

# Optical Engineering

OpticalEngineering.SPIEDigitalLibrary.org

## Energy flow: image correspondence approximation for motion analysis

Liangliang Wang  
Ruifeng Li  
Yajun Fang

**SPIE.**

Liangliang Wang, Ruifeng Li, Yajun Fang, "Energy flow: image correspondence approximation for motion analysis," *Opt. Eng.* **55**(4), 043109 (2016), doi: 10.1117/1.OE.55.4.043109.

# Energy flow: image correspondence approximation for motion analysis

Liangliang Wang,<sup>a,\*</sup> Ruifeng Li,<sup>a</sup> and Yajun Fang<sup>b</sup>

<sup>a</sup>Harbin Institute of Technology, State Key Laboratory of Robotics and System, 92 Xidazhi Street, Harbin, Heilongjiang 150001, China

<sup>b</sup>Massachusetts Institute of Technology, Computer Science and Artificial Intelligence Laboratory, 77 Massachusetts Avenue, Cambridge, Massachusetts 02139-4307, United States

**Abstract.** We propose a correspondence approximation approach between temporally adjacent frames for motion analysis. First, energy map is established to represent image spatial features on multiple scales using Gaussian convolution. On this basis, energy flow at each layer is estimated using Gauss–Seidel iteration according to the energy invariance constraint. More specifically, at the core of energy invariance constraint is “energy conservation law” assuming that the spatial energy distribution of an image does not change significantly with time. Finally, energy flow field at different layers is reconstructed by considering different smoothness degrees. Due to the multiresolution origin and energy-based implementation, our algorithm is able to quickly address correspondence searching issues in spite of background noise or illumination variation. We apply our correspondence approximation method to motion analysis, and experimental results demonstrate its applicability. © The Authors. Published by SPIE under a Creative Commons Attribution 3.0 Unported License. Distribution or reproduction of this work in whole or in part requires full attribution of the original publication, including its DOI. [DOI: [10.1117/1.OE.55.4.043109](https://doi.org/10.1117/1.OE.55.4.043109)]

Keywords: image correspondence approximation; energy map; energy flow; motion analysis.

Paper 151744 received Dec. 8, 2015; accepted for publication Apr. 5, 2016; published online Apr. 29, 2016.

## 1 Introduction

Motion analysis is a very significant topic in computer vision because of its demand in the area of human–computer interaction, video surveillance, intelligent transportation system, and others. As motion is a time-varying quantity reflecting the variation of an object’s status, in contrast to static image analysis, more useful changing information is available via spatial feature comparison between frames for motion analysis.<sup>1,2</sup> Therefore, at the center of motion analysis is to represent different motions according to their dissimilarities in space-time. From this perspective, techniques for analyzing motion can be divided into two categories: spatial dissimilarity-oriented and temporal dissimilarity-oriented methods.

To be definite, we regard spatial dissimilarity-oriented methods as techniques focusing on exploring dissimilarities of image features, and then combine or extend by adding time labels for motion representation. As a good example, Gilbert and Bowden<sup>3</sup> proposed a dense interest points detection algorithm for human action feature extraction, which is further temporally grouped for classification. Recently, spatiotemporal shape template<sup>4–7</sup> for motion representation attracts much attention for its effectiveness; however, the templates rely strongly on spatial shape representation. Similarly, approaches based on bag of spatiotemporal interest points<sup>8–11</sup> has great success in the field of motion analysis for its space-time invariance. Generally, in spite of spatial dissimilarity-oriented methods being very suitable for motion representation where spatial characteristics are obvious, they often fail to extract adequate global relationships of motion.

In contrast, temporal dissimilarity-oriented methods tend to first extract image features, and then focus on exploiting the relationship and dissimilarities between motion frames. Frame

difference is a very direct and useful scheme to express motion temporal dissimilarities. For example, in Ref. 12, motion energy image (MEI) is built up through image difference, based on which motion history image is formulated by fusing MEI for human movement recognition. Moreover, optical flow<sup>13</sup> is another popular temporal dissimilarity-oriented scheme by assuming brightness is constant between adjacent frames. Inspired by optical flow, Liu and Torralba<sup>14</sup> developed scale-invariant feature transform (SIFT) flow using SIFT points substituting raw pixels for dense correspondence analysis, which is further applied for motion field prediction and face recognition. Furthermore, Huang et al.<sup>15</sup> presented a correspondence map-based algorithm which can be employed for object recognition. Generally speaking, temporal dissimilarity-oriented methods cover both global and local features of motion, and many attempts have been made to address the motion analysis problem from the perspective of image correspondence approximation, as it is more accessible and applicable than frame difference techniques in most cases.

