

# Bags of Affine Subspaces for Robust Object Tracking

Sareh Shirazi<sup>†‡</sup>, Conrad Sanderson<sup>◦\*</sup>, Chris McCool<sup>‡</sup>, Mehrtash T. Harandi<sup>◦∇</sup>

<sup>†</sup> Australian Centre for Robotic Vision (ACRV)

<sup>‡</sup> Queensland University of Technology, Australia

<sup>∇</sup> Australian National University, Australia

<sup>\*</sup> University of Queensland, Australia

<sup>◦</sup> NICTA, Australia

**Abstract**—We propose an adaptive tracking algorithm where the object is modelled as a continuously updated bag of affine subspaces, with each subspace constructed from the object’s appearance over several consecutive frames. In contrast to linear subspaces, affine subspaces explicitly model the origin of subspaces. Furthermore, instead of using a brittle point-to-subspace distance during the search for the object in a new frame, we propose to use a subspace-to-subspace distance by representing candidate image areas also as affine subspaces. Distances between subspaces are then obtained by exploiting the non-Euclidean geometry of Grassmann manifolds. Experiments on challenging videos (containing object occlusions, deformations, as well as variations in pose and illumination) indicate that the proposed method achieves higher tracking accuracy than several recent discriminative trackers.

## I. INTRODUCTION

Object tracking is a core task in applications such as automated surveillance, traffic monitoring and human behaviour analysis [27], [42]. Tracking algorithms need to be robust to intrinsic object variations (eg., shape deformation and pose changes) and extrinsic variations (eg., camera motion, occlusion and illumination changes) [42].

In general, tracking algorithms can be categorised into two main categories: (i) generative tracking [2], [30], [35], and (ii) discriminative tracking [4], [19], [28]. Generative methods represent the object as a particular appearance model and then focus on searching for the location that has the most similar appearance to the object model. Discriminative approaches treat tracking as a binary classification task, where a discriminative classifier is trained to explicitly separate the object from non-object areas such as the background. To achieve good performance, discriminative methods in general require a larger training dataset than generative methods.

A promising approach for generative tracking is to model object appearance via subspaces [15], [25], [30], [40]. A common approach in such trackers is to apply eigen-decomposition on a set of object images, with the resulting eigenvectors defining a linear subspace. These linear subspaces are able to capture perturbations of object appearance due to variations in viewpoint, illumination, spatial transformation, and articulation. However, there are two major shortcomings. First, a linear subspace does not model the mean of the image set (ie., origin of the subspace) which can potentially hold useful discriminatory information; all linear subspaces have a common origin. Second, subspace based trackers typically search for the object location by comparing candidate image areas to the object model (linear subspace) using a brittle point-to-subspace distance [24], [34] (also known as distance-from-feature-space [36]), which

can be easily affected by drastic appearance changes such as partial occlusions. For face recognition and clustering it has been shown that improved performance can be achieved when subspace-to-subspace distances are used instead [5], [12], [31].

To address the shortcomings of traditional subspace based trackers, in this work<sup>1</sup> we propose a tracker with the following four characteristics:

- (1) Instead of linear subspaces, we propose to model object appearance using affine subspaces, thereby taking into account the origin of each subspace.
- (2) Instead of using point-to-subspace distance, we propose to represent the candidate areas as affine subspaces and use a subspace-to-subspace distance; this allows for more robust modelling of the candidate areas and in effect increases the memory of the tracker.
- (3) To accurately measure distances between subspaces, we exploit the non-Euclidean geometry of Grassmann manifolds [14], [29], [31].
- (4) To take into account drastic appearance changes that are not well modelled by individual subspaces (such as occlusions) [41], the tracked object is represented by a continuously updated bag of affine subspaces; this is partly inspired by [4], where bags of object images are used.

To the best of our knowledge, this is the first time that appearance is modelled by affine subspaces for object tracking. The proposed approach is somewhat related to adaptive subspace tracking [15], [30], [38]. In [15], [30] an object is represented as a single low-dimensional linear subspace, which is constantly updated using recent tracking results. In [38], an online subspace learning scheme employing Grassmann manifolds is used to update the object model. In the above methods, only linear subspaces and point-to-subspace distances are considered. In contrast, the proposed method uses affine subspaces and a more robust subspace-to-subspace distance. Furthermore, instead of updating a single subspace, the proposed method keeps a bag of recent affine subspaces, where old subspaces are replaced with new ones.

We continue the paper as follows. An overview of related work is given in Section II. Section III presents the proposed tracking approach in detail. Comparative evaluations against several recent tracking methods are reported in Section IV. The main findings and possible future directions are given in Section V.

---

<sup>1</sup>This paper is a thoroughly revised and extended version of our earlier preliminary work [33].

