

A tensor motion descriptor based on histograms of gradients and optical flow

Mota V. F.^{1a,b}, Perez E. A.^a, Maciel L. M.^a, Vieira M. B.^a, Gosselin, P. H.^c

^a*DCC/ICE, Universidade Federal de Juiz de Fora, Juiz de Fora, Brazil*

^b*DCC/ICE, Universidade Federal de Minas Gerais, Belo Horizonte, Brazil*

^c*INRIA Rennes Bretagne Atlantique*

Abstract

This paper presents a new tensor motion descriptor only using optical flow and HOG3D information: no interest points are extracted and it is not based on a visual dictionary. We propose a new aggregation technique based on tensors. This is a double aggregation of tensor descriptors. The first one represents motion by using polynomial coefficients which approximates the optical flow. The other represents the accumulated data of all histograms of gradients of the video. The descriptor is evaluated by a classification of KTH, UCF11 and Hollywood2 datasets, using a SVM classifier. Our method reaches 93.2% of recognition rate with KTH, comparable to the best local approaches. For the UCF11 and Hollywood2 datasets, our recognition achieves fairly competitive results compared to local and learning based approaches. *Keywords:* Global motion descriptor, optical flow, histogram of gradients, action recognition

¹Corresponding author: Tel: +55 32 88856321
E-mail address: virginiaferm@dcc.ufmg.br

1. Introduction

Human action recognition is a very attractive field of research as it is a key part in several areas such as video indexing, surveillance, human-computer interfaces, among others. Most works address this problem by a motion analysis and a representation step. Several descriptors were proposed over the past years, most of them using some motion representation, because it is one of the main characteristics that describe the semantic information of videos. Some examples of motion representations are the histogram of gradients and optical flow.

Usually the optical flow itself is not used as a descriptor. Instead, its histogram is largely associated with other features in order to improve the recognition rate [1, 2]. In our preliminary work, presented in [3], we showed that the modeling of optical flow vector fields gives a consistent global motion descriptor. This descriptor is obtained using the parameters of a polynomial model for each frame of a video. The coefficients were found through the projection of the optical flow on Legendre polynomials, reducing the dimension of the motion estimation per frame. The sequence of coefficients were then combined using orientation tensors.

This work is motivated by the possibility of combining the tensor descriptor presented in [3] with other global features. Indeed, the optical flow projected onto Legendre polynomial basis captures a specific nuance of the underlying motion. Its combination with other motion representations can improve the results and drive a competitive recognition for the problem of human action recognition.

Our main contribution is a new motion descriptor based on orientation

26 tensor which uses only optical flow [3] and tridimensional histogram of gra-
27 dients (HOG3D) information [4]: no interest points are extracted and no
28 bag-of-features strategy is used. The global tensor descriptor created is eval-
29 uated by a classification of KTH [5], UCF11 (also known as UCF YouTube)
30 [6] and Hollywood2 [7] video datasets with a non-linear SVM classifier.

31 **2. Related work**

32 Laptev et al [2] present a combination of histograms of gradients (HOG)
33 with histogram of optical flow (HOF) to characterize local motion and shape.
34 Histograms of spatial gradient and optical flow are computed and accumu-
35 lated in space-time neighborhoods of detected interest points. Similarly to
36 the SIFT descriptor, normalized histograms are concatenated to HOG and
37 HOF vectors. Then, the signature of the video is computed through a bag-
38 of-features technique.

39 In [1], HOG, HOF, MBH (motion boundary histogram) and trajectory are
40 combined in order to create a better motion descriptor. For each descriptor
41 type, bag-of-features are computed thanks to a visual codebook. A SVM
42 classifier is then used in the context of action classification for the KTH,
43 Hollywood2, UCF11 and UCF sports datasets.

44 Also using a bag-of-features strategy, Zhen and Shao [8] presents a new
45 descriptor for action recognition based on Laplacian pyramid coding. The
46 idea is to represent the video by the combination of motion history images and
47 three orthogonal planes, obtained from a set of cuboids extracted from the
48 video sequence. Then, this information is encoded with a Laplacian pyramid
49 model and the final video representation is computed thanks to an improved

50 version of bag-of-features using the soft-assignment coding and max pooling.

51 Kobayashi and Otsu [9] propose motion features based on co-occurrence
52 histograms of the space-time 3D gradient orientations. They are employed
53 for frame based features to densely characterize the motion. These frame-
54 based features are extracted from sub-sequences densely sampled along the
55 time axis. Thus, they describe a bag-of-frame-features approach to create
56 the video feature.

