, which is featured with low resolution ,strong shadows ,and camera movement that renders clean silhouette extraction impossible
- Features extracted are unable to capture smooth and fast motions, and they are sparse .This also explains why they generate poor results.

However, to overcome the three problems, we propose a technique based on the space-time interest points , and which will help to have a good detection of moving objects and even reduce the execution time by proposing a parallel algorithm (see Fig 4).

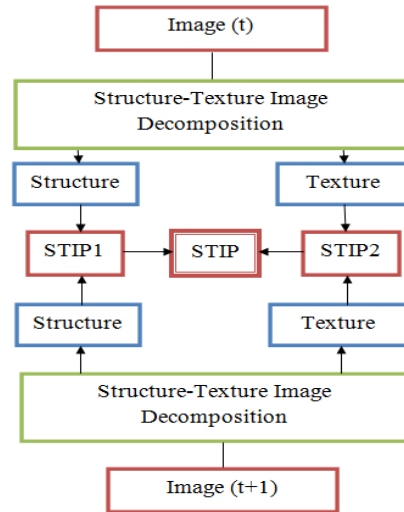


Fig. 3. The adapted Laptev algorithm

A complex scene can contain various information (noise, textures, shapes, background ...) , these ones influence the detection of moving objects. Our goal is to apply the algorithm of the detection of spatiotemporal interest points on both components

of the decomposition: a geometric structure component and a texture component.

Let STIP denote the final Space-Time Interest Points, STIP1 denotes the extracted Space-Time Interest Points between the components of the decomposition of Structure1 and Structure 2, also STIP2 denotes the extracted Space-Time Interest Points using Texture1 and Texture 2 components.

Our new Space-Time Interest Points will be calculated as the following:

$$(4) \quad STIP = STIP1 \cup STIP2$$

A. Enhanced Laptev algorithm

In the first step, the Structure-Texture Image Decomposition method is applied to the two consecutive frames of the video sequence. In the second step, two processes based on structures (Structure1 and Structure2) and textures (Texture1 and Texture2) had to be made equivalent to the two matching modes. Each process provides, as output result, the STIP1 extracted from the first mode and the STIP 2 extracted from the second mode. For the last step, the final STIP are computed by the equation 4. Figure 5 shows the steps of the enhanced Laptev algorithm.

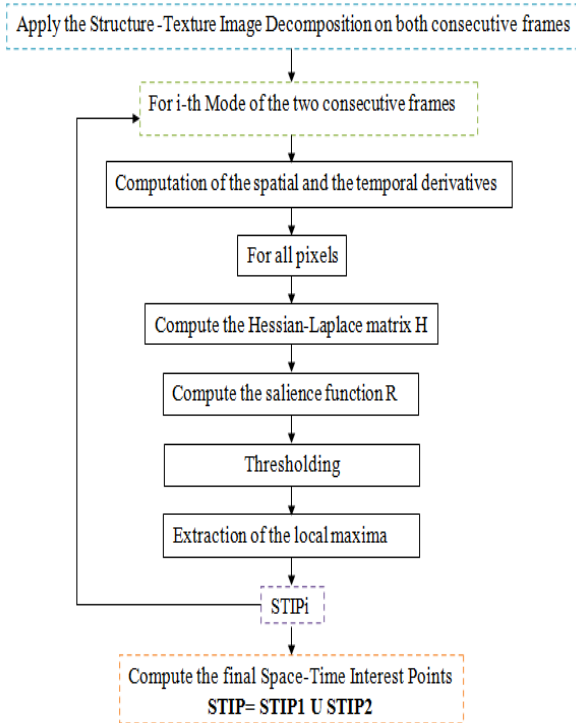


Fig. 4. Enhanced Laptev algorithm

The results illustrated in Fig 6, show that we come to locate moving objects in both sequences, we note that: the objects moving (the two players) are detected with our approach, for against just one player who has been detected with Laptev detector (Fig1 (a)).

