

An online spatio-temporal tensor learning model for visual tracking and its applications to facial expression recognition

Sheheryar Khan¹, Guoxia Xu¹, Raymond Chan², and Hong Yan¹

¹*Department of Electronic Engineering, City University of Hong Kong, Kowloon, Hong Kong*

²*Department of Mathematics, Chinese University of Hong Kong, Shatin, Hong Kong*

E-mail address: shehekhan2-c@my.cityu.edu.hk (S. Khan), guoxiaxu@cityu.edu.hk (G. Xu), rchan@math.cuhk.edu.hk (R. Chan), h.yan@cityu.edu.hk (H. Yan)

Abstract

Robust visual tracking remains a technical challenge in real-world applications, as an object may involve many appearance variations. In existing tracking frameworks, objects in an image are often represented as vector observations, which discounts the 2-D intrinsic structure of the image. By considering an image in its actual form as a matrix, we construct the 3rd order tensor based object representation to preserve the spatial correlation within the 2-D image and fully exploit the useful temporal information. We perform incremental update of the object template using the N-mode SVD to model the appearance variations, which reduces the influence of template drifting and object occlusions. The proposed scheme efficiently learns a low-dimensional tensor representation through adaptively updating the eigenbasis of the tensor. Tensor based Bayesian inference in the particle filter framework is then utilized to realize tracking. We present the validation of the proposed tracking system by conducting the real-time facial expression recognition with video data and a live camera. Experiment evaluation on challenging benchmark image sequences undergoing appearance variations demonstrates the significance and effectiveness of the proposed algorithm.

Keywords: Object tracking; appearance model; incremental N-mode SVD; facial expression recognition.

1 **1. Introduction**

2 Visual tracking in image sequences is amongst the dominant bottom-up units in computer
3 vision applications such as surveillance, robotics, intelligent transportation, and human
4 computer interaction (HCI). A critical requirement of these applications is to track the desired
5 target region of interest for a long period in unconstrained environments. For instance, in face
6 based-HCI (facial expression or identity recognition), the need for accurate face tracking
7 along with head orientations has been widely acknowledged (Jo, Lee, Park, Kim & Jaihie,
8 2014). Despite much progress in visual tracking and endeavours to improve face based-HCI,
9 modelling the appearance variability of target remains imperative due to the intrinsic (e.g.
10 pose variation deformations in shape, scale, and out-of-plane rotations) and extrinsic (e.g.
11 occlusions, illumination changes, and different camera viewpoint) variations.

12 Many researchers have attempted to address these issues in visual tracking and proposed
13 various complex models to deal with target appearance variations (Wu, Yi, Lim, & Yang,
14 2015). These studies revealed the performances of existing methods, each of which has its
15 own advantages and drawbacks (Smeulders, Arnold, Dung, Rita, Simone, Afshin & Mubarak
16 2014). In visual tracking for face based-HCI such as facial expression recognition (FER),
17 template drift (Matthews, Ishikawa & Baker, 2004) is one of the common issues because of the
18 accumulation of small errors in the template updating process. For effective appearance
19 representation of a target face, frequent template update is usually necessary to cope with
20 varying pose and head orientations. An inadequate updating strategy will ruin the purpose of
21 appearance representation. In order to obtain a good trade-off between the processing time
22 and accuracy of tracker, the template-updating process must be developed carefully.

23 Secondly, the template-updating strategies based on the image-as-vector form (Ross,
24 David, Jongwoo, Lin, & Yang, 2008; Ning, Yang, Jiang, Zhang & Yang, 2016) ignore the

25 fact that image is intrinsically a matrix, or a 2nd order tensor. A significant amount of spatial
26 correlation within the original structure of a 2-D image remained unexploited in the image-as-
27 vector form, which makes the appearance model less discriminative in tracking against
28 occlusions. Alternately, the multiway or tensor based image representation can provide a
29 better appearance structure for tracking by preserving the actual 2-D structure of an image to
30 facilitate visual tracking.

31 This paper focuses on building a tracking system based on the tensor framework and
32 presents a real-time application for facial expression recognition. An important aspect of
33 object motion in videos is their local similarity among several regions of the same frame.
34 More importantly, an object in video also possesses the strong temporal correlation among
35 succeeding frames. Based on these spatial and temporal correlation priors of a video, **we**
36 **construct** the *spatio-temporal* tensor appearance model for object tracking. The proposed
37 method effectively combines the dynamic model with a robust online tensor based eigen-basis
38 updating strategy to better cope with scaling and geometric normalization issues of human
39 faces, which is a key step in face-based HCI systems.

40 Performing the proposed learning procedure using the tensor representation **will not only**
41 **preserve the 2D** structure of an image but also significantly circumvent the large
42 dimensionality problem. For example, **it can convert a** 3rd order tensor of size $30 \times 30 \times 30$ **to a**
43 **smaller dimension** of $10 \times 10 \times 10$. The image-as-vector form would require a 27000×1000
44 basis matrix, but the tensor formulation requires only three sets of 30×10 basis matrices.
45 Intuitively, the tensor based learning along with effective appearance updating yields fast and
46 more reliable target localization, which is useful for building a robust real-time FER system.

47 **1.1 Related works and context**

48 A variety of tracking approaches have been proposed **to achieve improved robustness,**

49 accuracy and computational efficiency. However, the performance of most trackers is
50 constrained with certain conditions which makes them inapt for real applications (Wu, Yi,
51 Lim, & Yang, 2015). Discriminative and generative appearance models are widely accepted
52 tracking approaches that effectively model the target appearance based on spatial information.
53 Discriminative trackers treat a tracking task as a classification problem and discriminates the
54 target from surroundings. Following the discriminative framework, an online supervised
55 boosting method was proposed (Grabner & Bischof, 2006), whereas the semi-supervised
56 tracker was introduced in (Grabner, Leistner & Bischof, 2008), in which only the initial frame
57 label was provided. Later (Babenko, Yang, & Belongie, 2009) introduced the Multiple
58 Instance Learning (MIL) tracker to deal with unreliable positive and negative labels in an
59 online manner and uncovered the drift issues in tracking. Recently the work presented in
60 (Yang, Jiang, Zhang & Yang, 2016) proposed the online structured support vector machine
61 based discriminative tracking framework with fast learning and addressed the drift issues. The
62 author in (Bae, Kang, Liu, & Chung, 2016) presented real-time object tracking framework
63 based on discrete swarm optimization. However, the proposed strategies cannot deliver the
64 orientation information of the target, or the degree of rotation of the tracking window.

65 On the other hand, several methods made use of the appearance modelling of an
66 object based on the generative framework (Black & Jepson, 1998; Ross, David, Jongwoo, Lin,
67 & Yang, 2008; Hu, Li, Zhang, Shi, Maybank & Zhang, 2011). Among them, sub-space learning-
68 based models have gained much attention in visual tracking against model drifting due to the
69 constant subspace assumption instead of the constant brightness assumption. Moreover, the
70 task of subspace learning is memory efficient, thus yielding comparatively faster processing.
71 For instance, authors in (Black & Jepson, 1998) proposed view-based appearance models for
72 tracking with sub-space learning. A view-based eigenbasis model of the target is trained off-
73 line and tracking is performed on matching sequential views of a target. However, the lack of

74 training samples along with maximum possible viewing conditions is still a challenging task.
75 The work presented by (Lim, Ross, Lin, & Yang, 2004) accomplished tracking by incremental
76 subspace learning, in which target subspace is updated during the tracking process to deal
77 with appearance changes. Their work developed an updating strategy by extending the SKL
78 (sequential Karhunen–Loeve) algorithm (Levy & Lindenbaum, 2000) but focused only on the
79 similarity between candidate and target subspace.

80 Later, (Ross, David, Jongwoo, Lin, & Yang, 2008) introduced an adaptive image-as-
81 vector subspace learning model Incremental Visual Tracker (IVT), which gained much
82 popularity. The IVT introduced the eigenbasis and mean updating strategy in tracking for
83 updating the appearance variations sequentially. The template-updating strategy followed in
84 IVT flattens the target regions to retain the vectorised shape, which yields the extrinsic
85 information about the target subregions. Several other improved versions (Hu, Li, Zhang, Shi,
86 Maybank & Zhang, 2011) of this model were further proposed. However, their performances
87 degraded in unrestricted conditions. Extended version of IVT was proposed by (Wang, Lu, &
88 Yang, 2013) under the Gaussian-Laplacian noise assumption to enhance the robustness in the
89 presence of outliers. The representation of a target by the flattened intensity vector can result
90 in a large dimensionality.

91 In visual tracking, the multilinear extension of object tracking based on online learning
92 is also introduced to capture spatio-temporal appearance and has gained much success. For
93 instance, an online tensor decomposition based tracking framework was reported by (Hu, Li,
94 Zhang, Shi, Maybank & Zhang, 2011), in which target image-as-matrix is proposed for a better
95 representation of spatial layout. For incremental updating, the R-SVD (Van Loan, 1996)
96 approach is utilized to update the subspace along with sample mean against each tensor mode.
97 However, the process of updating only considered the top eigenvalues and eigenvectors and
98 small weights were discarded in order to meet the real-time requirements, which may cause

99 error accumulation and model drift. Recently, the author in (Ma, Huang, Shen & Shao, 2016)
100 proposed the incremental tensor learning based pooling strategy and considered the target and
101 template as sparse coding tensors. Although good results are achieved on tracking image
102 sequences by using tensor pooling, however, a major concern of TPT is its extensive
103 computing procedure consisting of: (1) the 4th order online tensor learning along with the 3rd
104 order direct tensor decomposition on every upcoming frame. (2) K-means based dictionary
105 learning in tensor pooling. These factors result in slower tracking frame rate and therefore
106 make it feasible mainly for offline tracking.