57 The use of local features for human action recognition is more exploited,
58 as they provide higher recognition rates. In general, these approaches use
59 bag-of-features (BoF) strategy. Hence, there are few references about global
60 descriptors which do not rely on a visual dictionary and are uniquely depen-
61 dent on the video, instead of the whole training set as such in BoF method.
62 Global approaches, however, are much simpler to compute and can achieve
63 fast and fairly high recognition rates.

64 Zelnik et al presents a global descriptor based on histogram of gradi-
65 ents [10]. This descriptor is applied on the Weizmann video database and
66 is obtained with the extraction of multiple temporal scales through the con-
67 struction of a temporal pyramid. To calculate this pyramid, they apply a
68 lowpass filter on the video and sample it. For each scale, the intensity of
69 each pixel gradient is calculated. Then, a histogram of gradients is created
70 for each video and compared with others histograms to classify the database.

71 In order to obtain a global descriptor on the KTH dataset, Laptev et al
72 [11] apply the Zelnik descriptor [10] in two different ways: using multiple
73 temporal scales like the original and using multiple temporal and spatial
74 scales.

75 Solmaz et al [12] present a global descriptor based on bank of 68 Gabor
76 filters. For each video, they extract a fixed number of clips and compute the
77 3-D Discrete Fourier Transform. Applying each filter of the 3-D filter bank
78 separately to the frequency spectrum, the output is quantized in fixed sub-
79 volumes. They concatenate the outputs and perform dimension reduction
80 using PCA and classification by a SVM.

81 **3. Proposed Method**

82 *3.1. Tensor based on optical flow approximation*

83 The basic idea of a polynomial based model is to approximate a vector
84 field with a linear combination of orthogonal polynomials [13, 14]. Let us
85 define F an optical flow:

$$F : \Omega \subset R^2 \rightarrow R^2 \\ (x_1, x_2) \mapsto (V^1(x_1, x_2), V^2(x_1, x_2))$$

86 where the functions $V^1(x_1, x_2)$ and $V^2(x_1, x_2)$ corresponds to the horizontal
87 and vertical displacement of the point $(x_1, x_2) \in \Omega$.

88 This optical flow is then approximated by projecting the displacement
89 functions onto each polynomial $P_{i,j}$, which belong to an orthogonal basis, as
90 such Legendre basis.

91 In that way, it reduces the dimension of the optical flow field. Thus, we
92 can express $\tilde{F} = (\tilde{V}^1(x_1, x_2), \tilde{V}^2(x_1, x_2))$, using a basis of degree g , as:

$$\begin{cases} \tilde{V}^1(x_1, x_2) = \sum_{i=0}^g \sum_{j=0}^{g-1} \tilde{v}_{i,j}^1 P_{i,j} \\ \tilde{V}^2(x_1, x_2) = \sum_{i=0}^g \sum_{j=0}^{g-1} \tilde{v}_{i,j}^2 P_{i,j} \end{cases}$$

93 where

$$\begin{cases} \tilde{v}_{i,j}^1 = \int \int_{\Omega} V^1(x_1, x_2) P_{i,j} \omega(x_1, x_2) dx_1 dx_2 \\ \tilde{v}_{i,j}^2 = \int \int_{\Omega} V^2(x_1, x_2) P_{i,j} \omega(x_1, x_2) dx_1 dx_2 \end{cases} \quad (1)$$

94 It is important to note that the number of polynomials which composes
95 a basis of degree g is:

$$n_g = \frac{(g+1)(g+2)}{2}$$

96 *3.1.1. Orientation tensor: coding frame coefficients*

97 An orientation tensor is a representation of local orientation which takes
98 the form of an $m \times m$ real symmetric matrix for m -dimensional signals [15].

99 Given the vector \vec{v} with m elements, it can be represented by the tensor
100 $T = \vec{v}\vec{v}^T$. It is desired that the eigenvector with the largest eigenvalue of
101 the tensor points out the dominant direction of the signal. A signal with
102 no dominant direction is represented by an isotropic tensor, i.e. the three
103 eigenvalues are approximately equal. It is important to note that the well
104 known structure tensor is a specific case of orientation tensor [16].