B. Results and discussion

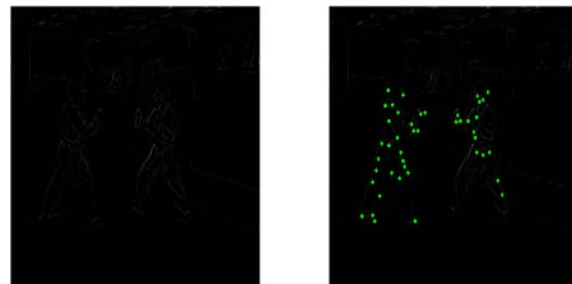
1) Experimental results.

Let  $f$  be an observed image which contains texture and/or noise. The rest of an image, which is called geometric structure, contains object hues and sharp edges (boundaries). Thus an image  $f$  can be decomposed as  $f = u + v$ , where  $u$  represents structure image and  $v$  is texture and/or noise.

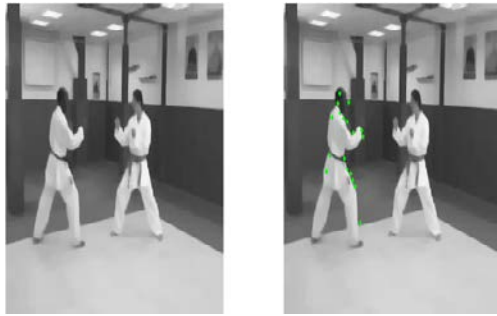
We propose to apply the algorithm of the detection of spatiotemporal interest points on both components of the decomposition: a geometric structure component and a texture component. In what follows, we represent some examples of clouds of space-time interest points detected in the first sequence (Fig 1):



i) age (t=44) / Image (t=45)



j) The green points are the detected Space-Time Interest Points on the texture components



k) The green points are the detected Space-Time Interest Points on the structure components

$$\text{Precision} = \frac{\text{NTP}}{\text{NTP} + \text{NFP}} \quad (5)$$

NTP is the number of the true positives (good detections), NFP the number of the false positives (false detections).



l) The red points represents the extracted Space-Time Interest Points with Our Proposed Approach

Fig. 5. Examples of clouds of space-time interest points, in each frame, we use the parameters (  $\sigma = 1.5$ ,  $k = 0.04$  and  $\sigma_s = 1.5$ ).

### C. Comparison and discussion

In order to correctly gauge performance of our proposed approach, we will proceed with a Comparative study to use the mask of the moving object and the precision.

The mask of the moving object is obtained by the Markovian approach [34] (see Fig 7). We distinguish two cases:

- True positive: the space-time interest point is in the mask, so it is on a moving object
- False positive: the space-time interest point isn't in the mask, so it isn't on a moving object

For each moving object, we have a number of the space-time interest points detected in the moving object (NTP) and a number of the space-time interest points extracted off the moving object (NFP).

The precision is defined by:

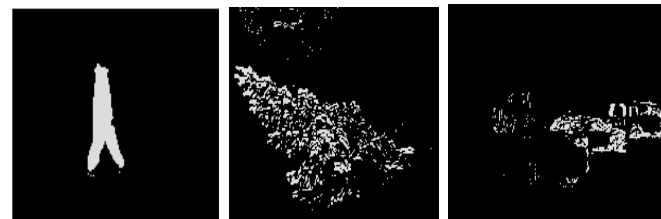
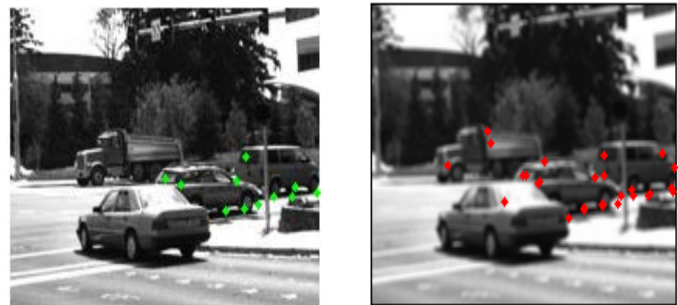
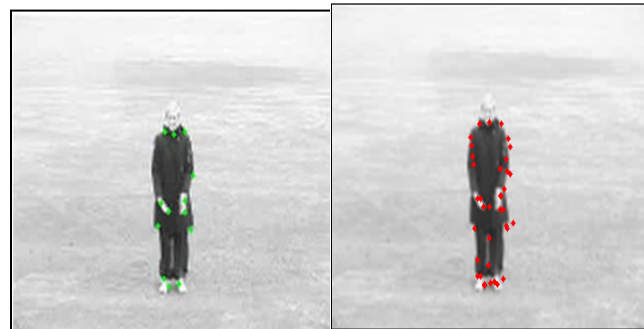