107 An incremental N-mode SVD was proposed in (Lee & Choi, 2014) and tested on 3D
108 face reconstruction to update the result of N-mode SVD with the arrival of new training data.
109 The incremental N-mode SVD presented full factorization and more accurate calculation of
110 the eigenstructure of the training tensor. In this work, we focus on N-mode SVD for
111 incremental learning for online visual tracking. Unlike TPT, which uses the standard R-SVD
112 to calculate the entire N-mode immediately, we conduct separate spatial and temporal
113 factorization of the appearance tensor and adopt the incremental N-mode SVD for calculating
114 the eigenstructure of the unfolded matrices. Our method captures variants of every mode
115 independently for calculating the residue error prediction of Bayesian posterior probability.
116 Compared with other incremental tensor subspace learning and vector-based methods, N-
117 mode SVD delivers more accurate approximation of the unfolded tensors in each mode and
118 updates the target appearance variations more effectively. Thus, it makes the tracking more
119 robust against variations in pose, geometry and illumination.

120 For the task of expression recognition, the face detectors are usually employed to extract
121 the facial region first. A well-known detection algorithm by (Viola P & Jones MJ, 2004) is
122 extensively used over the years, which works on learning the classifiers based on Haar
123 features. Despite the high detection rate of this technique, the performance degrades

124 noticeably for occluded and profile faces. A real-time facial expression approach was presented
125 (Geetha, Ramalingam, Palanivel, & Palaniappan, 2009), in which head contours were extracted to
126 locate the face region in images. Color space information is further utilized to extract the location of
127 face parts and then expression recognition is performed. This method heavily depends on
128 morphological image operations such as thresholding of image pixel values and is therefore not
129 suitable for low intensity and occluded images. (Wan, Shaohua & Aggarwal, 2014) proposed a robust
130 metric learning approach for spontaneous facial expression recognition, and (Owusu, Zhan, & Mao,
131 2014) developed a facial expression analysis system based on neural-AdaBoost. Recently the authors
132 in (Ali, Hasimah, et al, 2015) used empirical mode decomposition to conduct facial expression
133 recognition. In these reported studies, face detection was performed individually on every frame. To
134 deal with the face orientation changes, generally several pre-processing techniques are combined
135 carefully with this face detection framework to register the faces based on the locations of
136 eyes and nose before the expression recognition. However, this stage demands accurate
137 detection of facial parts, which is generally not possible in real-time applications.

138 **Contributions:** Based on the above-mentioned discussion, the correlation between the
139 motion of object and its context is still hard to capture. We propose to address the visual
140 tracking problem with an effective online spatio-temporal tensor learning framework, which
141 not only takes into account the spatio-temporal information but also effectively combines the
142 updating procedure with appearance modelling to achieve real-time tracking. The proposed
143 tracking algorithm produces a low dimensional tensor representation of the target online by
144 following an incremental update procedure of the mean and eigen-basis using the N-mode
145 SVD of unfolded matrices. When estimating the target, the likelihood of a candidate is
146 evaluated on the learned tensor subspace repeatedly based on the reconstruction error to avoid
147 missing the target position. The spatio-temporal appearance feature and fast incremental
148 update provide an improved tracking performance along with better computational efficiency

149 when compared to vector based subspace methods. Then we specifically designed a real-time
150 FER system based on the proposed tracking strategy. The tracking window obtained from the
151 proposed tracker is further processed to align the face geometry and then the task of
152 expression recognition is performed. Experiments revealed that the integrated system
153 performs effectively in both recorded videos as well as on live camera enabled videos.

154 The remaining of this work is organized as follows. In Section 2, we describe the material
155 and methods utilized in our proposed framework with a complete outline of our tracking
156 approach. Section 3 provides experiment results of the proposed tracking algorithm. Section 4
157 discusses the applications of FER using our method. Finally, we present the conclusion in
158 Section 5.

159 **2. Proposed Framework for Tracking**

160 **2.1 Outline of our tracking method**

161 The proposed tracking framework is built on three main stages as shown in Figure 1: (a)
162 spatio-temporal tensor based target appearance model, (b) Bayesian inference coupled with
163 particle filter, and (c) N-mode SVD for incremental updating tensor. In Figure 1(a), we
164 propose a tensor based approach to construct the target template appearance as a reservoir of
165 3rd order tensors. A tensor is initially decomposed using Higher-Order Singular Value
166 Decomposition (HOSVD) (Kolda, Tamara G., & Brett, 2009), where the appearance model
167 only takes into consideration of the initial region of interest to be tracked in subsequent
168 frames. In Figure 1(b), we consider the tracking problem in a generative framework as an
169 online tensor learning task. An accurate subspace representation is learned online, and the
170 updating procedure is carried out in the temporal direction. The processes (b) and (c) are
171 combined such that the subspace of the target is computed by incremental N-mode SVD over
172 the target's intensity-value template and is stored in a leaking memory to gradually forget old

173 observations. The procedure is summarized in Table 1. Sampling of the candidate window,
 174 which is assumed to follow a Gaussian distribution around the preceding position, is carried
 175 out by using particle filtering. When predicting the target, the confidence of each sample in
 176 terms of distance from candidate to learned tensor subspace is computed. The sample with
 177 lowest error is then selected. Furthermore, the error reconstruction stage allows us to repeat
 178 the sampling when the confidence level is not sufficient enough to identify the candidate as
 179 the target. The whole process is repeated, where the new frame is added to the reservoir and
 180 the last frame is removed to provide sufficient spatiotemporal information and to avoid
 181 unnecessary storage.

182 2.2 Tensor decomposition

183 A higher order tensor can be viewed as a generalization of a vector (first-order tensor) and
 184 a matrix (second-order tensor). An N -th order tensor can be denoted as $\mathcal{A} \in \mathcal{R}^{I_1 \times I_2 \times \dots \times I_N}$,
 185 each element of which can be represented as $a_{i_1 \dots i_n \dots i_N}$ for $1 \leq i_n \leq n_N$. The n -mode (N -th
 186 dimension) matrix unfolding of a tensor \mathcal{A} , denoted as $\mathcal{A}_{(n)} \in \mathcal{R}^{n_n \times (\prod_{i \neq n} n_i)}$, is obtained by
 187 fixing the index term i_n while unfolding all other modes and combining all other indices into
 188 one index. For better visualization, let us consider the process of unfolding the tensor \mathcal{D} into
 189 its 1st, 2nd and 3rd modes. The mode- n folding process of a tensor is the reverse process of
 190 mode- n unfolding, which restores the actual tensor. The entries of the mode- n product of \mathcal{A}
 191 and a matrix $M \in \mathcal{R}^{n_n \times m_n^M}$ are:

$$192 (\mathcal{A} \times_n M)_{i_1 \dots i_{n-1} i_{m_n} i_{n+1} \dots i_N} = \sum_{i_n} a_{i_1 \dots i_n \dots i_N} M_{n_n m_n} \quad (1)$$

193 In a tensor and matrix multiplication, the resulting tensor C can also be computed by
 194 matrix multiplication $C_{(n)} = M \mathcal{A}_{(n)}$ followed by mode- n folding. Note that for tensors and
 195 matrices of the appropriate sizes, $\mathcal{A} \times_m M \times_n N = \mathcal{A} \times_n N \times_m M$

196 and $(\mathcal{A} \times_n M) \times_n N = \mathcal{A} \times_n (MN)$.

197 HOSVD is a multilinear extension of the conventional matrix singular value
198 decomposition (SVD). For an N -th order tensor, HOSVD produces N orthonormal
199 matrices, $U^{(1)}, U^{(2)}, \dots, U^{(N)}$, spanning N spaces. An N -th order tensor $\mathcal{A} \in \mathcal{R}^{I_1 \times I_2 \times \dots \times I_N}$ can be
200 decomposed as follows:

$$201 \min_{\mathcal{C}, U^{(1)}, U^{(2)}, \dots, U^{(N)}} \|\mathcal{A} - \mathcal{C} \times_1 U^{(1)} \times_2 U^{(2)} \dots \times_N U^{(N)}\|$$

202 (2)

203 where $\mathcal{C} \in \mathcal{R}^{R_1 \times R_2 \times \dots \times R_N}$ is called the core tensor, and $U^{(n)} \in \mathcal{R}^{I_n \times R_n}$ contain singular vectors.
204 The solution of the above equation can be given by Tucker Decomposition (L.R.Tucker,
205 1966), which is the preliminary step in obtaining the start-up tracking procedure and
206 incremental tensor learning in our tracking algorithm.