105 In order to capture the motion variation in time, we can use both the
106 polynomial coefficients $\tilde{v}_{i,j}^1$ and $\tilde{v}_{i,j}^2$ (Eq. 1) and an approximation of their
107 first temporal derivative $\partial \tilde{v}_{i,j}^q = \tilde{v}_{i,j}^q(f) - \tilde{v}_{i,j}^q(f-1)$ with $i+j < g$, to create
108 a vector \tilde{v}_f for each frame f of the video:

$$\tilde{v}_f = [\tilde{v}_{0,0}^1, \dots, \tilde{v}_{g,0}^1, \tilde{v}_{0,0}^2, \dots, \tilde{v}_{g,0}^2, \partial \tilde{v}_{0,0}^1, \dots, \partial \tilde{v}_{g,0}^1, \partial \tilde{v}_{0,0}^2, \dots, \partial \tilde{v}_{g,0}^2].$$

109 Using the vector \tilde{v}_f , we generate an orientation tensor $T_f = \tilde{v}_f \tilde{v}_f^T$ for
110 each frame f of the video, which is a $4n_g \times 4n_g$ matrix. This orientation
111 tensor captures the covariance information between $\tilde{v}_{i,j}^1$ and $\tilde{v}_{i,j}^2$. It carries
112 only the information of the polynomial of frame f and its rate of change in
113 time.

114 *3.1.2. Global tensor descriptor*

115 We have to express the motion average of consecutive frames using a series
116 of tensors. This can be achieved by $T^{OF} = \sum_a^b T_f$ using all video frames or
117 an interval of interest. By normalizing T_f with a L_2 norm, we are able to
118 compare different video clips or snapshots regardless their length or image
119 resolution.

120 If the accumulation series diverges, we obtain an isotropic tensor which
121 does not hold useful motion information. But, if the series converge as an
122 anisotropic tensor, it carries meaningful average motion information of the
123 frame sequence. The conditions of divergence and convergence need further
124 studies.

125 Instead of using the entire optical flow of the video frames, it is also
126 possible to use only the optical flow from a region with most representative
127 motion. Then, we tested a sliding window with fixed dimensions placed
128 around the subject who is doing the action. The center of mass of global
129 optical flow gives the center of the window.

130 The accumulated tensor is symmetric, therefore we can use only a trian-
131 gular superior (or inferior) matrix to represent the video, which reduces the
132 number of coefficients of the final tensor descriptor.

133 *3.2. Tensor based on histogram of gradients*

134 The partial derivatives of the j -th video frame at point p

$$\vec{g}_t(p) = [dx \ dy \ dt] = \left[\frac{\partial I_j(p)}{\partial x} \quad \frac{\partial I_j(p)}{\partial y} \quad \frac{\partial I_j(p)}{\partial t} \right],$$

135 or, equivalently, in spherical coordinates $\vec{s}_t(p) = [\rho_p \ \theta_p \ \psi_p]$ with $\theta_p \in [0, \pi]$,
136 $\psi_p \in [0, 2\pi]$ and $\rho_p = \|\vec{g}_t(p)\|$, indicate brightness variation that might be

137 the result of local motion.

138 The gradient of all n points of the image I_j can be compactly represented
 139 by a tridimensional histogram of gradients $\vec{h}_j = \{h_{k,l}\}$, $k \in [1, nb_\theta]$ and
 140 $l \in [1, nb_\psi]$, where nb_θ and nb_ψ are the number of cells for θ and ψ coordi-
 141 nates respectively. There are several methods for computing the HOG3D. We
 142 performed some experiments with the icosahedron discretization [17], how-
 143 ever no significant enhancement was detected. Thus, we chose an uniform
 144 subdivision of the angle intervals to populate the $nb_\theta \cdot nb_\psi$ bins (Eq. 2), since
 145 it achieves good results and it is fast to compute.

$$h_{k,l} = \sum_p \rho_p \cdot w_p, \quad (2)$$

146 where $\{p \in I_j \mid k = 1 + \lfloor \frac{nb_\theta \cdot \theta_p}{\pi} \rfloor, l = 1 + \lfloor \frac{nb_\psi \cdot \psi_p}{2\pi} \rfloor\}$ are all points whose angles
 147 map to k and l bins, and w_p is a per pixel weighting factor which can be
 148 uniform or gaussian as in [18]. The whole gradient field is then represented
 149 by a vector \vec{h}_j with $nb_\theta \cdot nb_\psi$ elements.

150 3.2.1. Global tensor descriptor: coding HOG3D coefficients as tensors

151 Analogously to the previous descriptor (Sec. 3.1.2), the HOG3Ds with m
 152 bins \vec{h}_j , computed for j -th frames, can be combined in a tensor as following:

$$T^{HOG} = \sum_j \vec{h}_j \vec{h}_j^T,$$

153 using all video frames or an interval of interest. By normalizing T^{HOG} with a
 154 L_2 norm, we are able to compare different video clips or snapshots regardless
 155 their length or image resolution.