Motivated by the aforementioned observations, this paper solves the motion analysis problem by developing an image correspondence approximation scheme called energy flow, which can be used for dissimilarity searching in space-time between temporally adjacent frames. Particularly, our work first generates a multiscale energy map for image spatial effective representation, which allows for image detail preservation while extracting main features. Using energy map, energy flow at each scale is computed by Gauss–Seidel iteration based on the energy invariance constraint as well as global smoothness assumption.<sup>16</sup> Ultimately, we reconstruct an energy flow field on different scales for accurate image correspondence approximation.

The proposed scheme is capable of finding out dissimilarities between two images, which has great prospect in computer vision domain. Compared with optical flow techniques,<sup>13</sup> our algorithm is more reliable and has higher

\*Address all correspondence to: Liangliang Wang, E-mail: [yueyangmeng@163.com](mailto:yueyangmeng@163.com)

tolerance to illumination changes since multiscale energy rather than brightness is employed for pattern flow searching. As the application for motion analysis, our approach is very practical in contrast to SIFT flow<sup>14</sup> and other spatiotemporal representation methods, for its cheap and accessible characteristics.

The remainder of this paper is organized as follows. Sec. 2 gives an overview of related work. In Sec. 3, our energy flow concept is introduced. Section 4 shows the motion analysis results using energy flow. Finally, Sec. 5 concludes this paper.

## 2 Related Work

As energy flow is an image correspondence-based scheme, as well, motion analysis is a very broad topic allied closely with image segmentation, background modeling, tracking, object recognition, and others, we review previous work from three aspects: image correspondence, motion detection, and human action recognition.

### 2.1 Image Correspondence Approximation

Initially, Horn and Schunck<sup>16</sup> proposed an optical flow estimation method to find dense correspondence fields between images. Optical flow is very efficient for small motions, so a great deal of research<sup>13,17,18</sup> following this pipeline has been done for correspondence approximation. However, optical flow makes the brightness constancy assumption and therefore fails to deal with large lighting changes, it also cannot accurately describe the motion region if there is overlap or noise on the brightness layer.

Another popular image correspondence technique is SIFT,<sup>19</sup> which matches the images using sparse points that are robust to geometric and photometric variations on multiple scales. SIFT flow,<sup>14</sup> mentioned earlier, is actually an extension of SIFT by fusing it into optical flow formulation. Unfortunately, SIFT-based algorithms are either computationally consuming or too sparse to achieve precise correspondence approximation. To deal with these shortcomings, Tau and Hassner<sup>20</sup> further seek to propagate image scale information from detected interest points to its neighboring pixels context by considering locations where scales are detected, and then use the context for images separately and within correlated images, which results in more useful features for dense correspondence while keeping the computational burden low. Similarly, Zhang et al.<sup>21</sup> proposed an energy flow equation by replacing the brightness using image temperature features within the Horn–Schunck optical flow framework, which is employed for video segmentation.

Moreover, researchers present many approaches for approximating image correspondence from other points of view, such as Refs. 15 and 22, no matter if they work on pixels or interest points, the dilemma between accuracy and efficiency is challenging especially for wide-range practical applications.

### 2.2 Motion Detection

Broadly speaking, existing work for motion detection can be roughly divided into model-based and appearance-based detections. Model-based methods detect motions by comparing the target with a built model. It is ideal to directly use the

background image<sup>22</sup> without interference as the model if the scenario is static, but more often, using an estimated model from a priori knowledge is more actual, e.g., Gaussian mixture model (GMM)<sup>23</sup> is proposed for dynamic model estimation according to the Gaussian mixture distribution of pixels, which is widely applied for object tracking. In a very recent work, Haines and Xiang<sup>24</sup> further used a Dirichlet process GMM to provide a per-pixel density estimate for background computation. Model-based techniques are quick, but rely strongly on the established model. Appearance-based approaches pay more attention to learn a large number of sample features, and then accomplish motion detection by classification, e.g., histogram of oriented gradient (HOG)<sup>25</sup> is formulated to represent gradient features of an image, according to which, pedestrians can be detected via support vector machines (SVMs) framework.<sup>26</sup> In Ref. 4, a detector named action bank is presented for human motion detection, and on this basis, motion can be accurately localized through SVMs. Tamrakar et al.<sup>10</sup> introduced a bag of SIFT features for complex event detection.