## II. RELATED WORK

In this section, we first overview the evolution of subspace-based trackers. We then briefly describe two popular generative trackers: the mean shift tracker [9] and the fragments-based tracker [2]. Finally, we briefly cover two recent discriminative tracking methods: Multiple Instance Learning (MIL) tracker [4] and Tracking-Learning-Detection (TLD) [19].

### A. Subspace Based Trackers

As the main challenge in visual tracking is the difficulty in handling the appearance variability of a target object, it is imperative for a robust tracking algorithm to model such appearance variations. This can be difficult to accomplish when the object model is based on only a single image. Subspaces allow us to group images together and provide a single representation as a compact appearance model [30]. Subspace-based tracking originated with the work of Black and Jepson [7], where a subspace learning-based approach is proposed for tracking rigid and articulated objects. This approach uses a view-based eigenbasis representation with parameterised optical flow estimation. As the algorithm is based on iterative parameterised matching between the eigenspace and candidate image regions, it might have a relatively high computational load [22]. It also uses a single pre-trained subspace to provide the object appearance model across the entire video. As such, to achieve robust visual tracking with this method, it is necessary to first collect a large set of training images covering the range of possible appearance variations, which can be difficult to accomplish in practice.

Addressing the limitations of having a single representation for object appearance which is always learned off-line before tracking begins, Skocaj and Leonardis [34] developed a weighted incremental Principal Component Analysis (PCA) approach for sequentially updating the subspace. Although the method improves tracking accuracy, it has the limitation of being computationally intensive due to an optimisation problem that has to be computed iteratively. To address this issue, Li et al. [25] proposed an alternative incremental PCA-based algorithm for subspace learning. In this approach, the PCA model updating is performed directly using the previous eigenvectors and a new observation vector, thereby significantly decreasing the computational load of the update process.

Ho et al. [15] proposed an adaptive tracker using a uniform  $L_2$ -reconstruction error norm for subspace estimation, allowing explicit control on the approximation quality of the subspace. Empirical results show increases in tracking robustness and more swift reactions to environmental changes. However, as the method represents objects as a point in a linear subspace computed using only recent tracking results, the tracker may drift if large appearance changes occur [16].

Lim et al. [26] proposed a generalised tracking framework which constantly learns and updates a low dimensional subspace representation of the object. The updates are done using several observations at a time instead of a single observation. To estimate the object locations in consecutive frames, a sampling algorithm is used with robust likelihood estimates. The likelihood for each observed image being generated from a subspace is inversely proportional to the distance of that observation from the subspace. Ross et al. [30] improved the tracking framework in [26] by adding a forgetting factor to

focus more on recently acquired images and less on earlier observations during the learning and update stages.

Hu et al. [16] presented an incremental log-Euclidean Riemannian subspace learning algorithm in which covariance matrices of image features are mapped from a Riemannian manifold into a vector space, followed by linear subspace analysis. A block based appearance model is used to capture both global and local spatial layout information. Similar to traditional subspace based trackers, this method also uses a point-to-subspace distance.

### B. Other Generative Trackers

Among algorithms that do not use subspaces, two popular generative trackers are the mean shift tracker [9] and the fragments-based tracker [2]. The mean shift tracker models object appearance with colour histograms which can be applied to track non-rigid objects. Both the object model and candidate image areas are represented by colour pdfs, with the Bhattacharyya coefficient used as the similarity measure [18]. Tracking is accomplished by finding the local maxima of the similarity function using gradient information provided by the mean shift vector which always points toward the direction of maximum. While effective, the mean shift tracker is subject to several issues. First, the spatial information is lost, which precludes the application of more general motion models [2], [39]. Second, the Bhattacharyya coefficient may not be discriminative enough for tracking purposes [39]. Third, the method only maintains a single template to represent the object, leading to accuracy degradation if an object moves rapidly or if a significant occlusion occurs.

The fragments-based tracker [2] aims to handle partial occlusions via a parts-based model. The object is represented by multiple image fragments or patches. Spatial information is retained due to the use of spatial relationships between patches. Each patch votes on the possible positions and scales of the object in the current frame, by comparing its histogram with histograms of image patches in the frame. The tracking task is carried out by combining the vote maps of multiple patches by minimising a robust statistic. However, the object model is not updated and thereby it is not expected to handle tracking objects that exhibit significant appearance changes [37], [4].