Fig. 6. Examples of moving objects' masks

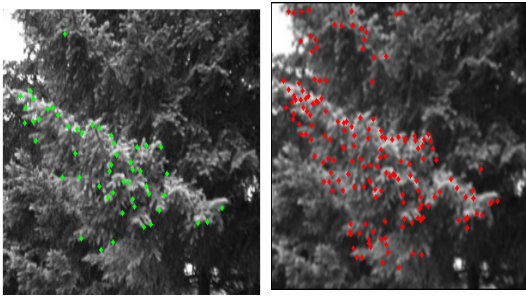
The test is performed on four examples of sequences and gives the following results:



a) Urban transport



b) Hand waving



Tree



Karate

Fig. 7. Examples of clouds of space-time interest points. We have  $\sigma_s = 1.5$  and  $\sigma_t = 1.5$  with  $k=0, 04$ . In each frame, the green points represent the extracted Space-Time Interest Points by Laptev [15] detector, and the red points are the extracted Space-Time Interest Points by Our Approach.

The results show that the real moving objects (the two players, the cars, the truck and branches of the tree) are better detected with the proposed approach than with laptev [15] detector

Still, the proposed approach is much less sensitive to noise, and the reconstruction of the image, it also extracts the densest features:

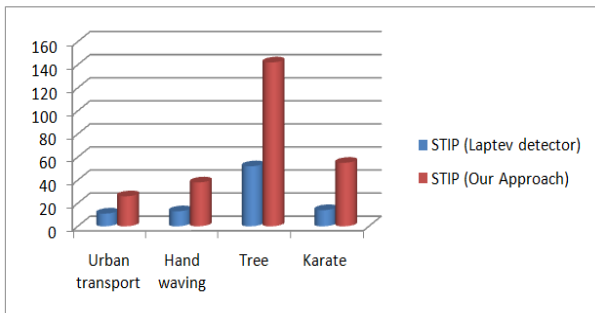


Fig. 8. Number of Space-Time Interest Points extracted in each frame.

TABLE I. Compared Results.

Videos	Precision (Laptev detector)	Precision (Our Approach)
Urban transport	81%	89%
Hand waving	92%	93%
Tree	96%	98%
Karate	92%	96%

The results, illustrated in Table 1, show that our approach allows a good detection of moving objects.

#### D. Tracking with the STIP

Object tracking module is responsible for estimating the location of each object in each new frame. This module ensures that the object is being tracked even if the motion detection module fails to detect it, either due to occlusion or stopped object.

We propose a technique based on the space-time interest points detected with our approach and Kanade-Lucas algorithm [33] (local method of the optical flow) which will help to have a tracking by proposing this algorithm (see Fig 8).