207 **2.3 Online tensor learning for tracking**

208 Online tensor learning model for tracking is built from the streaming data. The first K
209 target frames warped from sliding data are stored in terms of image gray levels in the initial
210 window as $\mathcal{A} \in \mathcal{R}^{I_1 \times I_2 \times I_3}$, where I_1 and I_2 are the width and height of target and I_3 is the
211 number of frames stacked in the tensor. Subsequently, the tensor \mathcal{A} is decomposed using
212 HOSVD into three orthogonal spaces by three orthonormal matrices U_1 , U_2 , and U_3 . When
213 a new frame comes from the video stream, the last frame from the sliding block is removed.
214 For each new frame, we may only need a portion of the mode matrices to further compute the
215 SVD, rather than re-computing the whole tensor. Incremental SVD, in this case, serves the
216 purpose to update the previous mode matrices with the arrival of new data.

217 The classic R-SVD algorithm operates on newly accessorial columns and rows in the
218 matrix, but is based on the zero mean assumption. In multi-linear generalization (Lee & Choi,

219 2014) introduced the N-mode SVD to compute the eigen-basis of a tensor with the mean
 220 update (Hall, Marshall, Martin, 2003) and therefore can keep tracking the subspace variations in
 221 each mode. For incremental updating of basis matrices, the incremental N-Mode SVD (Lee &
 222 Choi, 2014) is utilized. An extensive procedure for the incremental N-Mode SVD is followed
 223 to compute the eigen-basis with mean updating simultaneously. The process is summarized in
 224 Table 1. This algorithm can approximate the N-mode SVD efficiently with less memory and
 225 operate on a smaller portion of data each time from a relatively larger dataset. In this section,
 226 we provide an overview of N-Mode incremental update adapted for tensor based appearance
 227 model in context of visual tracking. The complete derivation of mode matrices can be found
 228 in (Lee & Choi, 2014).

229 After preparing the reservoir tensor, let $\mathcal{A} \in \mathcal{R}^{I_1 \times I_2 \times \dots \times I_n^a}$ be the current tensor, and
 230 $\mathcal{B} \in \mathcal{R}^{I_1 \times I_2 \times \dots \times I_n^b}$ be the data tensor with new video frame added, then the incremental
 231 procedure includes updating the mean and the total number of samples. \mathcal{A} and \mathcal{B} can be
 232 concatenated to form: $\mathcal{D} = [\mathcal{A} \ \mathcal{B}]_N$, where $\mathcal{D} \in \mathcal{R}^{I_1 \times I_2 \times \dots \times (I_n^a + I_n^b)}$. We only need to
 233 compute the k -mode projection matrix \mathcal{U}_k of \mathcal{D} and the unfolding matrix \mathcal{D}_k for $k =$
 234 $1, 2 \dots N - 1$. The process is illustrated in Figure 2, where three identical tensors are shown
 235 along with their respective unfolded matrices in three modes. The white region corresponds to
 236 the original tensor whereas the shaded portion represents the newly added tensor.

237 In order to update the appearance tensor based on previous projection matrices, when a
 238 new sample arrives, the incremental update procedure mainly consists of two parts: (1) spatial
 239 update, and (2) temporal update. In spatial update, we consider only mode-1 and mode-2
 240 unfolding of the tensor for updating. In temporal update, we consider the third mode for
 241 updating and the process is slightly different. The steps below formulate the factorization of
 242 both parts.

243 **(1) Spatial update (when $k \neq N$)**

244 In the incremental procedure for unfolded tensor in mode-1 and mode-2, we compute
 245 the new projection mode matrix $\mathcal{U}_k^* = [\mathcal{U}_k \quad \mathcal{Q}_k] \mathcal{U}'_k$, where $\mathcal{Q}_k = \text{orth}(B_{(k)} -$
 246 $\mathcal{U}_k \mathcal{U}_k^T B_{(k)})$ and $\text{orth}()$ is obtained by the standard QR decomposition (Hall, Marshall,
 247 Martin, 2003). The concatenation of previous tensor and newly arrived frame can be formed
 248 as: $[\mathcal{A} \quad \mathcal{B}]_k$, replacing the tensor with projection matrices and core in the first two modes
 249 as:

$$250 \quad \mathcal{D} = [\mathcal{U}_k \mathcal{C}_{(k)} (\otimes_{j \neq k} \mathcal{U}_j)^T \quad B_{(k)}] \quad (3)$$

$$251 \quad \mathcal{D} = [\mathcal{U}_k \quad \mathcal{Q}_k] \begin{bmatrix} \mathcal{C}_{(k)} & \mathcal{U}_k^T B_{(k)} \\ 0 & \mathcal{Q}_k^T B_{(k)} \end{bmatrix} [\mathcal{U}_k \mathcal{C}_{(k)} (\otimes_{j \neq k} \mathcal{U}_j)^T \quad B_{(k)}] \quad (4)$$

$$252 \quad \mathcal{U}'_k = \text{orth} \left(\begin{bmatrix} ff * \mathcal{C}_{(k)} & \mathcal{U}_k^T B_{(k)} \\ 0 & \mathcal{Q}_k^T B_{(k)} \end{bmatrix} \right) \quad (5)$$

$$253 \quad \mathcal{U}_k^* = \left[\mathcal{U}_k \quad \text{orth} \left([ff \mathcal{D}_{(k)} \quad \sqrt{\frac{ff * n_a * n_b}{ff * n_a + n_b}} (\bar{\mathcal{A}} - \bar{\mathcal{B}})] \right) \right] * \mathcal{U}'_k \quad (6)$$

254 where ff is the forgetting factor (Ross, David, Jongwoo, Lin, & Yang, 2008) for
 255 concentrating more effect on newly arrived sample.

256 (2) Temporal Update (where $k = N$)

257 For incremental update of 3rd mode, the updating structure for \mathcal{U}_N^* is different from the
 258 spatial structure in terms of concatenation and can be given as:

$$259 \quad \mathcal{D} = \begin{bmatrix} \mathcal{A} \\ \mathcal{B} \end{bmatrix}_{k=N} \quad (7)$$

$$260 \quad \mathcal{D} = \begin{bmatrix} \mathcal{U}_N & 0 \\ 0 & E \end{bmatrix} \begin{bmatrix} \mathcal{C}_N (\otimes_{j \neq N} \mathcal{U}_j)^T \\ \mathcal{B}_{(N)} \end{bmatrix} \quad (8)$$

$$261 \quad \mathcal{B}' = \mathcal{B}_{(N)} (\otimes_{j \neq N} \mathcal{U}_j) \mathcal{C}_N^T \quad (9)$$

$$262 \quad \mathcal{U}_N^* = \begin{bmatrix} \mathcal{U}_N & 0 \\ 0 & E \end{bmatrix} \text{orth} \left(\begin{bmatrix} C_N C_N^T & \mathcal{B}'^T \\ \mathcal{B}' & \mathcal{B}_{(N)} \mathcal{B}_{(N)}^T \end{bmatrix} \right) \begin{bmatrix} C_N (\otimes_{j \neq N} \mathcal{U}_j)^T \\ \mathcal{B}_{(N)} \end{bmatrix} \quad (10)$$

263 Instead of using the standard R-SVD (Hall, Marshall, Martin, 2003) that calculates the
 264 entire N-mode \mathcal{U}_{new} immediately, we propose to utilize spatial as well as temporal
 265 factorization of appearance tensor model that is updated dynamically using N-mode SVD.
 266 **Tracking is** achieved by the actual factorization of each unfolded matrix. The method
 267 provides an accurate approximation by keeping the dominant singular subspaces of current
 268 updating model. It incrementally builds Gaussian mixture models on each mode to describe
 269 the data falling into several classes in spatial domain and temporal domain incorporating the
 270 tensor multi-linear representations.

271 In the tracking framework, the newly arrived sample is categorized by evaluating its
 272 likelihood with the estimated subspace. The likelihood can be determined by the sum of the
 273 reconstruction residual error norms of the predictive Gaussian distribution. By the means of
 274 the notation of orthogonality in the tensor decomposition, we associate the dynamic Bayesian
 275 inference to approximate the distribution over the location of target. Let \bar{X} be the mean of
 276 observations, where X represents the center of the distribution for each class. Assume
 277 $g_k = U_k^T (X - \bar{X})$, which represents the shift of observations to the mean \bar{X} . The residual
 278 error vector H_k , orthogonal to every vector in U_k , is given by:

$$279 \quad H_k = (X - \bar{X}) - U_k U_k^T (X - \bar{X}) \quad (11)$$

280 The sum of residual error norms of a predictive state can be represented as:

$$281 \quad \text{error} = \sum_{k=1}^2 \left\| (X_{(k)} - \bar{X}_{(k)}) - (X_{(k)} - \bar{X}_{(k)}) \times_k (U_k U_k^T) \right\|^2 \\ 282 \quad \quad \quad + \left\| (X_{(N)} - \bar{X}_{(N)}) - (X_{(N)} - \bar{X}_{(N)}) \times_N (U_N U_N^T) \right\|^2 \quad (12)$$

283 2.4 Tensor likelihood Bayesian inference

284 We use online tensor learning to model the tracking process under the assumption that the
 285 object state in the tracking framework exhibits the Markov chain state transition process,
 286 where the present state can be effectively estimated from its past states. In this model, the
 287 motion of target among consecutive frames is usually considered as an affine motion. Assume
 288 that a given state of target is $X_t = \{x_t, y_t, \theta_t, s_t, \beta_t, \phi_t\}$ at time t , where the parameters are the
 289 x and y translations, rotation angle, scale, aspect ratio, and skewness, respectively. Now we
 290 consider the tracking formulation in Bayesian filtering framework, in which the hidden state
 291 of X_t of the target object at each time t is estimated with one of the k image observations
 292 $Z = \{Z_1, Z_2, \dots, Z_k\}, \{Z_t | t = 1, 2 \dots k\}$. Under this framework, the filtering Bayesian estimate
 293 of posterior $P(X_t|Z_t)$ can be given as:

$$294 P(X_t|Z_t) \propto P(Z_t|X_t) \int P(X_t | X_{t-1})P(X_{t-1} | Z_t)d_{X_{t-1}}$$