### C. Discriminative Trackers

Two recent discriminative methods are the Multiple Instance Learning tracker (MILTrack) [4] and the Tracking-Learning-Detection (TLD) approach [19]. In the MILTrack approach, instead of using a single positive image patch to update the classifier, a set of positive image patches is maintained and used to update a multiple instance learning classifier [10]. In multiple instance learning, training examples are presented in sets with class labels provided for entire sets rather than individual samples. The use of sets of images allows the MILTrack approach to achieve robustness to occlusions and other appearance changes. However, if the object location detected by the current classifier is imprecise, it may lead to a noisy positive sample and consequently a suboptimal classifier update. These noisy samples can accumulate and cause tracking drift or failure [43].

The TLD approach decomposes the tracking task into three separate tasks: tracking, learning and detection. It regards tracking results as unlabelled and exploits their underlying

structure using positive and negative experts to select positive and negative samples for update. This method makes a common assumption in tracking that the training samples follow the same distribution as the candidate samples. Such an assumption is problematic if the object's appearance or background changes drastically or continuously, which causes the underlying data distribution to keep changing [23].

### III. PROPOSED TRACKING APPROACH

The proposed tracking approach is comprised of four intertwined components, listed below. To ease understanding of the overall system, we first overview the components below, and then provide the details for each component in the following subsections.

- (A) *Particle Filtering Framework.* An object's location in consecutive frames is parameterised as a distribution in a particle filter framework [3], where a set of particles represents the distribution and each particle represents a location. The location history of the tracked object in previous frames is taken into account to create a set of candidate object locations in a new frame.
- (B) *Particle Representation.* We represent the  $i$ -th particle at time  $t$  using an affine subspace  $\mathcal{A}_i^{(t)}$ , which is constructed by taking into account the appearance of the  $i$ -th candidate location at time  $t$  as well as the appearance of the tracked object in several immediately preceding frames. Each affine subspace  $\mathcal{A}_i^{(t)}$  is comprised of mean  $\mu_i^{(t)}$  and basis  $U_i^{(t)}$ .
- (C) *Bag of Affine Subspaces.* To take into account drastic appearance changes, the tracked object is modelled by a set of affine subspaces, which we refer to as bag  $\mathcal{B}$ . During tracking the bag first grows to a pre-defined size, and then its size is kept fixed by replacing the oldest affine subspace with the latest affine subspace.
- (D) *Comparing Affine Subspaces.* Each candidate subspace  $\mathcal{A}_i^{(t)}$  from the pool of candidates is compared to the affine spaces in bag  $\mathcal{B}$ . The most likely candidate subspace is deemed to represent the best particle, which in turn indicates the new location of the tracked object. The distance between affine subspaces is comprised of the distance between their means and the Grassmann geodesic distance between their bases.

#### A. Particle Filtering Framework

We aim to obtain the location  $x \in \mathcal{X}$ ,  $y \in \mathcal{Y}$  and the scale  $s \in \mathcal{S}$  of an object in frame  $t$  based on information obtained from previous frames. A blind search in the space of location and scale is inefficient, since not all possible combinations of  $x$ ,  $y$  and  $s$  are plausible. To efficiently search the location and scale space, we adapt a particle filtering framework [3], [42], where the object's location in consecutive frames is parameterised as a distribution. The distribution is represented using a set of particles, with each particle representing a location and scale.

Let  $z_i^{(t)} = [x_i^{(t)}, y_i^{(t)}, s_i^{(t)}]^T$  denote the state of the  $i$ -th particle comprised of the location and scale at time  $t$ . Using importance sampling [3], the density of the location and scale space (or most probable candidates) at time  $t$  is estimated as a set of  $N$  particles  $\{z_i^{(t)}\}_{i=1}^N$  using particles from the previous

frame  $\{z_i^{(t-1)}\}_{i=1}^N$  and their associated weights  $\{w_i^{(t-1)}\}_{i=1}^N$  (with constraints  $\sum_{i=1}^N w_i^{(t-1)} = 1$  and each  $w_i \geq 0$ ). For now we assume the associated weights of particles are known and later discuss how they can be determined.

To generate  $\{z_i^{(t)}\}_{i=1}^N$ ,  $\{z_i^{(t-1)}\}_{i=1}^N$  is first sampled (with replacement)  $N$  times. The probability of choosing  $z_i^{(t-1)}$ , the  $i$ -th particle at time  $t-1$ , is equal to the associated weight  $w_i^{(t-1)}$ . Each chosen particle then undergoes an independent Brownian motion, which is modelled by a Gaussian distribution. As a result, for a chosen particle  $z_i^{(t-1)}$ , a new particle  $z_i^{(t)}$  is obtained as a random sample from  $\mathcal{N}(z_i^{(t-1)}, \Sigma)$ , where  $\mathcal{N}(\mu, \Sigma)$  denotes a Gaussian distribution with mean  $\mu$  and diagonal covariance matrix  $\Sigma$ . The latter governs the speed of motion by controlling the location and scale variances.