156 *3.2.2. Global tensor descriptor: subdividing the frame using a grid*

157 When the gradient histogram is computed using the whole image, the
 158 cells are populated with vectors regardless their position in the image. This
 159 implies in a loss of the correlation between the gradient vectors and their
 160 neighbors. As observed in several works [18], the subdivision of the video
 161 into cubes of frames enhances the recognition rate, using a gaussian weight
 162 for w_p .

163 Suppose the video frame f is uniformly subdivided in \vec{x} and \vec{y} directions
 164 by a grid with n_x and n_y non-overlapping blocks. Each block can be viewed
 165 as a distinct video varying in time. The smaller images result in gradient
 166 histograms $\vec{h}_j^{c,r}$, $c \in [1, n_x]$ and $r \in [1, n_y]$, with better position correlation.
 167 The tensor for the frame j is then computed as the addition of all block
 168 tensors:

$$T_j = \sum_{c,r} \vec{h}_j^{c,r} \vec{h}_j^{c,r}{}^T, \quad (3)$$

169 which captures the uncertainty of the direction of the m -dimensional vectors
 170 $\vec{h}_f^{a,b}$ in the frame j . This tensor is normalized using the L_2 norm. The
 171 image subdivision does not change the descriptor size and the accumulation
 172 described above is the same. The global descriptor with image subdivision
 173 and histograms of gradients is then

$$T^{HOG} = \sum_{j=1}^f T_j.$$

174 Another improvement is to accumulate the tensor obtained with the video
 175 frame flipped horizontally. Therefore, the HOG3D is computed for each
 176 block, the final tensor is computed (Eq. 3) and simply added to the original

177 frame tensor. This flipped version enforces horizontal gradient symmetries
178 that occur on the video, even those between multiple frames. In our experi-
179 ments (Sec. 4) all HOG3D descriptors are obtained using this improvement.

180 3.3. Combining orientation tensors

181 We propose to concatenate the individual tensors, computed with the
182 optical flow approximation (Sec. 3.1) and HOG3D (Sec. 3.2), to form the
183 final descriptor for the input video:

$$T = \{T^{OF}, T^{HOG}\}. \quad (4)$$

184 Despite other combination methods are possible, concatenation preserves
185 the motion information extracted by each individual descriptor. The informa-
186 tion of those descriptors are complementary and can improve the recognition
187 rate.

188 This descriptor depends only on the video itself, not requiring any re-
189 computation of the previously computed descriptors after the addition of
190 new videos and/or new action categories to the dataset.

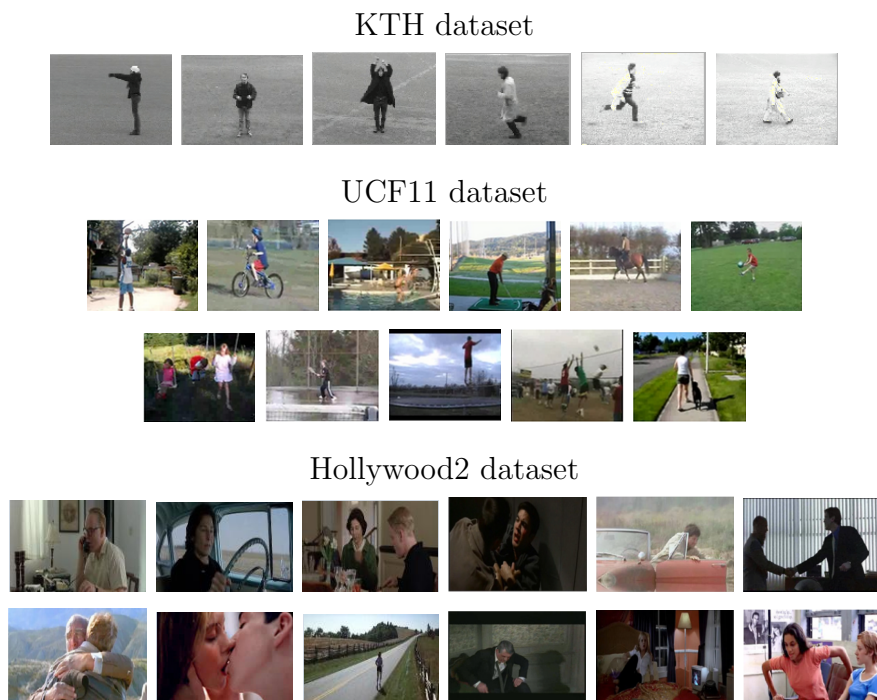
191 It is important to note that the nature of these two tensors is different and,
192 as such, need to be equalized. One possible way is to use a power normal-
193 ization in one of the descriptors. Experimentally, best results were obtained
194 by normalizing the HOG3D tensor: the T^{HOG} descriptor in Equation 4 has
195 all of its elements a_k adjusted to a_k^γ , $\gamma \in]0, 1]$.