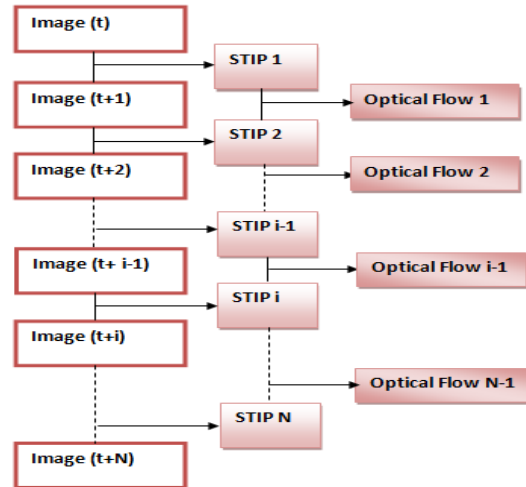
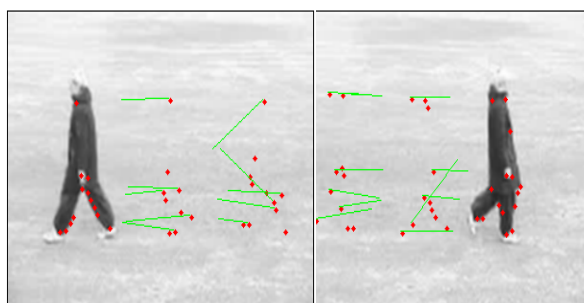


Fig. 9. The Algorithm Of Tracking

In what follows (see Fig 9), we represent the tracking by giving some examples of clouds of space-time interest points detected in the KTH dataset: walking.



a) Images (t=59, t=60, t=61, t=62)



B) Walking On The Left C) Walking On The Right

Fig. 10. An Example Of Our Tracking Algorithm In Operation. Features Being Tracked Are Shown As A Red Dot, With The Inter-Frame Displacement Shown As A Green Line. For This Example, We Used Four Images.

## 5. CONCLUSION

In the experimental part, the results are obtained from very different types of videos, namely sport videos and animation movies.

Our Approach improved the sparse STIP detector by applying the algorithm of the detection of spatiotemporal on both components of the decomposition: a geometric structure component and a texture component. This Approach is less sensitive to the noise effects and the parallel implementation requires low computation time.

## REFERENCES

- [1] A. Bugeau. Détection et suivi d'objets en mouvement dans des scènes complexes, application à la surveillance des conducteurs. Thèse de doctorat, IRISA, 2007.
- [2] Eric Galmar et Benoit Huet : Analysis of vector space model and spatiotemporal segmentation for video indexing and retrieval. In CIVR 2007, ACM International Conference on Image and Video Retrieved, July 9-11 2007, Amsterdam, The Netherlands, 07 2007.
- [3] R. Jain, H.H. Nagel. – On the analysis of accumulative difference pictures from image sequence of real world scenes. IEEE Trans. Pattern Anal. Machine Intel. 1(2) :20214, 1979.
- [4] V. Nicolas. - Suivi d'objets en mouvement dans une séquence vidéo. Thèse de doctorat, Université Paris Descartes. 2007
- [5] S. Pundlik, S. Birchfield. – Motion segmentation at any speed. Proc. of the British Machine Vision Conf., 2006.
- [6] M. Bregonzio n, TaoXiang,ShaogangGong , Fusing appearance and distribution information of interest points for action recognition
- [7] R. Laganière, R. Bacco, A. Hocevar, P. Lambert, G. Païs, et B.E. Ionescu. Video summarization from spatio-temporel features. ACM, 2008.
- [8] I. Laptev. On space-time interest points. International Journal of Computer Vision, 64(2/3) :107–123, 2005.
- [9] G. Willems, T. Tuytelaars, and L. Van Gool. An efficient dense and scale-invariant spatio-temporal interest point detector. European Conference on Computer Vision, 5303(2) :650-663, 2008.
- [10] P. Dollar, V. Rabaud, G. Cottrell, S. Belongie, Behavior recognition via sparse spatio-temporal features, in: VS-PETS, 2005.
- [11] A. Simac. Optical-flow based on an edge-avoidance procedure. Computer Vision and Image Understanding 113(2009): 511–531, 2008.
- [12] H.Wang. Evaluation of local spatio-temporal features for action recognition. BMVC '09 London, 2009.
- [13] A. Simac. Modélisation et gestion de concepts, en particulier temporels, pour l'assistance à. Thèse de Doctorat, Université de Grenoble, 2006.
- [14] C. Harris et M.J. Stephens. A combined corner and edge detector. In Alvey Vision Conference, 1988.
- [15] Laptev et T. Lindeberg. Space-time interest points. ICCV'03, pages 432–439, 2003.
- [16] . O. Riouland, M. Vetterli. Wavelets and signal processing. IEEE Signal Processing, 8(4) :14-38, 1991.
- [17] . M. E. Zervakis, V. Sundararajan, K. K. Parhi. Vector processing of wavelet coefficients for robust image denoising. Image and Vision Computing, 19 (1) :435-450, 2001.
- [18] J. C. Nunes, Y. Bouaoune, E. Delechelle, O. Niang, and Ph. Bunel. Image analysis by bidimensional empirical mode decomposition. Image and Vision Computing, 21:1019-1026, 2003.
- [19] J.C. Nunes, S. Guyot, and E. Deléchelle. Texture analysis based on local analysis of the bidimensional empirical mode decomposition. Journal of Machine Vision and Applications, 16 (3) :177-188, 2005.