295 (13)

296 where $P(Z_t|X_t)$ refers to the observation likelihood at time t , and $P(X_t | X_{t-1})$ corresponds to
 297 the motion model. In order to approximate the distribution over the target position and to
 298 draw the set of samples $X = \{X_t^{(i)}\}_{i=1}^{N_s}$, particle filter (Isard & Blake, 1996) is utilized. The
 299 optimal object state of the target \hat{X}_t in the present frame can be inferred by the maximum a
 300 posteriori (MAP) criterion:

$$301 \hat{X}_t = \operatorname{argmax}_{X_t \in X} P(X_t|Z_t)$$

302 (14)

303 In our implementation, the residual error determines the measure of similarity among
 304 candidate and the learned subspace. Thus, $P(Z_{t_i}|X_t)$ in our case is given by:

$$305 P(Z_{t_i}|X_t) \propto \exp(-\text{error}) \tag{15}$$

306 **3. Experiments**

307 For performance evaluation, the proposed online spatio-temporal tensor learning based
308 tracker is validated on seven challenging videos. The chosen sequences comprise of various
309 variations, including partial occlusion, pose and scale variations, illumination changes,
310 background cluttering, rotations and impulse motions. For comparison, we conducted
311 experiments on **several related tracking methods**. (1) Fragments-based tracking (FRAG
312 Tracker) (Adam, Rivlin, & Shimshoni, 2006) (2) Vector subspace learning-based tracking
313 algorithm (IVT tracker) (Ross, David, Jongwoo, Lin, & Yang, 2008). (3) Adaptive structural
314 local sparse appearance model (ASLA Tracker)(Jia, Lu & Yang, 2012), (4) **Discriminant tracker**
315 **based on circulant structure with kernels (CSK Tracker)** Henriques, J. F., Caseiro, R.,
316 **Martins, P., & Batista, J. 2012)**, and (4) **Sparsity based collaborative model for tracking**
317 **(SCM Tracker)** (Zhong, W., Lu, H., & Yang, M. H. 2012). For a fair comparison, the
318 proposed tracker is evaluated against these methods using the results provided by authors in
319 the benchmark (Wu, Lim & Yang, 2015). In our experiments, the target region obtained from
320 the video frames is normalized to the size of 32×32 pixels, and an initial tensor of length 15 in
321 the 3rd mode is built for the representation of object appearance. The forgetting factor in the
322 incremental N-mode SVD is set to 0.9, where the number of particles in the particle filter is
323 set to 500. The assigned affine parameter values are [9,5,0.05,0.005,0,0].

324 **3.1 Evaluation**

325 As discussed above, the proposed methodology based on the image as a 2nd order tensor
326 and appearance model as a 3rd order tensor can preserve more compact and useful
327 information as compared to an image represented as a vector. In addition, the incremental N-
328 mode SVD delivers a more robust tensor updating procedure. To evaluate the effectiveness of
329 our proposed schemes, we present both quantitative and qualitative comparison with related

330 methods against dominant challenges, such as occlusion, illumination changes, target scaling
331 (deformation) and rotations.

332 **3.1.1 Quantitative analysis**

333 We used the conventional metric position error, precision plots of one pass evaluation
334 (OPE) and success plot of OPE (Wu Y, Lim J, Yang, 2015) to evaluate the tracking
335 performance. The tracking windows obtained from IVT, FRAG, CSK, ASLA, SCM and our
336 proposed algorithm are compared with the available ground truth to generate the mean square
337 error. The average location error of each method on each video is listed in Table 2. Figure 3
338 shows the screen shots of the proposed tracker under several challenging conditions. The
339 relative position error per frame (in pixels) between the tracking result and ground truth is
340 reported in Figure 4, whereas the visual comparison on key frames is presented in Figure 5.

341 It is evident from Table 2 that the proposed tracking method is more effective than other
342 vector based methods. The advantage of our method is much notable with face based videos,
343 for which our tracker achieved better performance. For other videos, our method provides the
344 second best results with insignificant margins from the best results. The reason behind is that,
345 other videos have fewer challenging conditions and do not contain abrupt motions. However,
346 when the impulse motion and orientation changes occur, tensor based image representation
347 can provide a better performance. SCM achieved the second best result following the tensor
348 based tracker due to its frequent model updating strategy.

349 *Overall Performance:* The overall performance of the related trackers is evaluated using the
350 precision plots and success plots. The precision plot evaluates the robustness in terms of
351 percentage of video frames whose recorded location is within the provided threshold distance
352 to the ground-truth. Whereas the success plot is the measure of area under the curve (AUC) of
353 each tracker. Success plot indicates the ratio among correctly tracked frames whose overlap

354 threshold is larger than the given threshold. In terms of accuracy and overlap precision, the
355 overall performance of the tracker was also found to be better. Figure 6 shows precision plot
356 ranking with the threshold of 20 pixels. The proposed tracker achieves 6% better precision as
357 compared to SCM, whereas in success plot the proposed tracker achieves 5% better ranking
358 over CSK in terms of AUC.

359 **3.1.2 Qualitative analysis**

360 The visual analysis of the proposed tracking method is also carried out under several
361 challenging conditions.

362 *Occlusion:* Figure 3(a) shows the results of tracking key frames with occlusion. The subject
363 face in this sequence is being severely occluded by a book. The target in frame 50 covered
364 major part of her face by the book, but accurate tracking is still achieved using our method as
365 the tensorial information of the target is retained to compute the similarity, which makes it
366 less susceptible to noise.

367 In terms of locating tracked objects, our method also provides good accuracy. Figure 4
368 shows the result on the faceOCC1 sequence, in which the proposed method performed better
369 than ASLA and FRAG due to effective updating. The error rates of CSK and IVT are
370 comparable to ours. FRAG performs poorly in the occlusion scenario, due to the lack of
371 mechanism to deal with the appearance changes.

372 *Illumination variation and scale changes:* The tracking results for videos with severe
373 illumination changes is depicted in the video sequence named *Mhyang* as shown in Figure
374 3(d). In frame 540, the target moves backward from its position and resulting in scale
375 changes, whereas in frames 760 and 1200, the illumination effect is more significant, which
376 makes the tracking process more challenging. Frame 1410 is chosen to indicate the target

377 rotation. The proposed method performs effectively under these conditions.

378 The video sequence provided in Figure 3(b) shows a person walking out of the dark
379 place into an area with spot lights, where the motion of target and camera is also found. In
380 this sequence, our tracker and IVT successfully retained the tracking region throughout the
381 sequence, whereas FRAG lost the tracking process and drifted away from its position. As
382 evident from Figures 5(c) and 5(d), our tracker best places the tracking window on the target
383 and accurately captures the face rotation along with the side pose, while the ASLA and IVT
384 windows are scaled slightly larger than the target.

385 The sample results shown in Figure 3(c) are taken from video sequence *Dudek*, which
386 involves illumination, scale and pose variations. Frame 209 shows the tracking result under
387 quick face occlusion, whereas frames 500 and 997 show the expression variation. The target
388 is also under the motion with changing background and scale. The proposed tracker quickly
389 adapts the changes in scale and poses as evidenced in the results.

390 The effect of head pose changing on tracking result in the comparison to others is
391 more evident in Figure 5(b). In frame 250, all other trackers captured the face location but the
392 head rotation is more accurately located by **TTracker**. Moreover, the scale changes in
393 subsequent frames of Figure 5(b) is modelled effectively with **TTracker**, where ASLA and
394 CSK localized a larger region and FRAG tracker completely lost the target. This is due to the
395 proposed incremental template updating mechanism, which empowers the tracker to cope
396 with gradual appearance changes in our method, whereas the classical approaches were not
397 able to do this.

398 **The good performance** of the proposed tracking framework **can be** credited by the fact
399 that the tensor framework **delivers** more structural motion information of the target as
400 compared to vector based strategies. In terms of adaptability, the incremental tensor based

401 learning through the N-mode SVD could learn more specific appearance changes of the target
402 by capturing the variations with the passage of time. Meanwhile, the tensor reservoir in our
403 proposed model serves the purpose of retaining the information in each mode. In case the
404 information content related to one mode of the tensor is affected, the remaining modes are
405 still enough to restore the subspaces contents. Hence, the minor imprecisions in the template
406 location do not accumulate further, and thus the procedure assists the tracker to tolerate the
407 template drift. We presented the visual results related to face based images, so as to
408 emphasize the tracking performance under the natural head movements and orientation
409 changes. Our method is particularly useful for face based HCI, such as expression
410 recognition.

411 **4. Application to human facial expression recognition**

412 Conducting real-time face tracking to examine the facial expressions is precluded by
413 problems such as head pose orientation, scale variation and at the same time the paucity of
414 fast processing framework. **On the other hand**, face tracking provides a promising tool to
415 track the initial face position of subject for subsequent frames with less computational
416 complexity. Meanwhile, the scale and rotation information can also be effectively **determined**
417 along with face **tracking**, which can sufficiently avoid the need of face registration.