#### B. Particle Representation via Affine Subspaces

To accommodate a degree of variations in object appearance, particle  $z_i^{(t)}$  is represented by an affine subspace  $\mathcal{A}_i^{(t)}$ , constructed from the appearance of the  $i$ -th candidate location at time  $t$  as well as the appearance of the tracked object in several immediately preceding frames. Each affine subspace  $\mathcal{A}_i^{(t)}$  can be described by a 2-tuple:

$$\mathcal{A}_i^{(t)} = \left\{ \mu_i^{(t)}, U_i^{(t)} \right\} \quad (1)$$

where  $\mu_i^{(t)} \in \mathbb{R}^D$  is the origin (mean) of the subspace and  $U_i^{(t)} \in \mathbb{R}^{D \times n}$  is the basis of the subspace. The parameter  $n$  is the number of basis vectors.

The subspace is obtained as follows. Let  $v(z_i^{(t)})$  represent the vectorised form of the  $i$ -th candidate image patch at time  $t$ . The top-left corner of the patch is indicated by  $(x_i^{(t)}, y_i^{(t)})$  and its size by  $s_i^{(t)}$ . The patch is resized to a fixed size of  $H_1 \times H_2$  pixels and represented as a column vector of size  $D = H_1 \times H_2$ . In the same manner, let  $v(z_*^{(t-1)})$  denote the vectorised form of the appearance of the tracked object at time  $(t-1)$ , with  $z_*^{(t-1)}$  denoting the particle that was deemed at time  $(t-1)$  to represent the tracked object. The vectorised forms of the candidate image patch as well as the patches containing the tracked object in the previous  $P$  frames are used to construct the following  $D \times (P+1)$  sized matrix:

$$\mathbf{V}_i^{(t)} = \left[ v(z_i^{(t)}), v(z_*^{(t-1)}), \dots, v(z_*^{(t-P)}) \right] \quad (2)$$

The subspace origin  $\mu_i^{(t)}$  is the mean of  $\mathbf{V}_i^{(t)}$ . The subspace basis  $U_i^{(t)}$  is obtained by performing singular value decomposition (SVD) of  $\mathbf{V}_i^{(t)}$  and choosing the  $n$  dominant left eigenvectors corresponding to the  $n$  largest eigenvalues.

#### C. Bag of Affine Subspaces

To take into account drastic appearance changes that might not be well modelled by subspaces, we propose to adapt the approach of keeping a history of object appearance variations [4], by modelling the tracked object via a set of affine subspaces obtained during the tracking process. We refer to such a set as a *bag* of affine subspaces, defined as:

$$\mathcal{B} = \{ \mathcal{A}_1, \dots, \mathcal{A}_K \} \quad (3)$$

where  $K$  is the number of subspaces to keep. The bag is updated every  $W$  frames by replacing the oldest affine subspace with the latest. The size of bag determines the memory of the tracking system.

To demonstrate the benefit of the bag approach, consider the following scenario. A person is being tracked, with the appearance of their whole body modelled as a single subspace. At some point a partial occlusion occurs, and only the upper body is visible for several frames. The tracker then learns the new occluded appearance of the person. If the tracker is only aware of the very last seen appearance (ie., the upper body), the tracker is likely to lose the object upon termination of the occlusion. Keeping a set of affine subspaces (ie., both upper body and whole body) increases memory of the tracked object and hence can help to overcome the confounding effect of drastic appearance changes.

#### D. Comparing Affine Subspaces

Each candidate subspace  $\mathcal{A}_i^{(t)}$  from the pool of candidates is compared to the affine spaces in bag  $\mathcal{B}$ . The most likely candidate subspace is deemed to represent the best particle, which in turn indicates the new location and scale of the tracked object.

The simplest distance measure between two affine subspaces is the minimal Euclidean distance, ie., the minimum distance of any pair of points of the two subspaces. However, such a measure does not form a metric [5] and it does not consider the angular distance between affine subspaces, which can be a useful discriminator [20]. On the other hand, using only the angular distance ignores the origin of affine subspaces and reduces the problem to a linear subspace case, which we wish to avoid.

To address the above limitations, we propose a distance measure with the following form:

$$\text{dist}(\mathcal{A}_i, \mathcal{A}_j) = \alpha \hat{d}_o(\boldsymbol{\mu}_i, \boldsymbol{\mu}_j) + (1 - \alpha) \hat{d}_g(\mathbf{U}_i, \mathbf{U}_j) \quad (4)$$

where  $\alpha \in [0, 1]$  is a mixing weight, while  $\hat{d}_o(\cdot, \cdot) \in [0, 1]$  is a normalised distance between the origins of the subspaces and  $\hat{d}_g(\cdot, \cdot) \in [0, 1]$  is a normalised Grassmann geodesic distance between bases of the subspaces.