196 4. Experimental results

197 We compute the optical flow using the method described by Augereau
198 et al [19]. This method was chosen because we found experimentally that it

199 computes a more regular optical flow than the one computed by the standard
 200 Lucas-Kanade [20]. The descriptor is evaluated by a classification of KTH,
 201 UCF11 and Hollywood2 datasets (Figure 1), using a SVM classifier.

Figure 1: Examples of videos from KTH, UCF11 and Hollywood2 datasets.



202 *4.1. KTH Dataset*

203 The KTH actions dataset [5] consists of six human action classes: walk-
 204 ing, jogging, running, boxing, waving, and clapping. Each action class is
 205 performed several times by 25 subjects. The sequences were recorded in four
 206 different scenarios: outdoors, outdoors with scale variation, outdoors with
 207 different clothes and indoors. The background is homogeneous and static in

208 most sequences. In total, the data consists of 2391 video samples. We use
 209 the same evaluation protocol of the original paper [5], as [1].

210 The best optical flow descriptor for KTH dataset was obtained with a
 211 sliding window with fixed dimensions put around the subject who is doing
 212 the action. The center of mass of global optical flow gives the center of the
 213 window. It works for KTH scenes because they have only one person acting
 214 and a nearly static background. Table 1 shows the recognition rates for this
 215 descriptor using a sliding window of 60x100 pixels. The best recognition rate
 216 was 87.8% with polynomials of degree 8, leading to a descriptor with 16290
 217 elements.

218 In [4] is reported that the best result is achieved with a grid 8x8 and
 219 128 bins obtaining 92.0% of recognition rate. Thus, we concatenate our
 220 optical flow tensor descriptor with this HOG3D to form a new global motion
 221 descriptor. Table 2 shows the recognition rates for several degrees. The best
 222 recognition rate was 93.2% with polynomials of degree 5 (3670 elements)
 223 concatenated with a HOG3D of 128 bins (8256 elements). The confusion
 224 matrix is presented in Table 3.

Table 1: Recognition rates of KTH dataset for several degrees using only the optical flow descriptor with a sliding window with dimensions 60x100.

Degree	1	2	3	4	5	6	7	8	9	10
Rate (%)	81.8	85.5	85.9	86.7	85.8	87.6	87.6	87.8	87.5	87.1

225 4.2. UCF11 Dataset

226 The UCF11 dataset [6] consists in 11 action categories: basketball shoot-
 227 ing, biking/cycling, diving, golf swinging, horse back riding, soccer juggling,

Table 2: Recognition rates of KTH dataset for several degrees using a sliding window with dimensions 60x100 concatenated with a HOG3D [4] of 128 bins with $\gamma = 1$.

Degree	1	2	3	4	5	6	7	8	9	10
Rate (%)	92.6	92.1	92.5	92.7	93.2	92.1	92.0	91.3	91.8	90.6

Table 3: Confusion matrix of KTH dataset for the best result, 93.2% using polynomials of degree 5 and a HOG3D of 8x16 with grid 8x8 [4].

	Box	HClap	HWay	Jog	Run	Walk
Box	95.8	0.0	0.0	0.0	0.0	4.2
HClap	6.2	93.8	0.0	0.0	0.0	0.0
HWay	0.7	4.2	95.1	0.0	0.0	0.0
Jog	0.0	0.0	0.0	90.3	6.9	2.8
Run	0.0	0.0	0.0	14.6	84.7	0.7
Walk	0.0	0.0	0.0	0.7	0.00	99.3

228 swinging, tennis swinging, trampoline jumping, volleyball spiking, and walk-
 229 ing with a dog. We use the same evaluation protocol of the original paper
 230 [6].

231 The sliding window is not interesting for this dataset because its actions
 232 are more complex than those in KTH dataset. Table 4 shows the recognition
 233 rates for several degrees using the optical flow tensor descriptor. The best
 234 recognition rate was 57.8% with polynomials of degree 12 (66430 elements).

Table 4: Recognition rates of UCF11 dataset for several degrees using only the optical flow descriptor.

Degree	1	2	3	4	5
Rate (%)	34.4	45.4	49.9	50.7	53.7
Degree	6	7	8	9	10
Rate (%)	56.1	56.3	56.5	56.3	57.4
Degree	11	12	13	14	15
Rate (%)	57.3	57.8	56.9	57.8	56.6

235 Table 5 presents recognition rates for several parameter sets for HOG3D
 236 descriptor (3.2). The best recognition rate was 68.9% using a grid 32x32 and
 237 a HOG3D of 128 bins.