### 2.3 Human Action Recognition

As human action is a very large-volume data digitally, the heart of action recognition is to extract spatiotemporal features<sup>3</sup> to represent actions. Considering the characteristics of action, many action descriptors have been presented, e.g., Derpanis et al.<sup>6</sup> developed a spatial-temporal orientation template generated via three-dimensional Gaussian filtering on raw raw image intensity features for reflecting the dynamics of actions. In Ref. 7, action videos are segmented into spatiotemporal graphs expressing hierarchical, temporal, and spatial relationships of actions, and then a matching algorithm is formulated for action recognition. Additionally, a lot of techniques originated from image correspondence and motion detection are widely applied for action recognition, e.g., Laptev et al.<sup>8</sup> build a spatiotemporal bag of words (BoW) model to represent action interest points consisting of HOG and optical flow features. Furthermore, context of interest points is able to be used for action representation, e.g., in Ref. 27, the action context feature is defined as the relative coordinates of pairwise interest points in space-time, and then GMMs are used to describe the context distributions of interest points.

## 3 Methodology

Our goal is to explore correspondence between images for motion analysis. In this work, a temporal dissimilarity-oriented scheme is presented while the spatial features of images are deep extracted. Given two temporally adjacent frames, we start from building multilayer Laplacian stacks for both, respectively, using Gaussian kernel convolution implementation, and energy map is further established for image feature extraction. We compute the energy flow between two energy maps based on the energy invariance constraint, and energy flow field is reconstructed to approximate the correspondence.

### 3.1 Energy Map

To exploit the local features of an image, the first step of our algorithm is to represent an image  $I$  on multiple scales employing Laplacian stacks. Let  $G(\sigma)$  denote a two-dimensional normalized Gaussian kernel with standard deviation  $\sigma$ ,

and let  $*$  denote the convolution operator, the image  $I$  can be decomposed into a  $m$ -scale ( $m \geq 1$ ) descriptor  $\{L_S(I) | 0 \leq S \leq m\}$ , where

$$L_S(I) = \begin{cases} I - I * G(\sigma) & \text{if } S = 0 \\ I * G(\sigma^S) - I * G(\sigma^{S+1}) & \text{if } S > 0 \end{cases} \quad (1)$$

Despite the fact that Laplacian stacks are able to find out full details as its origin at multiresolution processing, for each subband, it is band limited.<sup>28</sup> Therefore, in order to describe an image more accurately with fewer noises by considering the dissimilarity between different scales, a rectification process is implemented in our work. Based on the Laplacian stacks, and inspired by power maps proposed in Refs. 22 and 28, we establish our energy map according to the absolute value of Laplacian coefficients because the variation produced by difference of Laplacian stacks rather than its orientation is the point of our concern. For  $I$  on the  $S$ 'th scale, we define the transfer energy as

$$T_S(I) = \ln |L_S(I)| * G(\sigma^{S+1}). \quad (2)$$

Here, we transform the absolute value of Laplacian coefficients into logarithmic domain. Since the value of  $|L_S(I)|$  at many pixels is 0, which brings infinitely small quantity impacting the following computation, we make the following revision:

$$T'_S(I) = \ln |L'_S(I)| * G(\sigma^{S+1}), \quad (3)$$

where

$$|L'_S(I)| = \begin{cases} |L_S(I)| & \text{if } |L_S(I)| \neq 0 \\ 1 & \text{if } |L_S(I)| = 0 \end{cases} \quad (4)$$

Then we continue to define the energy map considering both the absolute value of  $L_S(I)$  and the exponent of weighted transfer energy:

$$E_S(I) = |L_S(I)| e^{\lambda T'_S(I)}, \quad (5)$$

where  $\lambda$  is an adaptable parameter. Since the revision process adds noises to  $e^{\lambda T'_S(I)}$  by conserving zeros of  $|L_S(I)|$ , we further modify it using  $P_S(I)$ :

$$P_S(I) = \begin{cases} e^{\lambda T'_S(I)} & \text{if } |e^{\lambda T'_S(I)} - e^{\lambda \rho}| > \epsilon \\ 0 & \text{else} \end{cases}, \quad (6)$$

where  $\epsilon$  is the infinitely small quantity, and  $\rho$  is a parameter determined by image quality.