We define the distance between the origins of  $\mathcal{A}_i$  and  $\mathcal{A}_j$  as:

$$\hat{d}_o(\boldsymbol{\mu}_i, \boldsymbol{\mu}_j) = \gamma \|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2 \quad (5)$$

where  $\gamma$  is a scaling parameter. Under the assumption that normalised images are used so that each pixel value is in the  $[0, 1]$  interval, the elements of  $\boldsymbol{\mu} \in \mathbb{R}^D$  are also in the  $[0, 1]$  interval. As such, the maximum value of the  $\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2$  component in Eqn. (5) is equal to  $D$ , and hence  $\gamma = 1/D$ .

A Grassmann manifold (a special type of Riemannian manifold) is defined as the space of all  $n$ -dimensional linear subspaces of  $\mathbb{R}^D$  for  $0 < n < D$  [1], [11], [13], [14], [29]. A point on Grassmann manifold  $\mathcal{G}_{D,n}$  is represented by an orthonormal basis through a  $D \times n$  matrix. The length of the shortest smooth curve connecting two points on a manifold is known as the geodesic distance. For Grassmann manifolds, the squared geodesic distance between subspaces  $\mathbf{E}$  and  $\mathbf{F}$  is given by:

$$d_g(\mathbf{E}, \mathbf{F}) = \|\Theta\|^2 \quad (6)$$

where  $\Theta = [\theta_1, \theta_2, \dots, \theta_n]$  is the principal angle vector, ie.

$$\cos(\theta_k) = \max_{\mathbf{e}_k \in \mathbf{E}, \mathbf{f}_k \in \mathbf{F}} \mathbf{e}_k^T \mathbf{f}_k \quad (7)$$

subject to  $\|\mathbf{e}_k\| = \|\mathbf{f}_k\| = 1$ ,  $\mathbf{e}_k^T \mathbf{e}_l = \mathbf{f}_k^T \mathbf{f}_l = 0$ ,  $l = 1, \dots, k-1$ . In other words, the first principal angle  $\theta_1$  is the smallest angle between all pairs of unit vectors in the two subspaces, with the remaining principal angles defined similarly. The principal angles can be computed through the SVD of  $\mathbf{E}^T \mathbf{F}$ , with the  $k$ -th singular value corresponding to  $\cos(\theta_k)$  [11], [1]. The principal angles have the property of  $\theta_i \in [0, \pi/2]$ . As such, the maximum value of  $d_g(\mathbf{E}, \mathbf{F})$  is  $n\pi^2/4$ . Therefore, we define the normalised squared Grassmann geodesic distance as

$$\hat{d}_g(\mathbf{E}, \mathbf{F}) = \beta d_g(\mathbf{E}, \mathbf{F}) \quad (8)$$

where  $\beta = 4/(n\pi^2)$ .

To measure the overall likelihood of a candidate affine subspace  $\mathcal{A}_i^{(t)}$  according to bag  $\mathcal{B}$ , the individual likelihoods of  $\mathcal{A}_i^{(t)}$  according to each affine subspace in  $\mathcal{B}$  are integrated using a straightforward sum rule [21], [32]:

$$p(\mathcal{A}_i^{(t)}|\mathcal{B}) = \sum_{k=1}^K \hat{p}(\mathcal{A}_i^{(t)}|\mathcal{B}[k]) \quad (9)$$

where  $\hat{p}(\mathcal{A}_i^{(t)}|\mathcal{B}[k])$  is the normalised likelihood and  $\mathcal{B}[k]$  indicates the  $k$ -th affine subspace in bag  $\mathcal{B}$ . In order to generate the new set of particles for a new frame, the overall likelihood for each particle is considered as the particle's weight. The likelihoods are normalised to sum to 1 using:

$$\hat{p}(\mathcal{A}_i^{(t)}|\mathcal{B}[k]) = \frac{p(\mathcal{A}_i^{(t)}|\mathcal{B}[k])}{\sum_{j=1}^N p(\mathcal{A}_j^{(t)}|\mathcal{B}[k])} \quad (10)$$

where  $N$  is the number of particles. The individual likelihoods are obtained using:

$$p(\mathcal{A}_i^{(t)}|\mathcal{B}[k]) = \exp\left(\frac{-\text{dist}(\mathcal{A}_i^{(t)}, \mathcal{B}[k])}{\sigma}\right) \quad (11)$$

where  $\sigma$  is a fixed parameter used to ensure that large distances result in low likelihoods. The most likely candidate subspace is deemed to represent the best particle, which in turn indicates the new location of the tracked object:

$$\mathbf{z}_*^{(t)} = \mathbf{z}_j^{(t)}, \quad \text{where } j = \underset{i}{\text{argmax}} p(\mathcal{A}_i^{(t)}|\mathcal{B}) \quad (12)$$

#### E. Computational Complexity

The computational complexity of the proposed tracking framework is dominated by generating a new affine subspace and comparing two subspaces. The subspace generation step requires  $O(Dn^2)$  operations by performing thin SVD [8]. Computing the geodesic distance between two points on Grassmann manifold  $\mathcal{G}_{D,n}$ , requires  $O(n^3 + Dn^2)$  operations for calculating the principal angles.