418 In this section, we present the real-time facial expression recognition system based on the
419 proposed tracker. **We have considered a seven-class recognition problem including the neutral**
420 **expression, and six prototype expressions, Happy, Sad, Fear, Disgust, Anger and Surprise.**
421 We also present the performance evaluation on publically available facial expression datasets
422 and also on self-collected videos.

423 **4.1 Facial expression recognition**

424 Generally, a facial expression recognition system involves two key phases: effective facial

425 representation and accurate classifier design.

426 **4.1.1 Feature extraction**

427 In our work, we used Gabor filters to extract the facial information from the tracked
428 face region. Gabor filters are widely accepted in many face based HCI systems, owing to their
429 robustness against photometric disturbances and invariance to image registration issues
430 (M.Amin & H. Yan, 2009; Haghigat, Zonouz, & Abdel-Mottaleb, 2015). The Gabor function
431 in the spatial domain characterizes a Gaussian-shaped envelop modulated by a complex
432 sinusoidal signal:

$$433 \quad g(x, y) = \frac{1}{2\pi\delta_x\delta_y} \exp\left\{-\frac{1}{2} \left[\left(\frac{x}{\delta_x}\right)^2 + \left(\frac{y}{\delta_y}\right)^2\right]\right\} + i(ux + vy) \quad (16)$$

434 A set of Gabor functions with multi-scales and orientations are employed for extracting more
435 compact and effective image representation. We used 3 scales and 5 orientations of Gabor
436 function in our experiment to extract the frequency contents of the image. Furthermore, 8-bit
437 down-sampling is carried out to reduce the neighbourhood pixel redundancy.

438 **4.1.2 Classification**

439 The final stage of the FER system is based on classifier design. We used support
440 vector machines (SVMs) for generalized performance. SVM is a binary discriminant
441 classifier which is based on structural risk minimization principle that produces the maximum
442 margin hyperplane between two classes. As we considered 7-class recognition problem, a
443 multiclass SVM classifier can be constructed by using the one-against-all strategy (J. Weston
444 and C. Watkins 1998). Here we briefly present the multiclass SVMs used in our experiments.
445 Given the training data of size N as; $(g_1, l_1), \dots, (g_N, l_k)$ where, $g_j \in \mathfrak{R}^R$ feature vector and
446 $l_j \in \{1, \dots, 7\}$ represent the corresponding expression labels. Multiclass SVMs defines only
447 one optimization problem but constructs seven class rules. So to follow the k th function
448 $\mathbf{w}_k^T \phi(g_j) + b_k$ to partition training vectors of class k from remaining feature vectors, we

449 minimize the following objective function:

$$450 \min_{\mathbf{w}, \mathbf{b}, \xi} \frac{1}{2} \sum_{k=1}^7 \mathbf{w}_k^T \mathbf{w}_k + \sum_{j=1}^N \sum_{k \neq l_j} \xi_j^k \quad (17)$$

451 Subject to the constraints:

$$452 \mathbf{w}_{l_j}^T \phi(g_j) + b_{l_j} \geq \mathbf{w}_k^T \phi(g_j) + b_k + 2 - \xi_j^k \quad (18)$$

$$\xi_j^k \geq 0, \quad j = 1, \dots, N, \{k = 1, \dots, 7 \mid l_j\}$$

453 where, ϕ is the mapping function, C penalizes the training errors, $\mathbf{b} = [b_1 \dots b_7]^T$ is the bias

454 vector and ξ is a slack variable, and $\xi = [\xi_1^1, \dots, \xi_i^k, \dots, \xi_N^7]^T$, whereas the decision function

455 can be given by:

$$456 h(g) = \operatorname{argmax}_{k=1, \dots, 7} (\mathbf{w}_k^T \phi(g_j) + b_k) \quad (19)$$

457 After training the multiclass SVMs, a new feature vector from test image is classified

458 using the equation above to recognize the facial expression.

459 4.2 Experiments on facial expression recognition

460 In this work, we used the proposed tracker to track the face location over several benchmark

461 facial expression video datasets and performed the online facial expression recognition. We

462 considered three widely used FER datasets, extended Cohn–Kanade dataset (CK+)

463 (Lucey, Cohn, Kanade, Saragih, Ambadar, & Matthews, 2010), MMI dataset (Pantic, Valstar,

464 Rademaker & Maat, 2005) and MUG dataset (Aifanti, Niki, Christos, & Anastasios

465 Delopoulos, 2010). In order to guarantee the generalization performance of the system, one of

466 the dataset was used to prepare the training images, and the other two datasets were used for

467 testing. As CK+ contains more subjects and videos, we used it for training. Challenging

468 subject videos that contain sufficient head rotations from other datasets were used for testing.

469 Figure 7 shows the tracking result obtained from the proposed tracker on images taken from

470 the MMI and MUG datasets with apex frames. The second row of Figure 7 demonstrates the

471 cropping result from the tracker window. It is evident that despite the presence of scale

472 variation and head rotations, the cropped images are well aligned in terms of face geometry.

473 Table 3 presents the result of SVMs in terms of recognition rate against different kernels
474 used. Apart from Gabor features, we also recorded the recognition rate for histogram based
475 features based on local binary patterns (LBP) (Zavaschi, Britto, Oliveira & Koerich, 2013).
476 LBP features can be computed more efficiently than Gabor features, but are more sensitive to
477 illumination variations and rotations. Gabor features with linear SVMs have the highest
478 recognition rate compared with other kernels and LBP.

479 Tracking results from the proposed tracker were used to compare the performance of FER
480 with the results obtained by other methods. We evaluated the results from several approaches,
481 including face detector based method (Rahulamathavan, Y., Phan, R. C. W., Chambers, J. A.,
482 & Parish, D. J. 2013), active appearance model based feature detector, (Aifanti, N., &
483 Delopoulos, A. 2014) and feature point based face model (Kumar, S., Bhuyan, M. K., &
484 Chakraborty, B. K. 2016). Figure 8 compares the results of 6 expressions on the MUG
485 database using the leave-one-out validation strategy. The reason for using only 6 expressions
486 here is that several methods above do not consider the neutral expression. It can be seen that
487 the results from the proposed method were consistently higher for each expression. For the
488 surprise expression, our recognition rate is comparable with the feature detection method,
489 whereas for the fear expression, the margin between our method and the feature detection
490 method is higher. The overall average recognition rate of our method is also better than all
491 other algorithms. This fact can be credited by the accurate face registration using our
492 proposed method, which yields comparatively better feature representation for FER.

493 **4.3 Online experiments**

494 A graphic user interface (GUI) is designed and utilized for conducting online experiments
495 with videos taken by a camera. Figure 9 shows a snapshot of the GUI. A subject is asked to

496 perform expressions in front of the camera and real-time expression recognition is performed.
497 The figure shows the input frame with the tracking window and corresponding the aligned
498 and cropped face region. The bottom graph indicates the confidence level of each expression
499 for a test frame. A higher score indicates the presence of particular expression. The tracking
500 window is obtained on every frame and is then processed further to obtain the affine
501 transformed tracked image. Classification is done later and the expression label is generated.
502 **Figure 10** shows the result of proposed tracker on **250** frames of a recorded video. The subject
503 starts from neutral expression and plays **7** expressions. The top image shows the key frames
504 of video with expression variation and the tracking window. It can be seen that the tracker
505 follows the face region accurately over all **250** frames, despite head rotations. The bottom
506 curves show the normalized score of each expression. It is interesting to note that, the neutral
507 class, which is a transition expression between two universal expressions and is played for a
508 very short time, is also being effectively recognized by the system. The main reason behind
509 these refined results is the perfect face tracking that reduces the problems of miss alignment
510 and registration of the face.

511 Apart from benchmark videos from MMI and MUG datasets, we also tested the proposed
512 tracker on self-collected videos in order to quantify the recognition performance. Table 4
513 shows the **confusion matrix for 7 expressions** associated with these test data of 13 persons in
514 total 30 videos. As can be seen from the table, the average recognition rate of 94.16 % is
515 achieved, where diagonal entries indicate the recognition rate of each expression. The surprise
516 and happy expressions can be recognized with the highest accuracy, while we noticed that sad
517 and neutral expressions are sometimes difficult to recognize, due to the fact that the sad
518 expression is highly subject dependent, which sometimes resembles the neutral state.

519 **5. Conclusion**

520 In this paper, we propose a tensor based method to construct the target template appearance

521 as a reservoir of 3rd order tensors and consider the object tracking problem in a generative
522 framework as an online tensor learning task. An effective N-mode SVD based tensor
523 eigenspace representation is learned online, and the updating procedure is carried out over the
524 time span. The proposed multi-mode model is demonstrated for object tracking to better deal
525 with large appearance variations caused by shape deformations, occlusions and drifts.
526 Experiment comparisons with existing tracking strategies revealed the effectiveness of the
527 proposed method, especially when the target motion is under rotational changes. Finally, the
528 task of facial expression recognition is investigated by integrating the proposed tracking
529 strategy with an expression recognition module. A GUI is developed and used for evaluation
530 of our method on public and self-prepared videos. Our system can effectively recognize
531 human facial expressions from videos and streaming camera with encouraging recognition
532 rate of 94.16% for 7 classes of basic expressions. We believe that the proposed method will
533 facilitate future development of face based-HCI applications and find useful applications to
534 other object tracking and recognition systems.