Finally, energy map is built up as follows:

$$E_S(I) = |L_S(I)| P_S(I). \quad (7)$$

Thus, we can conclude that our energy map is essentially the multilayer Laplacian energy stacks for action spatial feature extraction. Figure 1 shows an example of energy map, it is worth noting that the four layers of energy map are displayed with the same size in spite of actually every backward layer decreases into one-fourth with respect to its forward layer. Additionally, it is worth noting that  $\sigma$  is set as 2,  $m$  is chosen as 4,  $\lambda$  is selected as  $-0.3$ , and  $\rho$  ranges from  $-2$  to  $-0.5$  in our work which are practically proven to work well.

### 3.2 Energy Flow

To extract temporal features between frames, we regard motion as the apparent motion of the energy. Therefore, as we know, there are two smoothness assumptions<sup>13</sup> for optical flow computation: global smoothness<sup>16</sup> which can produce dense optical flow field but fail to describe boundaries and local smoothness<sup>17</sup> which is more robust but often results in sparse motion description. Considering the advantage of the energy map on depicting boundaries, and motivated by Horn-Schunck optical flow formulation,<sup>16</sup> we make the assumption that the spatial energy at two continuous times on the same scale is equal using global smoothness assumption. Moreover and likewise, we define "energy conservation law" as follows: let  $E_S(x, y, t)$  denote the energy of a pixel  $(x, y)$  of an image  $I$  at time  $t$  on the  $S$ 'th scale, after a small time interval  $\delta t$  at the point  $(x + \delta x, y + \delta y)$ , we thus define

$$E_S(x, y, t) = E_S(x + \delta x, y + \delta y, t + \delta t). \quad (8)$$

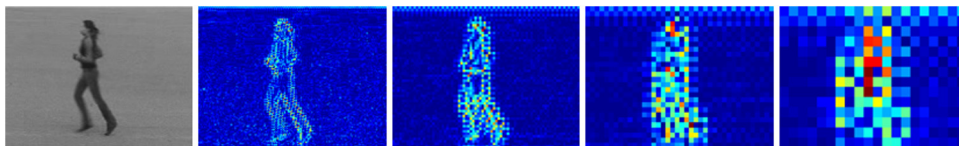
Based on this assumption, we expand the above equation using Taylor series:

$$E_S(x, y, t) + \delta x \frac{\partial E_S}{\partial x} + \delta y \frac{\partial E_S}{\partial y} + \delta t \frac{\partial E_S}{\partial t} + o(2) = E_S(x, y, t), \quad (9)$$

where  $o(2)$  denotes the first-order of infinitely small quantity. Then dividing  $\delta t$  on both sides of Eq. (9), and as  $\delta t \rightarrow 0$ , we can get

$$\frac{\partial E_S}{\partial x} \frac{dx}{dt} + \frac{\partial E_S}{\partial y} \frac{dy}{dt} + \frac{\partial E_S}{\partial t} = 0. \quad (10)$$

Here, we define the velocity of a pixel as  $\nu_S = (\nu_{Sx}, \nu_{Sy})$  and  $\nu_{Sx} = (dx/dt)$ ,  $\nu_{Sy} = (dy/dt)$ , so we can get the energy flow constraint equation:



**Fig. 1** An example of energy map. The first column is the initial action image, and from the second to the fifth columns are the energy maps on four layers, respectively.

$$\frac{\partial E_S}{\partial x} \nu_{Sx} + \frac{\partial E_S}{\partial y} \nu_{Sy} + \frac{\partial E_S}{\partial t} = 0. \quad (11)$$

Then we describe energy flow using the energy flow field descriptor  $\nu_S = (\nu_{Sx}, \nu_{Sy})$ , which can be computed by minimizing the following objective function:

$$\nu_S = \underset{\nu_S}{\operatorname{argmin}} \sum_x \left[ \alpha_1 \left\| \frac{\partial E_S}{\partial x} \nu_{Sx} + \frac{\partial E_S}{\partial y} \nu_{Sy} + \frac{\partial E_S}{\partial t} \right\|^2 + \alpha_2 \left( \left\| \frac{\partial \nu_S}{\partial x} \right\|^2 + \left\| \frac{\partial \nu_S}{\partial y} \right\|^2 \right) \right], \quad (12)$$

where  $\alpha_1$  ( $\alpha_1 \neq 0$ ) and  $\alpha_2$  are respectively the weights for data and smoothness terms indicating the energy invariance and global smoothness assumption.<sup>13</sup> Likewise, the ratio  $\alpha_2/\alpha_1$  is determined by the image quality.<sup>29</sup>