## IV. EXPERIMENTS

We evaluated the accuracy of the proposed method on eight commonly used challenging videos that have ground truth<sup>2</sup> for object locations: *Girl* [6], *Occluded Face* [2], *Occluded Face 2*, *Tiger 1*, *Tiger 2*, *Coke Can*, *Surfer*, and *Coupon Book* [4]. The videos contain various challenges such as object occlusions, impostor objects, pose variations, long-term appearance changes, illumination variations and non-stationary cameras. Example frames are shown in Fig. 3.

*Occluded Face* contains a face to be tracked with an occlusion challenge due to a book covering various parts of the face. *Occluded Face 2* also contains a face tracking task with occlusions, but includes long-term appearance changes due to the addition of a hat. The *Girl* sequence involves tracking a face with challenges such as severe pose variations and occlusion caused by another face, acting as a distractor. *Tiger 1* and *Tiger 2* contain a moving toy with many challenges such as frequent occlusions, pose variations, fast motion (which causes motion blur) and illumination changes. *Coupon Book* contains a book being moved around, with a very similar impostor book introduced to distract the tracker. *Coke Can* contains a specular object being moved around by hand, which is subject to occlusions, fast motion as well as severe illumination variations due to a lamp. *Surfer* involves tracking of the face of a surfer with many challenges such as non-stationary camera, pose variations and occlusion caused by waves.

Each video is composed of 8-bit grayscale images, resized to  $320 \times 240$  pixels. We used normalised pixel values (between 0 and 1) as image features. For the sake of computational efficiency in the affine subspace representation, we resized each candidate image region to  $32 \times 32$ , with the number of eigenvectors ( $n$ ) and number of previous frames ( $P$ ) set to 3 and 5, respectively. The number of particles ( $N$ ) is set to 100. Furthermore, we only consider 2D translation and scaling in the motion modelling component.

Based on preliminary experiments, a bag of size  $K = 10$  with the update rate  $W = 5$  is used. For the Brownian motion covariance matrix (Section III-A), the diagonal variances corresponding to the  $x$  location,  $y$  location and scale are set to  $5^2$ ,  $5^2$  and  $0.01^2$ , respectively. The parameter  $\sigma$  in Eqn. (11) is set to 0.01. We have kept the parameters fixed for all videos, to deliberately avoid optimising for any specific video. This is reflective of real-life conditions, where a tracker must work in various environments.

The source code for the proposed tracking algorithm is available at <http://arma.sourceforge.net/subspacetracker/>

### A. Quantitative Comparison

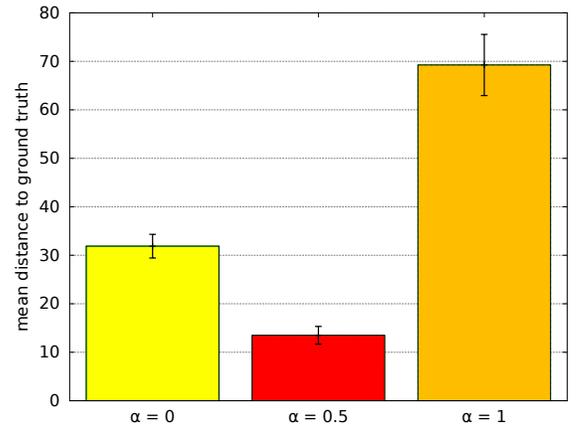
Following [4], we evaluated tracking error using the distance (in pixels) between the center of the bounding box around the tracked object and the ground truth. The mean of the distances over the eight videos is taken as the overall tracking error.

Fig. 1 shows the tracking error for three settings of  $\alpha$  in Eqn. (4).  $\alpha = 0$  ignores the origins and only uses the linear subspaces (ie.,  $\mu = 0$  for all models);  $\alpha = 0.5$  combines the origins and subspaces;  $\alpha = 1$  uses only the origins. Using  $\alpha = 0.5$  leads to considerably lower error than the other two settings, thereby indicating that use of the mean in conjunction with the subspace basis is effective.