535 **Acknowledgements**

536 This work is supported by Hong Kong Research Grants Council (Project C1007-15G)

References

- Adam, A., Rivlin, E., & Shimshoni, I. (2006, June). Robust fragments-based tracking using the integral histogram. In *Computer vision and pattern recognition, 2006 IEEE Computer Society Conference on* (Vol. 1, pp. 798-805). IEEE.
- Aifanti, N., & Delopoulos, A. (2014). Linear subspaces for facial expression recognition. *Signal Processing: Image Communication*, 29(1), 177-188.
- Aifanti, N., Papachristou, C., & Delopoulos, A. (2010, April). The MUG facial expression database. In *Image Analysis for Multimedia Interactive Services (WIAMIS), 2010 11th International Workshop on* (pp. 1-4). IEEE.
- Ali, H., Hariharan, M., Yaacob, S., & Adom, A. H. (2015). Facial emotion recognition using empirical mode decomposition. *Expert Systems with Applications*, 42(3), 1261-1277.
- Babenko, B., Yang, M. H., & Belongie, S. (2009, June). Visual tracking with online multiple instance learning. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on* (pp. 983-990). IEEE.
- Bae, C., Kang, K., Liu, G., & Chung, Y. Y. (2016). A novel real time video tracking framework using adaptive discrete swarm optimization. *Expert Systems with Applications*, 64, 385-399.
- Black, M. J., & Jepson, A. D. (1998). Eigentracking: Robust matching and tracking of articulated objects using a view-based representation. *International Journal of Computer Vision*, 26(1), 63-84.
- Geetha, A., Ramalingam, V., Palanivel, S., & Palaniappan, B. (2009). Facial expression recognition—A real time approach. *Expert Systems with Applications*, 36(1), 303-308.
- Grabner, H., & Bischof, H. (2006, June). On-line boosting and vision. In *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on* (Vol. 1, pp. 260-267). IEEE.
- Grabner, H., Leistner, C., & Bischof, H. (2008). Semi-supervised on-line boosting for robust tracking. *Computer Vision—ECCV 2008*, 234-247.
- Haghighat, M., Zonouz, S., & Abdel-Mottaleb, M. (2015). CloudID: Trustworthy cloud-based and cross-enterprise biometric identification. *Expert Systems with Applications*, 42(21), 7905-7916.
- Hall, P., Marshall, D., & Martin, R. (2002). Adding and subtracting eigenspaces with eigenvalue decomposition and singular value decomposition. *Image and Vision Computing*, 20(13), 1009-1016.
- Henriques, J. F., Caseiro, R., Martins, P., & Batista, J. (2012, October). Exploiting the circulant structure of tracking-by-detection with kernels. In *European conference on computer vision* (pp. 702-715). Springer, Berlin, Heidelberg.
- Hu, W., Li, X., Zhang, X., Shi, X., Maybank, S., & Zhang, Z. (2011). Incremental tensor subspace learning and its applications to foreground segmentation and tracking. *International Journal of Computer Vision*, 91(3), 303-

327.

Isard, M., & Blake, A. (1996, April). Contour tracking by stochastic propagation of conditional density. In European conference on computer vision (pp. 343-356). Springer Berlin Heidelberg.

Jia, X., Lu, H., & Yang, M. H. (2012, June). Visual tracking via adaptive structural local sparse appearance model. In Computer vision and pattern recognition (CVPR), 2012 IEEE Conference on (pp. 1822-1829). IEEE.

Jo, J., Lee, S. J., Park, K. R., Kim, I. J., & Kim, J. (2014). Detecting driver drowsiness using feature-level fusion and user-specific classification. *Expert Systems with Applications*, 41(4), 1139-1152.

Kolda, T. G., & Bader, B. W. (2009). Tensor decompositions and applications. *SIAM review*, 51(3), 455-500.

Kumar, S., Bhuyan, M. K., & Chakraborty, B. K. (2016, March). An efficient face model for facial expression recognition. In Communication (NCC), 2016 Twenty Second National Conference on (pp. 1-6). IEEE.

L. R. Tucker, "Some mathematical notes on three-mode factor analysis," *Psychometrika*, vol. 31, no. 3, pp. 279–311, Sep. 1966.

Lee, M., & Choi, C. H. (2014). Incremental N-Mode SVD for Large-Scale Multilinear Generative Models. *IEEE Transactions on Image Processing*, 23(10), 4255-4269.

Levey, A., & Lindenbaum, M. (2000). Sequential Karhunen-Loeve basis extraction and its application to images. *IEEE Transactions on Image processing*, 9(8), 1371-1374.

Lim, J., Ross, D. A., Lin, R. S., & Yang, M. H. (2004, December). Incremental Learning for Visual Tracking. In *Nips* (Vol. 17, pp. 793-800).

Lucey, P., Cohn, J. F., Kanade, T., Saragih, J., Ambadar, Z., & Matthews, I. (2010, June). The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression. In *Computer Vision and Pattern Recognition Workshops (CVPRW), 2010 IEEE Computer Society Conference on* (pp. 94-101). IEEE.

M. A. Amin and H. Yan, (2009) An empirical study on the characteristics of Gabor representations for face recognition, *International Journal of Pattern Recognition and Artificial Intelligence*, 23(3): pp.401-431.

Ma, B., Huang, L., Shen, J., & Shao, L. (2016). Discriminative tracking using tensor pooling. *IEEE transactions on cybernetics*, 46(11), 2411-2422.

Matthews I, Ishikawa T, Baker S (2004). The template update problem. *IEEE Trans Pattern Anal Mach Intell* 26(6):810–815

Ning, J., Yang, J., Jiang, S., Zhang, L., & Yang, M. H. (2016). Object tracking via dual linear structured SVM and explicit feature map. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 4266-4274).

- Owusu, E., Zhan, Y., & Mao, Q. R. (2014). A neural-AdaBoost based facial expression recognition system. *Expert Systems with Applications*, 41(7), 3383-3390.
- Pantic, M., Valstar, M., Rademaker, R., & Maat, L. (2005, July). Web-based database for facial expression analysis. In *Multimedia and Expo, 2005. ICME 2005. IEEE International Conference on* (pp. 5-pp). IEEE.
- Rahulamathavan, Y., Phan, R. C. W., Chambers, J. A., & Parish, D. J. (2013). Facial expression recognition in the encrypted domain based on local fisher discriminant analysis. *IEEE Transactions on Affective Computing*, 4(1), 83-92.
- Ross, D. A., Lim, J., Lin, R. S., & Yang, M. H. (2008). Incremental learning for robust visual tracking. *International journal of computer vision*, 77(1), 125-141.
- Smeulders, A. W., Chu, D. M., Cucchiara, R., Calderara, S., Dehghan, A., & Shah, M. (2014). Visual tracking: An experimental survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(7), 1442-1468.
- Tucker, L. R. (1966). Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31(3), 279-311.
- Van Loan, C. F. (1996). *Matrix computations* (Johns Hopkins studies in mathematical sciences).
- Viola P, Jones MJ (2004). Robust real-time face detection. *Int J Comput Vis* 57(2):137–154
- Wan, S., & Aggarwal, J. K. (2014). Spontaneous facial expression recognition: A robust metric learning approach. *Pattern Recognition*, 47(5), 1859-1868.
- Wang, D., Lu, H., & Yang, M. H. (2013). Online object tracking with sparse prototypes. *IEEE transactions on image processing*, 22(1), 314-325.
- Weston, J., & Watkins, C. (1998). Multi-class support vector machines. Technical Report CSD-TR-98-04, Department of Computer Science, Royal Holloway, University of London, May.
- Wu, Y., Lim, J., & Yang, M. H. (2015). Object tracking benchmark. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9), 1834-1848.
- Zavaschi, T. H., Britto, A. S., Oliveira, L. E., & Koerich, A. L. (2013). Fusion of feature sets and classifiers for facial expression recognition. *Expert Systems with Applications*, 40(2), 646-655.
- Zhong, W., Lu, H., & Yang, M. H. (2012). Robust object tracking via sparsity-based collaborative model. In *Computer vision and pattern recognition (CVPR), 2012 IEEE Conference on* (pp. 1838-1845). IEEE.

Highlights

- Visual tracking in videos is an essential component in human computer interaction.
- An online tensor based learning strategy is proposed for visual tracking.
- The tracking method show superior tracking performance in challenging conditions.
- The proposed tracker delivers the scale and orientation information of the target.
- Real time facial expression recognition system is presented using proposed tracker.

Table 1: The incremental N-mode SVD algorithm for updating the tensor based appearance model.

Algorithm 1: Incremental N-Mode SVD

Function: Incremental N-Mode SVD($C_{(k)}, \mathcal{U}_k, A_{(k)}, B_{(k)}, M_k, ff, t$)

Input: Unfolded core tensor $C_{(k)}$, projection matrix \mathcal{U}_k , unfolded stacked tensor $A_{(k)}$, unfolded newly added tensor $B_{(k)}$, forgetting factor ff , updated mean M_k , time t .