238 Table 6 shows the recognition rates obtained with the proposed descriptor
 239 with a grid of 32x32 and a HOG3D of 8x16 [4]. A power normalization with
 240 $\gamma = 0.2$ was applied on the final HOG3D tensor. The best recognition rate
 241 was 72.7%, concatenating the HOG3D with polynomials of degree 13 (88410
 242 elements). The confusion matrix for this result is presented in Table 7.

Table 5: Recognition rates of UCF11 dataset for several parameter sets for the HOG3D descriptor [4].

Parameters	Rate (%)
Grid 4x4 HOG3D 8x16	65.0
Grid 8x8 HOG3D 8x16	67.5
Grid 16x16 HOG2D 8x16	68.4
Grid 32x32 HOG3D 8x16	68.9

Table 6: Concatenating the optical flow tensor descriptor with a grid of 32x32 and a HOG3D of 8x16 [4] for UCF11 dataset. A power normalization was applied on the final HOG3D tensor with $\gamma = 0.2$.

Degree	1	2	3	4	5
Rate (%)	69.3	68.0	70.0	70.0	71.2
Degree	6	7	8	9	10
Rate (%)	70.7	71.8	71.5	71.4	72.2
Degree	11	12	13	14	15
Rate (%)	71.9	72.5	72.7	72.4	71.8

243 4.3. Hollywood2 Dataset

244 The Hollywood2 dataset [7] consists of 12 action classes: answering the
 245 phone, driving car, eating, fighting, getting out of the car, hand shaking,
 246 hugging, kissing, running, sitting down, sitting up, and standing up. We
 247 use the same evaluation protocol of the original paper [7]. The performance
 248 is evaluated by computing the average precision (AP) for each of the ac-
 249 tion classes. For the individual optical flow descriptor, the best results are

Table 7: Confusion matrix of UCF11 dataset for the best result, 72.7% using polynomials of degree 13 and a HOG3D of 8x16 with grid 32x32 [4].

	Bike	Dive	Golf	Juggle	Jump	Ride	Shoot	Spike	Swing	Tennis	WDOg
Bike	72.4	0.7	1.0	0.0	0.7	7.5	0.0	1.0	3.0	1.0	12.8
Dive	0.0	89.2	3.2	0.0	0.0	1.8	0.6	1.2	0.6	0.7	2.7
Golf	0.0	4.9	86.8	3.0	0.0	0.0	0.0	0.0	2.7	1.0	1.6
Juggle	0.0	0.7	6.1	62.9	7.2	1.6	2.8	1.9	6.8	5.2	4.7
Jump	0.0	0.0	3.0	8.3	74.1	0.0	0.7	0.7	10.6	0.0	2.7
Ride	3.2	0.0	0.0	0.0	0.0	81.1	0.7	0.9	1.0	0.7	12.4
Shoot	7.0	3.4	0.0	8.6	2.2	0.0	55.4	14.5	1.5	3.8	3.6
Spike	0.0	3.0	0.0	1.7	0.0	1.0	5.5	84.9	3.0	0.0	1.0
Swing	6.8	0.8	3.9	2.5	4.5	2.3	2.0	1.1	65.9	1.0	9.2
Tennis	1.9	0.0	5.5	4.1	5.1	1.2	7.0	2.6	0.6	69.7	2.3
WDOg	7.0	1.0	1.8	3.4	0.0	17.9	1.8	0.7	5.2	3.9	57.4

250 achieved with a Gaussian kernel. For the combination, the triangular kernel
 251 shows the best recognition rates.

252 Table 8 shows the recognition rates for several degrees using the optical
 253 flow tensor descriptor. Similar to UCF11 dataset, the sliding window is not
 254 interesting for this dataset because the actions are more complex than on
 255 KTH dataset. We can see that the recognition rates achieved are very low.
 256 In fact, the summation of tensors will tend to be an isotropic tensor because
 257 there are a lot of different motions happening at the same time in the scenes.
 258 The best recognition rate was only 15% with polynomials of degree 2 (300
 259 elements).

260 In [4] is reported that the best result is achieved with a grid 4x4 and
 261 128 bins obtaining 34.03% of recognition rate. Thus, we concatenate our
 262 optical flow tensor descriptor with this HOG3D to form a new global motion
 263 descriptor. Table 9 shows the recognition rates for several degrees. The best
 264 recognition rate was 40.3% concatenating the HOG3D with polynomials of
 265 degree 3 (820 elements). The average precision for each class is presented in
 266 Table 10.

Table 8: Recognition rates of Hollywood2 dataset for several degrees using only the optical flow descriptor.