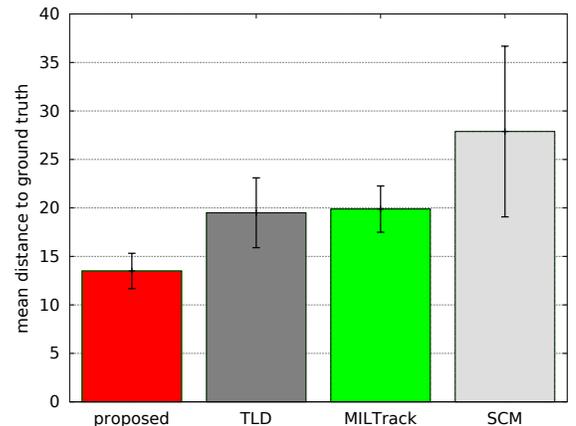
<sup>2</sup>The videos and the corresponding ground truth were obtained from [http://vision.ucsd.edu/~bbabenco/project\\_miltrack.html](http://vision.ucsd.edu/~bbabenco/project_miltrack.html)

Fig. 2 compares the tracking error of proposed tracker against three recent methods: Tracking-Learning-Detection (TLD) [19], Multiple Instance Learning Tracker (MILTrack) [4], and Sparsity-based Collaborative Model (SCM) [44]. For simplicity, the proposed tracker used  $\alpha = 0.5$  in Eqn. (4). Fig. 3 shows the resulting bounding boxes for several frames from the *Coupon Book*, *Surfer*, *Coke Can*, *Occluded Face 2*, and *Girl* videos. We use the publicly available source codes for MILTrack<sup>2</sup>, TLD<sup>3</sup>, and SCM<sup>4</sup>.

The proposed method obtains notably lower tracking error than TLD, MILTrack and SCM. Compared to TLD (the second best tracker), the mean distance to ground truth has decreased by more than 30%. Furthermore, the standard error of the mean [17] for the proposed tracker is considerably lower, indicating more consistent performance.



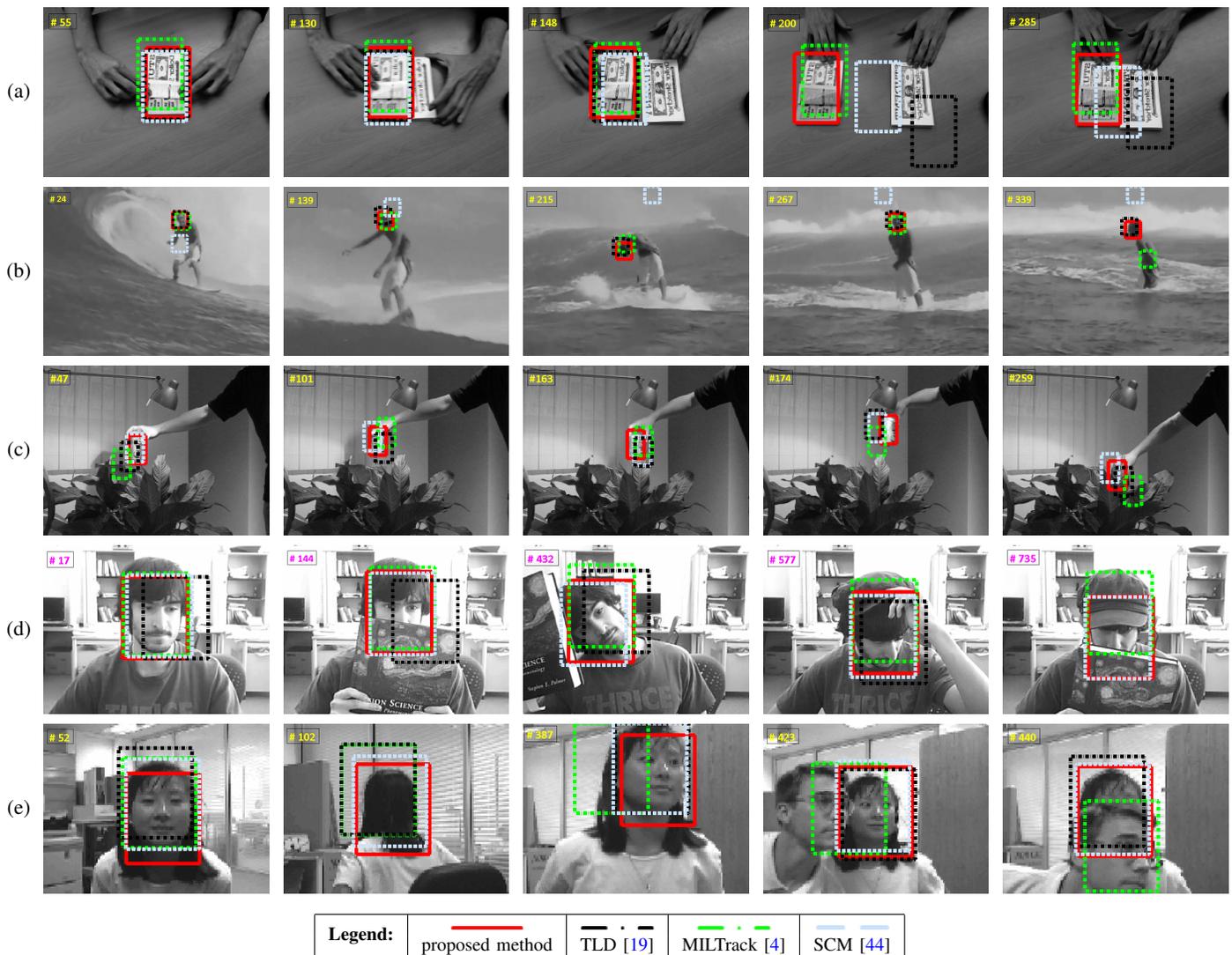
**Fig. 1:** Tracking error for various settings of  $\alpha$  in Eqn. (4). Tracking error is measured as the distance (in pixels) between the center of the bounding box around the tracked object and the ground truth. For each setting of  $\alpha$ , the mean of the distances over the eight videos is reported. The bars indicate the standard error of the mean [17].  $\alpha = 0$ : only the eigenbasis is used (ie. linear subspace),  $\alpha = 0.5$ : eigenbasis and mean (ie. affine subspace),  $\alpha = 1$ : mean only (origins of subspaces).