Repeat

If $t = 0$, then:

- 1: Apply *HOSVD* using eq.(2) to get $[u_n, c_{(n)}]$
- 2: $M_{(t)} = \frac{n_a}{n_a+n_b} \overline{\mathcal{A}_{(k)}} + \frac{n_b}{n_a+n_b} \overline{\mathcal{B}_{(k)}}$ the mean of concatenated matrices $\mathcal{A}_{(k)}$ and $\mathcal{B}_{(k)}$
- 3: $t=t+1$

Else:

- 1: Apply N-mode updating strategy:
 - (a) Spatial update (when $k \neq N$); Eq.(4-7)
 - (b) Temporal update (where $k = N$); Eq.(8-11)
- 2: $M_{(t)} = \frac{n_a}{n_a+n_b} \overline{\mathcal{A}_n} + \frac{n_b}{n_a+n_b} \overline{\mathcal{B}_n}$
- 3: $t=t+1$

End if

Iteration until the end of video.

End

Table 2: Comparison of tracking results obtained using existing trackers and the proposed tensor based tracker (TTracker) in terms of position error on seven videos. The best and second best results are shown in bold and italic fonts respectively

Sequences	Tracking Methods					
	IVT	CSK	ASLA	SCM	Frag	TTracker
David	11.44	38.52	5.59	22.13	91.57	5.17
Fish	18.22	7.32	3.40	6.32	25.21	5.08
Mhyang	7.44	9.12	<i>2.03</i>	8.85	15.51	1.91
Twinings	10.75	10.23	16.73	8.04	22.47	<i>9.24</i>
Clifbar	59.46	47.54	57.51	31.67	40.83	44.87
Dudek	<i>10.03</i>	19.76	14.95	27.61	87.70	9.49
Faceoccl	<i>18.74</i>	17.45	78.16	22.06	51.88	15.98
Average	19.44	21.42	25.48	<i>18.09</i>	47.88	13.09

Table 3: Comparison between Gabor wavelet based features and LBP features in terms of average recognition rate for different types of kernels used in the SVM classifier.

Feature Extraction	Recognition Rate (%)		
	SVM (linear)	SVM (polynomial)	SVM (RBF)
Gabor	94.16	92.31	91.04
LBP	90.87	89.54	92.24

Table 4: Confusion matrix of recognition rates obtained using linear SVM for seven facial expressions.

Average recognition rate = 94.16%							
	<i>Angry</i>	<i>Disgust</i>	<i>Fear</i>	<i>Happy</i>	<i>Sad</i>	<i>Surprise</i>	<i>Neutral</i>
<i>Angry</i>	91.16	5.51	0	0	0	3.33	0
<i>Disgust</i>	7.84	92.16	0	0	0	0	0
<i>Fear</i>	0	0	89.67	0	0	6.64	3.69
<i>Happy</i>	0	0	0	97.18	0	0	2.82
<i>Sad</i>	0	0	0	0	95.83	0	5.17
<i>Surprise</i>	0	0	2.66	0	0	97.34	0
<i>Neutral</i>	0	0	0	0	4.21	0	95.79

Figures

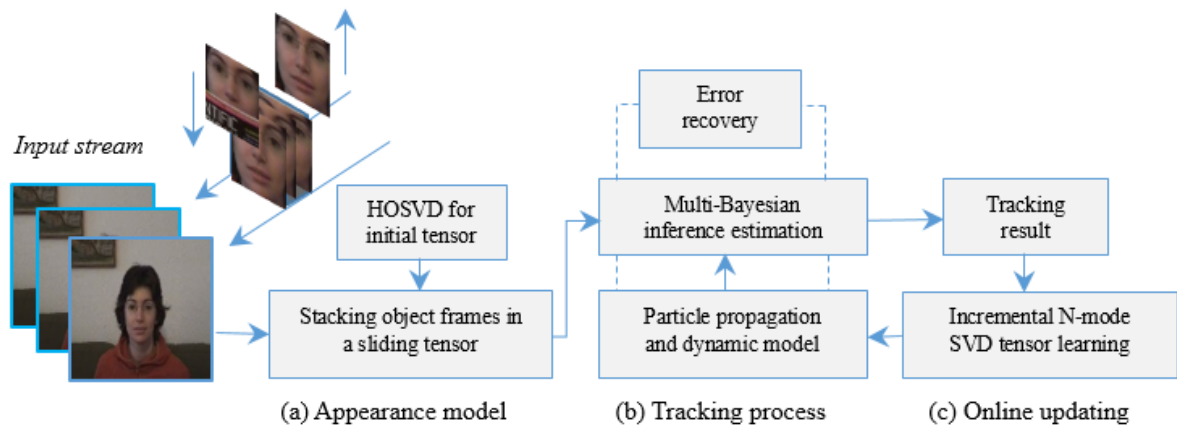


Fig 1: The architecture of the proposed online spatio-temporal tensor based learning model for visual tracking.

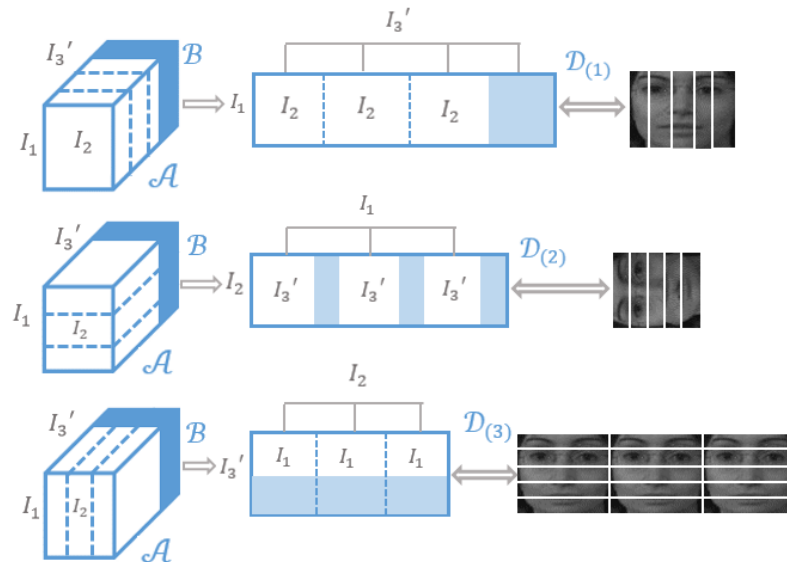


Fig 2: Tensor unfolding in respective modes and the analogous association with an image in terms of slices.

Figures



Fig 3: Screen shots of tracking results obtained by our method from key frames on face based videos in challenging sceneries.

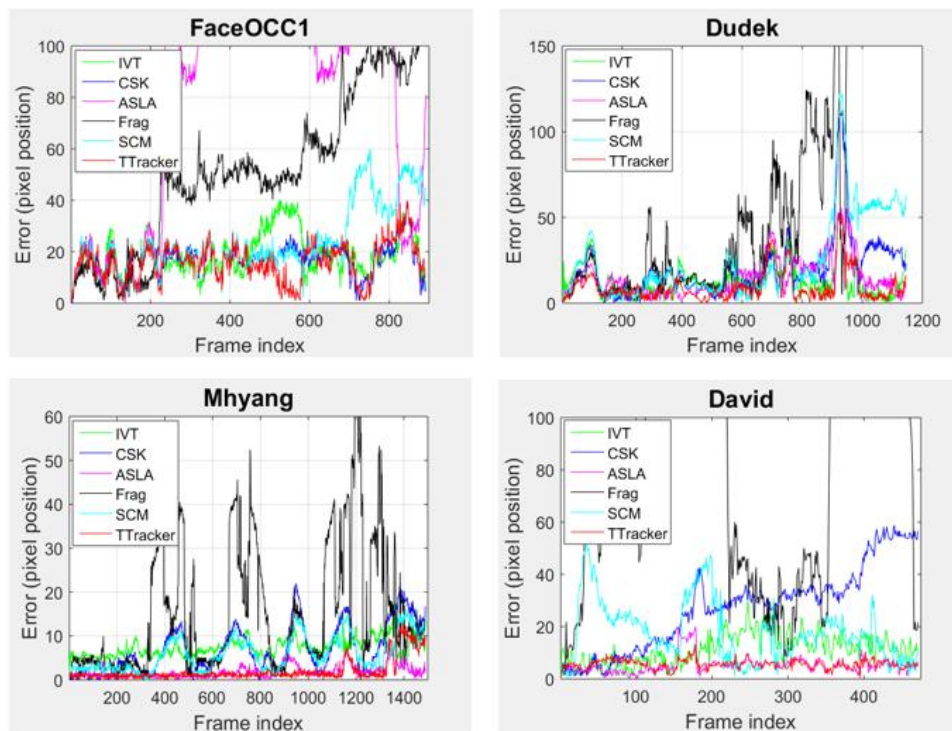


Fig 4: Position errors (pixels) with respect to ground truth and comparison of different trackers on face based videos.

