

# Inferring Ongoing Human Activities Based on Recurrent Self-Organizing Map Trajectory

Qianru Sun  
qianrusun@sz.pku.edu.cn

Hong Liu  
hongliu@pku.edu.cn

Engineering Lab on Intelligent Perception for Internet of Things(ELIP), Shenzhen Graduate School  
Peking University, CHINA  
Key Laboratory of Machine Perception(Ministry of Education)  
Peking University, CHINA

Automatically inferring ongoing activities is to enable the early recognition of unfinished activities, which is quite meaningful for applications, such as online human-machine interaction and security monitoring. State-of-the-art methods use the spatio-temporal interest point (STIP) based features as the low-level video description to handle complex scenes [1, 2, 3]. While the existing problem is that typical bag-of-visual words (BoVW) focuses on feature distribution but ignores the inherent contexts in sequences, resulting in low discrimination when directly dealing with limited observations. To solve this problem, the Recurrent Self-Organizing Map (RSOM) [4], which was designed to process sequential data, is novelly adopted in this paper for the high-level representation of ongoing activities. The innovation lies that observed features and their spatio-temporal contexts are encoded in a trajectory of the pre-trained RSOM units. Additionally, a combination of Dynamic Time Warping (DTW) distance and Edit distance, named DTW-E, is specially proposed to measure the structural dissimilarity between RSOM trajectories.

**RSOM Trajectory:** Since the RSOM constitutes a direct extension of SOM, we start from SOM. SOM is to map the data from an input space  $V_I$  onto a lower dimensional space  $V_L$  (a map) in such way that the topological relationships in  $V_I$  are preserved and the SOM units approximate closely the probability density function of  $V_I$ . Suppose each unit  $i$  in SOM is associated with a weight vector  $\mathbf{w}_i = [w_{i1}, w_{i2}, \dots, w_{in}]^T \in \mathfrak{R}^n$  with the same dimension as the input vector  $\mathbf{x} = [x_1, x_2, \dots, x_n]^T \in \mathfrak{R}^n$ . Learning process that leads to self-organization on a map can be summarized as,

(i) The feature vector  $\mathbf{x}(t)$  is input, then its best matching unit ( $bm_u$ ) on the map is found by computing the minimum distance as:

$$bm_u = \arg \min_{i \in V_L} \{\|\mathbf{x}(t) - \mathbf{w}_i(t)\|\} \quad (1)$$

(ii) The winner  $bm_u$  and its neighbors on the map have their weights  $\mathbf{w}_i(t)$  updated towards  $\mathbf{x}(t)$  as:

$$\mathbf{w}_i(t+1) = \mathbf{w}_i(t) + \alpha(t) \cdot N_{bm_u, i} \cdot \|\mathbf{x}(t) - \mathbf{w}_i(t)\| \quad (2)$$

where  $\|\cdot\|$  denotes the Euclidean norm,  $\alpha(t) = \alpha_i \cdot (\alpha_f / \alpha_i)^{T(i)/T_{max}} \in [0, 1]$  is the learning rate, where the  $\alpha_i$  and  $\alpha_f$  denote the initial rate and final rate.  $T(i) = \{1, 2, \dots, T_{max}\}$  where  $T_{max}$  is iteration number.  $N_{bm_u, i}$  is called neighborhood function and defined over the units on the map. Typically,  $N_{bm_u, i} = \exp\{-\|r_{bm_u} - r_i\|^2 / 2\sigma^2\}$ , where  $r_{bm_u} \in \mathfrak{R}^2$  and  $r_i \in \mathfrak{R}^2$  are the location vectors of unit  $bm_u$  and  $i$  on the map, and  $\sigma$  defines the Gaussian kernel width.

SOM is not originally designed to accommodate the time series, its temporal extension RSOM is hence adopted here to learn the temporal contexts in activity sequences. It is to utilize both the feature vectors before  $\mathbf{x}(t)$  and  $\mathbf{x}(t)$  itself to search the best matching unit of  $\mathbf{x}(t)$ . This is done by associating the following recursive equation to each unit  $i$  to compute the difference vector  $\mathbf{y}_i(t)$ :

$$\mathbf{y}_i(t) = \lambda \cdot \|\mathbf{x}(t) - \mathbf{w}_i(t)\| + (1 - \lambda) \cdot \mathbf{y}_i(t-1) \quad (3)$$

where  $0 < \lambda < 1$  is a factor determining the influence of earlier difference vectors on the current  $\mathbf{x}(t)$ . When  $\lambda$  is close to 0, the system of Eq. (3) involves a heavy backward memory, whereas when  $\lambda$  is near to 1, Eq. (3) describes a slight memory. Now equations of RSOM for searching  $bm_u$  and adapting weights are as follows,

$$bm_u = \arg \min_{i \in V_L} \{\|\mathbf{y}_i(t)\|\} \quad (4)$$

$$\mathbf{w}_i(t+1) = \mathbf{w}_i(t) + \alpha(t) \cdot N_{bm_u, i} \cdot \|\mathbf{y}_i(t)\| \quad (5)$$