Utilizing the Gauss–Seidel iteration, Eq. (12) can be solved as follows:

$$\nu_{Sx}^{k+1} = \bar{\nu}_{Sx}^k - \frac{\frac{\partial E_S}{\partial x} \bar{\nu}_{Sx}^k + \frac{\partial E_S}{\partial y} \bar{\nu}_{Sy}^k + \frac{\partial E_S}{\partial t}}{\left(\frac{\alpha_2}{\alpha_1}\right)^2 + \left(\frac{\partial E_S}{\partial x}\right)^2 + \left(\frac{\partial E_S}{\partial y}\right)^2} \frac{\partial E_S}{\partial x}, \quad (13)$$

$$\nu_{Sy}^{k+1} = \bar{\nu}_{Sy}^k - \frac{\frac{\partial E_S}{\partial x} \bar{\nu}_{Sx}^k + \frac{\partial E_S}{\partial y} \bar{\nu}_{Sy}^k + \frac{\partial E_S}{\partial t}}{\left(\frac{\alpha_2}{\alpha_1}\right)^2 + \left(\frac{\partial E_S}{\partial x}\right)^2 + \left(\frac{\partial E_S}{\partial y}\right)^2} \frac{\partial E_S}{\partial y}, \quad (14)$$

where  $k$  ( $k \geq 0$ ) denotes the iteration number, and in our work,  $k$  is set as 100 to guarantee both efficiency and accuracy.

### 3.3 Energy Flow Field Reconstruction

Therefore, after iteration via Eqs. (13) and (14), from the macropoint of view, for two frames, we can get a final energy flow field sequence abbreviated as  $\{V_S = (\nu_{Sx}^{k+1}, \nu_{Sy}^{k+1}) | 0 \leq S \leq m\}$  on multiple scales. Because for high-pass scales, the energy map averages response over a larger region of the image,<sup>28</sup> to represent the details produced by tiny variation during the time interval  $\delta t$  and to guarantee the avoidance of noise simultaneously, we reconstruct energy flow field on the velocity layer rather than on the energy map layer for expressing image correspondence relationship using  $V_0$ , which can be computed by iteration as follows:

$$V_S = \begin{cases} V_S & \text{if } S = m \\ V_S + \frac{S+1}{S+2} V_{S+1} * G(\sigma^{S+1}) & \text{if } S < m \end{cases}. \quad (15)$$

## 4 Experiments

As our algorithm is an image correspondence-based scheme for dissimilarity searching between adjacent frames, to better reveal its performance, we test our algorithm for motion analysis from two facets: motion detection and human action recognition. Also, we believe that our method can be used in more areas.

### 4.1 Motion Detection

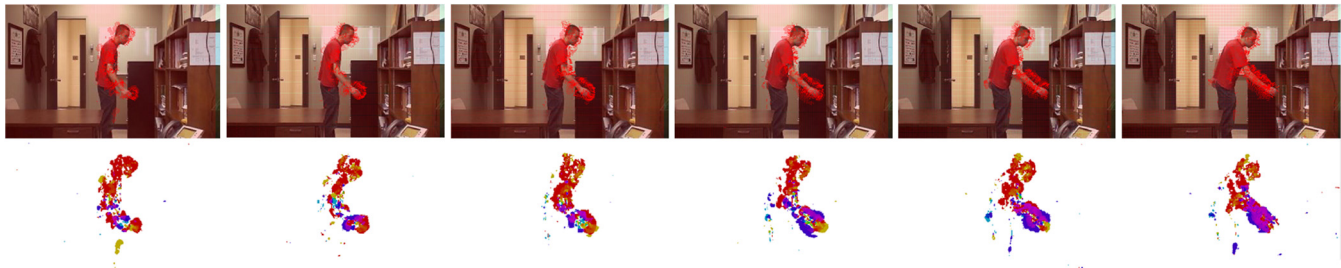
We verify our algorithm for motion field prediction using frames from ChangeDetection.NET 2014 change detection database<sup>30</sup> without additional processing. ChangeDetection.Net 2014 is a very complex benchmark for event and motion detection consisting of 31 videos depicting indoor and outdoor scenes with boats, cars, trucks, and pedestrians.

To visualize energy flow velocities, we display oriented arrows of energy flow field from the previous frame to the current status, and one velocity vector in  $2 \times 2$  or  $5 \times 5$  pixels is set to be visible and the magnifying scale factor of arrows is 5 or 10 determined by image quality. As well, we utilize color maps to show energy flow field regions according to the value of  $\arctan(\nu_{0x}^{101}/\nu_{0y}^{101})$  at each pixel, it is worth noting that the previous frames are often not given but can be inferred from our visualizations which reflect motion variations.