Figures

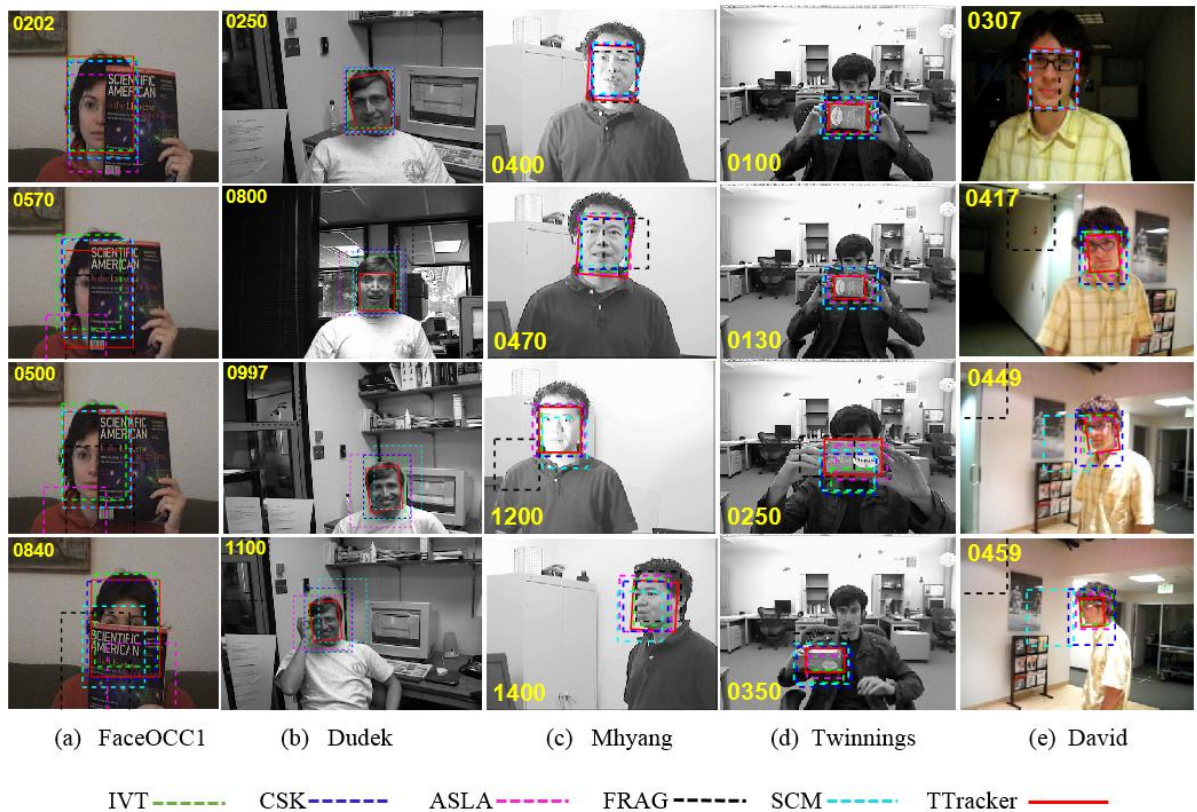


Fig 5: Screen shots of tracking results obtained by existing trackers in comparison with the proposed tensor based tracker (TTracker).

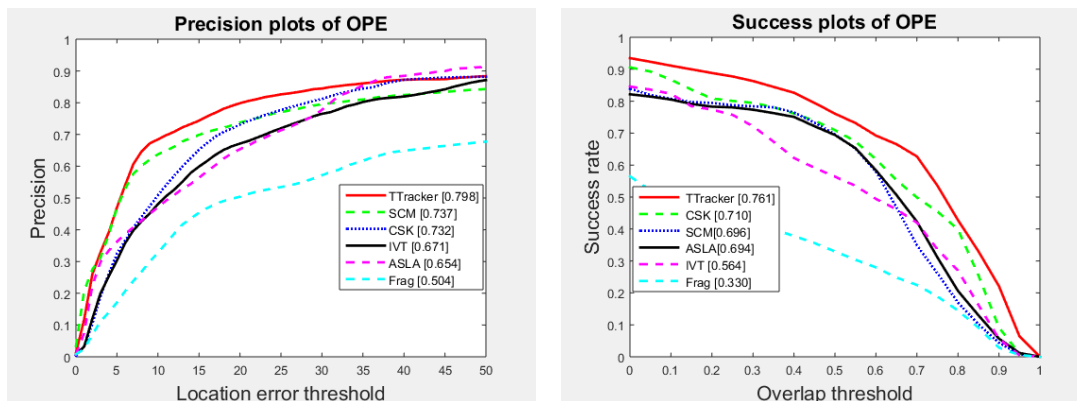


Fig 6: Comparison of different trackers in terms of the precision and success plots of the one pass evaluation (OPE) measure.

Figures

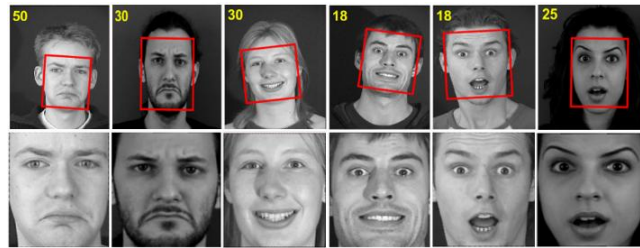


Fig 7: Results obtained by the proposed tracker on key frames taken from sample videos of MMI and MUG facial expression database under head rotations and varying expressions (upper row). Images on the lower row show cropped faces with rotation corrected face geometry.

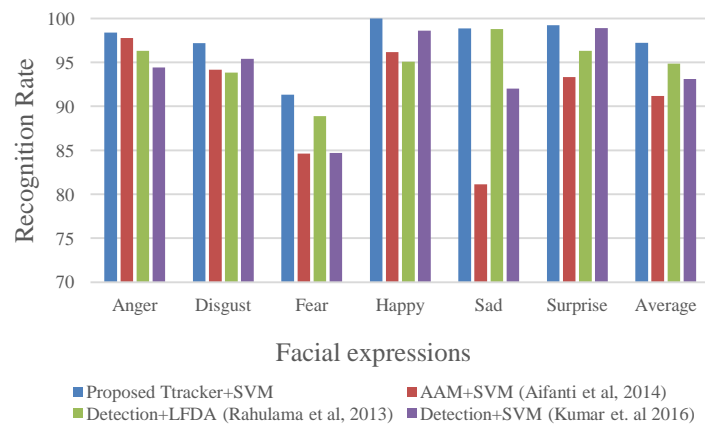


Fig 8: Performance comparison of the proposed TTracker based FER with existing facial expression recognition methods on MUG dataset.



Fig 9: Screen shot of GUI for expression recognition using the proposed tracker for offline videos as well as online ones taken by a camera.

Figures

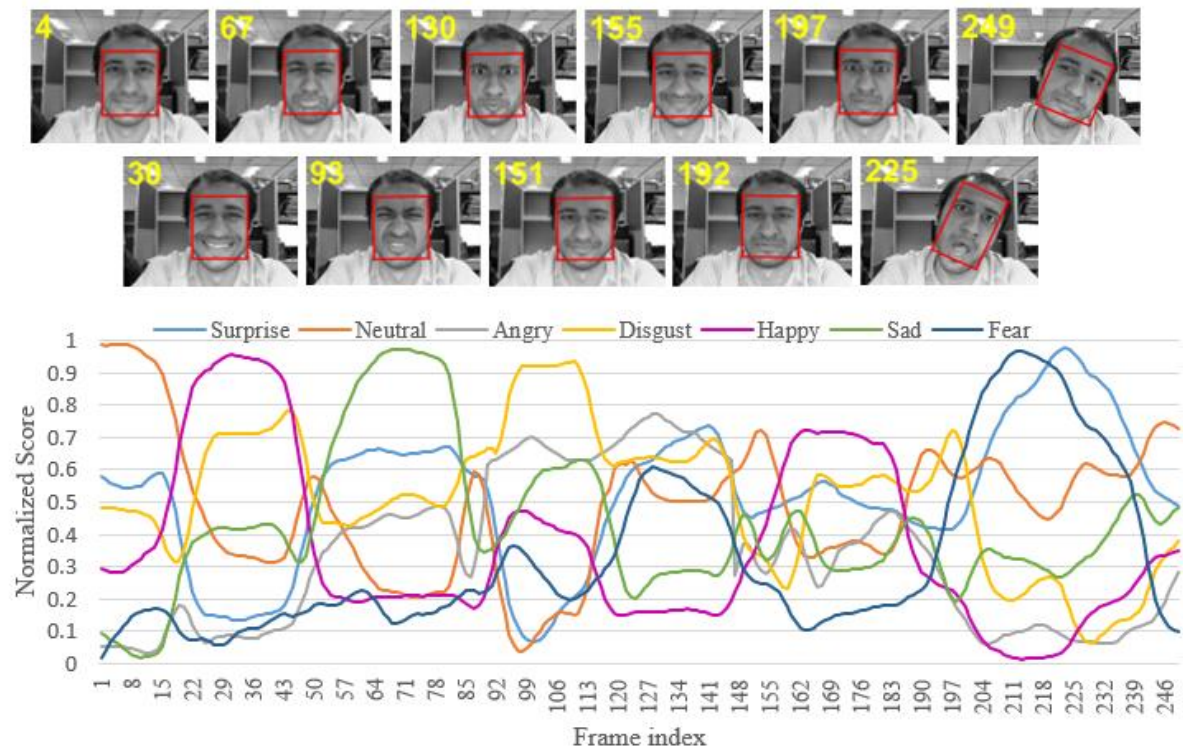


Fig 10: Test results of a subject posing facial expressions in front of a live camera. The waveform indicates the expression confidence in terms of normalized score for each frame, while the test expressions belonging to corresponding key frames along with tracking results are displayed in the images above.