Degree	1	2	3	4	5	6	7	8	9	10
Rate (%)	12.0	15.0	13.2	13.3	12.1	13.0	14.5	13.6	13.0	13.3

267 4.4. Comparison with the state-of-the-art

268 A comparison with the state-of-the-art methods is presented in Table 11.

Table 9: Concatenating the optical flow descriptor with a grid of 4x4 and a HOG3D of 8x16 [4] for Hollywood2 dataset. A power normalization was applied on the final HOG3D tensor with $\gamma = 0.2$.

Degree	1	2	3	4	5	6	7	8	9	10
Rate (%)	39.5	39.9	40.3	40.2	40.3	39.8	40.1	39.7	39.9	40.3

Table 10: Average precision (AP) for each class for the best result on Hollywood2 dataset.

Action	AP (%)
AnswerPhone	26.3
DriveCar	71.4
Eat	37.5
FightPerson	50.4
GetOutCar	32.0
HandShake	18.6
HugPerson	29.1
Kiss	47.9
Run	54.9
SitDown	57.5
SitUp	8.8
StandUp	48.7
Mean	40.3

Table 11: Comparison with state-of-the-art for KTH, UCF11 and Hollywood2 datasets.

KTH		UCF11		Hollywood2	
Laptev et al [11]	72.0	Perez et al [4]	68.9	Perez et al [4]	34.0
Laptev et al [2]	91.8	Wang et al [1]	84.2	Laptev et al [2]	45.2
Solmaz et al [12]	92.0			Kobayashi and Otsu [9]	47.7
Zhen and Shao [8]	92.0			Wang et al [1]	58.3
Perez et al [4]	92.0				
Wang et al [1]	94.2				
Kobayashi and Otsu [9]	95.6				
Our Method	93.2	Our Method	72.7	Our Method	40.3

269 The proposed method achieves a competitive accuracy with a much sim-
 270 pler global approach, using only the information from optical flow and his-
 271 tograms of gradients, without any bag-of-features strategy [1, 2, 9].

272 In all datasets we improved the performance of the descriptors previously
 273 proposed in [3, 4] and showed better results than other global descriptors
 274 [11, 12].

275 When compared to bag-of-features strategy on KTH dataset, our descrip-
 276 tor shows better performance than those methods which uses HOF and HOG
 277 as features [2]. The addition of more information to the descriptor, as MBH
 278 and trajectory associated with HOF and HOG [1], induces a better recog-
 279 nition than our descriptor. Even though, the recognition rate is very close
 280 with a much simpler approach.

281 For UCF11 and Hollywood2 datasets, we note that using several features
 282 plays an important role and that learning methods improve overall recog-
 283 nition. The performance of our descriptor is lower than these approaches
 284 [1, 2, 9] but is fairly competitive.

285 Thereby, we can conclude that our descriptor aggregates useful informa-
286 tion of optical flow and HOG3D, enhancing the recognition rate. Moreover,
287 our descriptor only depends on the video, no learning method is required.
288 The addition of new videos and/or new action categories with our approach
289 does not require any re-computation or changes to the previously computed
290 descriptors.

291 **5. Conclusion**

292 In this paper, we presented a novel approach for motion description in
293 videos combining optical flow and HOG3D information. It is an effective ap-
294 proach reaching 93.2% of recognition rate with KTH, comparable to the best
295 local and learning-based methods. However, for the UCF11 and Hollywood2
296 datasets we note points of interest play an important role and that learning
297 methods improve overall recognition. Our recognition rate is lower than the
298 approaches based on codebook but is fairly competitive in both datasets.

299 The main advantage of our method is that it reaches good recognition
300 rates depending uniquely on the video. The addition of new videos and/or
301 new action categories with our approach does not require any re-computation
302 or changes to the previously computed descriptors. Finally, it might be
303 valuable in a scenario where no human action classification method solves all
304 application demands.

305 The drawback of our method is that larger and complex video datasets
306 require higher degree polynomials to give good classification results. As a
307 consequence, the number of coefficients increases exponentially leading to
308 high time complexity. In some cases, increasing the degree does not neces-

309 sarily leads to a better classification, such as in Hollywood2 dataset.

310 In order to improve the recognition rate of our descriptors, we intend to
311 further analyze the spectral characteristics of the proposed orientation tensor.
312 Furthermore, we need to study the conditions of divergence and convergence
313 of the tensor accumulation.