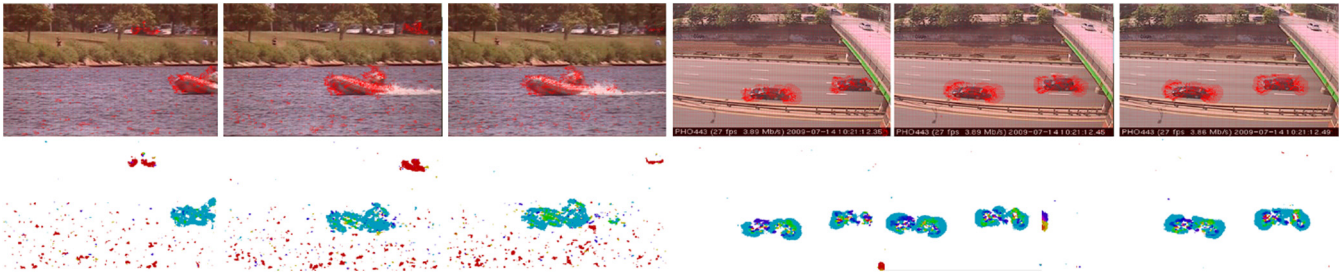
Figure 2 gives the example results of continuous human motion detection in a relatively static scenario, the grabbing motion is slow, a large part of the human body is not moving, and a small part moves slightly. From detection results, we can see that our algorithm is able to depict moving parts effectively with little noises and the boundaries are precisely detected. Also, the overlap within motions is successfully addressed.

Figure 3 gives the example results of motion detection in the lake and highway scenarios. The lake scenario is very challenging as it includes motions of a man driving a boat, a black car's motion far away from lens, and the lake water flow. However, we deal with the case well and the main motion variations are detected. For the highway scenario, the motion is very quick leading to big variations, and it is shown from the results that the motions are localized very accurately, but a part of the car's body is disregarded.

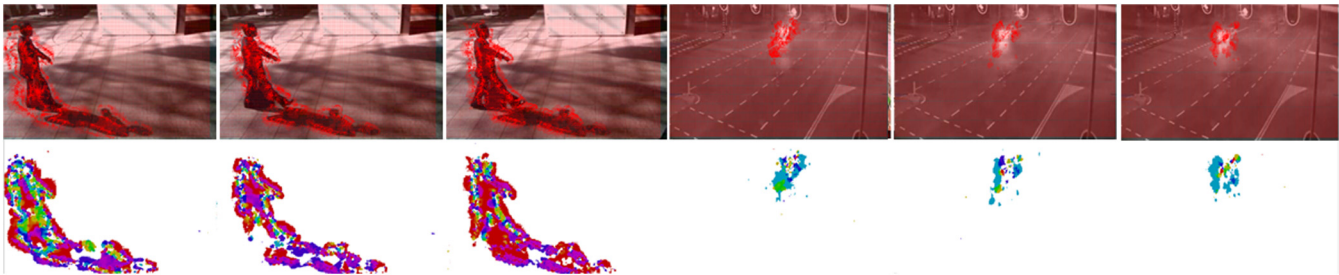
Figure 4 gives the example results of motion detection in a shadow scenario and at night. The results of pedestrian



**Fig. 2** Example results of human motion detection. Images in the top row are continuous frames with oriented arrows describing energy flow velocities from its previous frame to the current status, and the bottom row shows the color maps. The previous frame of the first image is not given.



**Fig. 3** Example results of motion detection in the lake and highway scenarios. Images in the top row are representative frames with oriented energy flow arrows, and the bottom row shows the corresponding color maps.



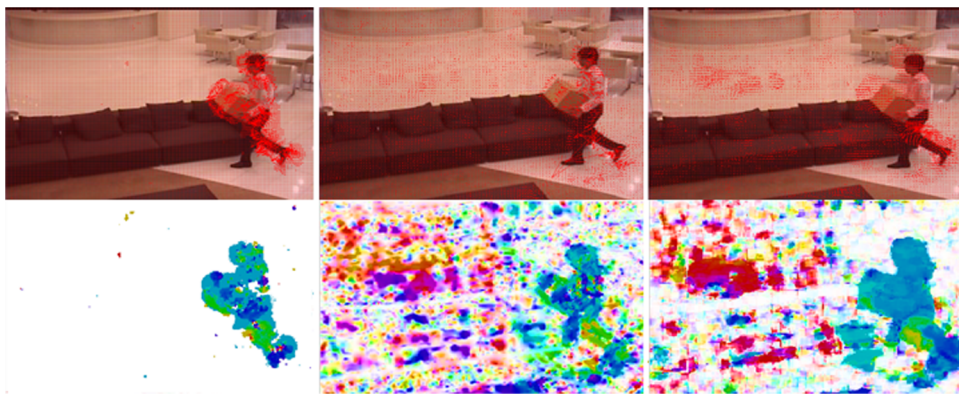
**Fig. 4** Example results of motion detection in a shadow scenario and at night. Images in the top row are representative frames with oriented energy flow arrows, and the bottom row shows the corresponding color maps.