Then, we introduce how to map the input feature sequence to a time-varying trajectory of  $bm_u$ . Supposing that an  $M \times M$  map is learned after a fixed number of iterations using Eq. (3)(4)(5). For simplification, the one-dimensional value  $b$  of the original coordinate  $bm_u \in \mathfrak{R}^2$ , i.e.,

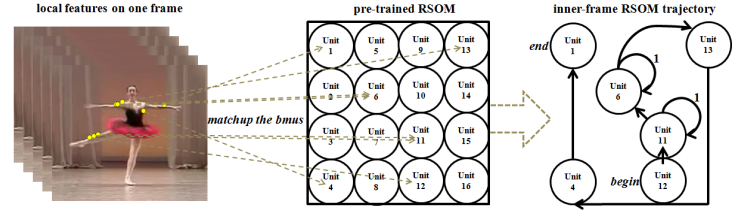


Figure 1: Illustrations of a pre-trained  $4 \times 4$  RSOM and the assumed trajectory (12, 11, 11, 6, 6, 13, 4, 1) of local features (yellow points). Note that the  $bm_u$  order in the trajectory is based on their image locations (detected order): from left to right, then from up to down.



Figure 2: Samples of UT-Interaction (a, b) and Rochester Activities dataset (c).

$b = bm_u(2) \times M + bm_u(1) \in [1, M^2]$ , is used as the location index in the final trajectory. During video input at time  $t$ , STIP based local features are first extracted on the current frame based on above equations. Then feature vectors search their  $bm_u$  on the map, and compose an inner-frame trajectory  $\mathbf{b}$ . Finally,  $\mathbf{b}$  is used to generate the inter-frame trajectory  $\mathbf{Trj}$ .

$$\mathbf{b}_f = \{b_k | k = 1, 2, \dots, K(f)\}; \mathbf{Trj}(t) = \{\mathbf{b}_f | f = 1, 2, \dots, F(t)\} \quad (6)$$

where  $K(f)$  is the number of local features on the  $f$ th frame and it changes with different frames.  $F(t)$  is the frame number until time  $t$ .  $\mathbf{Trj}(t)$  is thus used as a high-level representation of the current observed activity.

**DTW-E distance:** The structure of RSOM trajectory is clear and special that each subset on one frame (inner-frame) contains the human shape information and the whole sequence (inter-frame) contains the long-range temporal relationships. Therefore, how to reasonably measure the likelihood between RSOM trajectories for pattern classification arises another problem. To solve this, a hierarchical distance based on the combination of DTW distance and Edit distance, named DTW-E, is specially defined (implementation of this algorithm is described in our paper).

**Experiments and Conclusions:** Two real-world datasets (Figure 2) with different characteristics, complex scenes [5] and inter-class ambiguities [6], serve as sources of data for evaluation. Experimental results based on kNN classifiers confirm that our approach can infer ongoing human activities at any stage with high accuracies.

Our conclusion is that our approach can make efficient inferences even with complex scenes and inter-class ambiguities. It hence confirms the ability of RSOM trajectory to extract sufficient discrimination. Moreover, the RSOM trajectory with the advantage of before-after independence is proved more suitable to the recognition of unfinished patterns.

- [1] P. Dollár, V. Rabaud, G. Cottrell, S. Belongie. Behavior recognition via sparse spatio-temporal features. *VS-PETS*, pages 65-72, 2005.
- [2] H. Wang, M. M. Ullah, A. Klaser, I. Laptev, C. Schmid. Evaluation of local spatio-temporal features for action recognition. In *Proc. BMVC*, pages 124.1-124.11, 2009.
- [3] P. Scovanner, S. Ali, M. Shah. A 3-Dimensional SIFT Descriptor and its Application to Action Recognition. *ACM Conf. Multimedia*, pages 357-360, 2007.
- [4] M. Varsta, José del R. Millán, J. Heikkonen. A Recurrent Self-Organizing Map for Temporal Sequence Processing. *LNCS 1327*, pages 421-426, 1997.
- [5] M. S. Ryoo, J. K. Aggarwal. UT-Interaction Dataset. *ICPR Contest on Semantic Description of Human Activities (SDHA)*, 2010.
- [6] R. Messing, C. Pal, H. Kautz, Activity recognition using the velocity histories of tracked keypoints, In *Proc. ICCV*, pages 104-111, 2009.