314 **6. Acknowledgements**

315 Authors thank to FAPEMIG and CAPES for funding.

316 **References**

- 317 [1] H. Wang, A. Kläser, C. Schmid, L. Cheng-Lin, Action Recognition by
318 Dense Trajectories, in: IEEE Conference on Computer Vision & Pattern
319 Recognition, Colorado Springs, United States, 2011, pp. 3169–3176.
- 320 [2] I. Laptev, M. Marszałek, C. Schmid, B. Rozenfeld, Learning realistic hu-
321 man actions from movies, in: Computer Vision & Pattern Recognition,
322 2008.
- 323 [3] V. F. Mota, E. A. Perez, M. B. Vieira, L. M. Maciel, F. Precioso, P.-H.
324 Gosselin, A tensor based on optical flow for global description of motion
325 in videos, in: SIBGRAPI 2012 (XXV Conference on Graphics, Patterns
326 and Images), Ouro Preto, MG, Brazil, 2012, pp. 298–301.
- 327 [4] E. A. Perez, V. F. Mota, L. M. Maciel, D. Sad, M. B. Vieira, Combining
328 gradient histograms using orientation tensors for human action recog-
329 nition, in: International Conference on Pattern Recognition, 2012, pp.
330 3460–3463.

- 331 [5] C. Schuldt, I. Laptev, B. Caputo, Recognizing human actions: A local
332 svm approach, in: In Proc. ICPR, 2004, pp. 32–36.
- 333 [6] J. Liu, J. Luo, M. Shah, Recognizing realistic actions from videos in the
334 wild, IEEE CVPR, 2009.
- 335 [7] M. Marszałek, I. Laptev, C. Schmid, Actions in context, in: Conference
336 on Computer Vision & Pattern Recognition, 2009.
- 337 [8] X. Zhen, L. Shao, A local descriptor based on laplacian pyramid coding
338 for action recognition, Pattern Recognition Letters (0) (2012) –, in Press.
339 doi:10.1016/j.patrec.2012.10.021.
- 340 [9] T. Kobayashi, N. Otsu, Motion recognition using local auto-correlation
341 of spacetime gradients, Pattern Recognition Letters 33 (9) (2012) 1188
342 – 1195. doi:10.1016/j.patrec.2012.01.007.
- 343 [10] L. Zelnik-manor, M. Irani, Event-based analysis of video, in: In Proc.
344 CVPR, 2001, pp. 123–130.
- 345 [11] I. Laptev, B. Caputo, C. Schuldt, T. Lindeberg, Local velocity-adapted
346 motion events for spatio-temporal recognition, Comput. Vis. Image Un-
347 derst. 108 (2007) 207–229. doi:10.1016/j.cviu.2006.11.023.
- 348 [12] B. Solmaz, S. M. Assari, M. Shah, Classifying web videos using a
349 global video descriptor, Machine Vision and Applications (2012) 1–
350 13doi:10.1007/s00138-012-0449-x.
- 351 [13] M. Druon, Modélisation du mouvement par polynômes orthogonaux :

- 352 application à l'étude d'écoulements fluides, Ph.D. thesis, Université de
353 Poitiers (02 2009).
- 354 [14] O. Kihl, B. Tremblais, B. Augereau, M. Khoudeir, Human activities dis-
355 crimination with motion approximation in polynomial bases, in: IEEE
356 International Conference on Image Processing, Hong-Kong, 2010, pp.
357 2469–2472.
- 358 [15] C.-F. Westin, A tensor framework for multidimensional signal process-
359 ing, Ph.D. thesis, Linköping University, Sweden, S-581 83 Linköping,
360 Sweden, dissertation No 348, ISBN 91-7871-421-4 (1994).
- 361 [16] B. Johansson, G. Farnebeck, G. F. Ack, A theoretical comparison of
362 different orientation tensors, in: Symposium on Image Analysis, SSAB,
363 2002, pp. 69–73.
- 364 [17] A. Kläser, M. Marszałek, C. Schmid, A spatio-temporal descriptor based
365 on 3d-gradients, in: British Machine Vision Conference, 2008, pp. 995–
366 1004.
- 367 [18] D. G. Lowe, Object recognition from local scale-invariant features, in:
368 Proceedings of the International Conference on Computer Vision - Vol-
369 ume 2, ICCV '99, IEEE Computer Society, Washington, DC, USA, 1999.
- 370 [19] B. Augereau, B. Tremblais, C. Fernandez-Maloigne, Vectorial computa-
371 tion of the optical flow in color image sequences., in: Thirteenth Color
372 Imaging Conference, 2005, pp. 130–134.
- 373 [20] B. Lucas, T. Kanade, An iterative image registration technique with an

374 application to stereo vision (ijcai), in: Proceedings of the 7th Interna-
375 tional Joint Conference on Artificial Intelligence (IJCAI 81), 1981, pp.
376 674–679.