detection with shadow are promising since we are aimed at motion detection instead of detecting pedestrians. As motion detection at night with illumination changes, our approach is also very robust.

As a comparison, Fig. 5 compares our method with optical flow methods of Refs. 16 and 17 using color map on examples from ChangeDetection. Net 2014. Between two frames of human walking with a box, the main motion lies on wiggling of the foot behind and translation of the upper body, and from the results, we can see that the method of Ref. 16 cannot describe boundaries accurately and is heavily damaged by noise; the method of Ref. 17 enlarges the motion part and is not reliable in contrast to our approach. Moreover, the average running time of our algorithm for 10 times is 0.039 s, compared with 17.143 and 0.918 s by

methods of Refs. 17 and 16. We implement all the experiments in MATLAB on an i5-core PC with a 6 GB RAM.

Moreover, to further validate our approach, we compare its overall results with another four methods for motion detection on ChangeDetection. Net 2014 shown in Table 1. We select three popular metrics for evaluation: recall ( $Re = N_{tp}/(N_{tp} + N_{fn})$ ), false positive rate [ $Fpr = N_{fp}/(N_{fp} + N_{tn})$ ], and precision [ $Pr = N_{tp}/(N_{tp} + N_{fp})$ ], which are determined by the number of true positives ( $N_{tp}$ ), true negatives ( $N_{tn}$ ), false positives ( $N_{fp}$ ), and false negatives ( $N_{fn}$ ). From the comparison, we can see that our method outperforms popular optical flow methods,<sup>16,17</sup> and can handle real-time action detection well in contrast to GMM<sup>23</sup> and background modeling<sup>24</sup> based algorithms.



**Fig. 5** Example results of pedestrian detection on ChangeDetection. Net 2014 database using energy flow and optical flow methods (respectively proposed by Horn and Mahmoudi).

**Table 1** Overall action detection results of different algorithms on ChangeDetection. Net 2014 database.

Method	Re	Fpr	Pr
Horn and Schunck <sup>16</sup>	0.68	0.030	0.56
Mahmoudi et al. <sup>17</sup>	0.69	0.022	0.65
Stauffer and Grimson <sup>23</sup>	0.62	0.025	0.60
Haines and Xiang <sup>24</sup>	0.78	0.013	0.74
Our algorithm	0.76	0.009	0.81

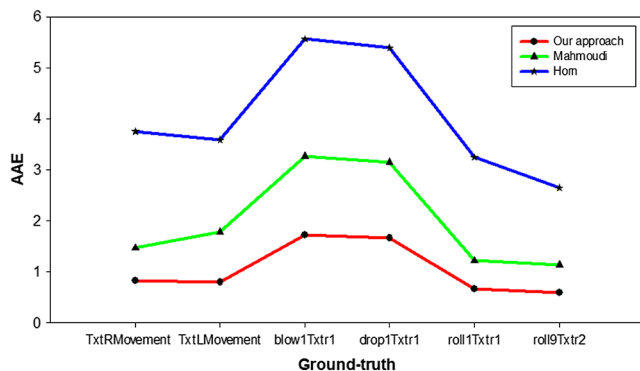
To evaluate our energy flow errors, we compute average angular errors (AAE) of energy flow using ground-truth sequences (“TxtrMovement,” “TxtrLMovement,” “blow1Txtr1,” “drop1Txtr1,” “roll1Txtr1,” and “roll9Txtr2”) from University College London (UCL) database<sup>18</sup> by averaging all the AE calculated by the following equation:

$$AE = \cos^{-1} \left[ \frac{1 + \nu_{0x}^{101} \times \nu_x + \nu_{0x}^{101} \times \nu_y}{\sqrt{1 + (\nu_{0x}^{101})^2 + (\nu_{0x}^{101})^2} \sqrt{1 + \nu_x^2 + \nu_y^2}} \right], \tag{16}$$

where  $(\nu_x, \nu_y)$  denotes the velocity of ground-truth at  $(x, y)$ . As a comparison, the AAE of Refs. 16 and 17 is also shown in Fig. 6.