- [20] S. M. A. Bhuiyan, R. R. Adhami, J. F. Khan. A novel approach of fast and adaptive bidimensional empirical mode decomposition. IEEE International Conference on Acoustics, Speech and Signal Processing, Las Vegas, 2008.
- [21] J. F. Aujol. Contribution à l'analyse de textures en traitement d'images par méthodes variationnelles et équations aux dérivées partielles. Thèse de Doctorat, Université de Nice Sophia Antipolis, 2004.
- [22] J. F. Aujol, G. Aubert, L. Blanc Féraud, and A. Chambolle. Decomposing an image. application to textured images and SAR images. Rapport technique, Université de Nice Sophia-Antipolis, 2003.
- [23] J. Gilles. Décomposition et détection de structures géométriques en imagerie. Thèse de Doctorat, Ecole Normale Supérieure de Cachan, 2006.
- [24] L. Rudin, S. Osher, End E. Fatimi. Nonlinear total variation based noise removal algorithms. Physica D, 60 : 259-268, 1992.
- [25] T. F. Cha, S. Osher, J. Shen. The digital TV filter and nonlinear denoising. IEEE Transactions Image Processing. 10 :231-241, 2001.
- [26] S. Osher, L. Rudin. Total variation based image restoration with free local constraints. I Proc. IEEE ICIP, volume I, Austin, TX : 31-35, 1994.
- [27] Y. Meyer. Oscillating patterns in image processing and in some nonlinear evolution equations. The Fifteenth Dean Jacqueline B. Lewis Memorial Lectures, American Mathematical Society, 2001.
- [28] L. Vese, S. Osher. Modeling textures with total variation minimization and oscillating patterns in image processing. Journal of scientific computing, 19 (1-3) :553-572, 2002.
- [29] J. F. Aujol. Contribution à l'analyse de textures en traitement d'images par méthodes variationnelles et équations aux dérivées partielles. Thèse de Doctorat, Université de Nice Sophia Antipolis, 2004.
- [30] A. Chambolle. An algorithm for Total Variation Minimization and application. Journal of Mathematical Imaging and vision, 20 (1-2) : 89-97, 2004.
- [31] Laptev, B. Caputo, Recognizing human actions: a local SVM approach, ICPR, vol.3, 2004, pp.32-36.
- [32] J. GILLES. Décomposition et détection de structures géométriques en imagerie. Thèse de Doctorat, ECOLE NORMALE SUPERIEURE DE CACHAN, 2006.
- [33] J.-Y. Bouguet, Pyramidal Implementation of the Lucas Kanade Feature Tracker, Intel Corporation, Microprocessor Research Labs (2000).
- [34] F. Luthon, Module de Traitement d'Image, Laboratoire d'Informatique de l'Université de Pau et des Pays de l'Adour (2001).