### 4.2 Human Action Recognition

For action recognition issue, we select sequences from Kungl Tekniska Högskolan (KTH) (2391 video clips including 6 actions performed by 25 persons)<sup>26</sup> and human metabolome database (HMDB) (6849 video clips divided into 51 action categories)<sup>31</sup> action databases. Using energy flow field between two frames as features, we cluster 100k features of the energy flow field descriptors using  $k$ -means algorithm by setting  $k$  as 4000, then encode them via a BoW as depicted in Ref. 10, and finally we classify actions under SVMs framework with radial basis function kernel which is practically demonstrated robust. For each action, same as in Refs. 26 and 31, we select 16 persons’ video clips for training and the rest for testing on KTH, while we choose 70 video clips for training and 30 video clips for testing on HMDB.



**Fig. 6** AAE of different methods using ground-truth from UCL database.

boxing	100	0	0	0	0	0
handclapping	0	98.9	1.1	0	0	0
handwaving	0	0.8	99.2	0	0	0
running	0	0	0	88.1	3.4	8.5
walking	0	0	0	2.0	90.2	7.8
jogging	0	0	0	3.8	10.7	85.5
	boxing	handclapping	handwaving	running	walking	jogging

**Fig. 7** Confusion matrix of our algorithm on KTH dataset.

**Table 2** ARR of different methods on KTH database.

Method	ARR
Schuldt and Caputo <sup>26</sup>	71.7
Derpanis et al. <sup>6</sup>	89.34
Laptev et al. <sup>8</sup>	91.8
Iosifidis and Pitas <sup>11</sup>	92.13
SIFT + BoW	85.46
Optical flow + BoW	79.59
Our algorithm	93.65

Figure 7 gives the confusion matrix using our method on KTH database, and the average recognition rate (ARR) reaches 93.65%. Table 1 compares our algorithm with other related works.<sup>6,11,23,26</sup> In the meanwhile, with the same settings except using SIFT<sup>31</sup> and optical flow features<sup>17</sup> replacing our energy flow features, we get the ARR which is shown in Table 2.

**Table 3** ARR of different methods on HMDB database.

Method	ARR
Kuehne <sup>23</sup>	23.18
Sadanand and Corso <sup>4</sup>	26.9
Cao et al. <sup>5</sup>	27.84
SIFT + BoW	21.56
Optical flow + BoW	16.87
Our algorithm	27.92





29. T. Xue et al., "Refraction wiggles for measuring fluid depth and velocity from video," in *Proc. Euro. Computer Vision* (2014).
30. Y. Wang et al., "Cdnnet 2014: an expanded change detection benchmark dataset," in *Proc. IEEE Conf. Computer Vision Pattern Recognition*, pp. 387–394 (2014).
31. H. Kuehne et al., "HMDB: a large video database for human motion recognition," in *Proc. IEEE Computer Vision*, pp. 2556–2563 (2011).
32. R. P. E. Rosten and T. Drummond, "Faster and better: a machine learning approach to corner detection," *IEEE Trans. Pattern Anal. Mach. Intell.* **32**, 105–119 (2010).
33. J. S. Beis and D. G. Lowe, "Shape indexing using approximate nearest-neighbour search in high-dimensional spaces," in *Proc. IEEE Conf. Computer Vision Pattern Recognition*, p. 1000 (1997).

**Liangliang Wang** is a PhD student at State Key Laboratory of Robotics and System, Harbin Institute of Technology. He received his BS and MS degrees in mechatronics engineering from Harbin Engineering University in 2009 and Harbin Institute of Technology in 2011, respectively. He visited CSAIL, Massachusetts Institute of

Technology, between 2014 and 2015 as a visiting student hosted by Prof. Berthold K. P. Horn. His current research interests include machine vision and pattern recognition.

**Ruifeng Li** is a professor and the vice director at State Key Laboratory of Robotics and System, Harbin Institute of Technology. He received his PhD from Harbin Institute of Technology in 1996. He is a member of the Chinese Association for Artificial Intelligence and the president of Heilongjiang Province Institute of Robotics. His current research interests include artificial intelligence and robotics.

**Yajun Fang** received her PhD from CSAIL, Massachusetts Institute of Technology in 2010. Subsequently, she joined the Martinos Image Research Center at Harvard University as a postdoc. Currently, she is a research scientist at Intelligent Transportation System Center at the Massachusetts Institute of Technology. Her research field is computer vision and intelligent transportation systems.