**Fig. 2:** Comparison of the proposed method against Tracking-Learning-Detection (TLD) [19], Multiple Instance Learning Tracking (MILTrack) [4], Sparsity-based Collaborative Model (SCM) [44]. Tracking error is measured as per Fig. 1.

<sup>3</sup><http://info.ee.surrey.ac.uk/Personal/Z.Kalal/tld.html>

<sup>4</sup>[http://ice.dlut.edu.cn/lu/Project/cvpr12\\_scm/cvpr12\\_scm.htm](http://ice.dlut.edu.cn/lu/Project/cvpr12_scm/cvpr12_scm.htm)



**Fig. 3:** Examples of bounding boxes resulting from tracking on several videos containing occlusions, distractors/impostors, pose variations and variable object illumination. Best viewed in colour. Frames from the following videos are shown: (a) *Coupon Book*, (b) *Surfer*, (c) *Coke Can*, (d) *Occluded Face 2* [4], and (e) *Girl* [6].

### B. Qualitative Comparison

On the *Coupon Book* video, TLD and SCM are confused by the distractor/impostor book. While MILTrack mostly stays with the original book, its accuracy is lower than the proposed method which consistently stays centered on the original book, unaffected by the impostor book. On the *Surfer* video, the proposed method and TLD consistently track the person’s face. This is in contrast to SCM which quickly loses track, and MILTrack which drifts towards the end of the video. On the *Coke Can* video, which contains dramatic illumination changes and rapid movement, MILTrack loses track after a part of the object is almost faded by the lamp light. SCM and TLD are affected to a lesser extent. In contrast, the proposed method consistently tracks the can, unaffected by the illumination variations. On the *Occluded Face 2* video, SCM and TLD lose accuracy due to confusion by occlusions, while SCM and the proposed method correctly track the face. On the *Girl* video, the proposed method and SCM manage to track the correct person throughout the whole video. TLD is affected by the severe pose variation (ie. the person turning around) but recovers when the

face appears frontal again. MILTrack loses track after the pose change and then tracks the distractor/impostor face. Overall, the qualitative observations agree with the quantitative results, with the proposed method achieving the lowest tracking error.

### V. MAIN FINDINGS AND FUTURE DIRECTIONS

In this paper we addressed the problem of object tracking subject to appearance changes due to occlusions as well as variations in illumination and pose. We proposed an adaptive tracking approach where the object is modelled as a continuously updated bag of affine subspaces, with each subspace constructed from the object’s appearance over several consecutive frames. The bag of affine subspaces takes into account drastic appearance changes that are not well modelled by individual subspaces, such as occlusions. Furthermore, during the search for the object’s location in a new frame, we proposed to represent the candidate image areas also as affine subspaces, by including the immediate tracking history over several frames. Distances between affine subspaces from the object model and candidate areas are obtained by exploiting

the non-Euclidean geometry of Grassmann manifolds. The use of bags of affine subspaces was embedded in a particle filtering framework.

Comparative evaluations on challenging videos against several recent discriminative trackers, such as Tracking-Learning-Detection [19] and Multiple Instance Learning Tracking [4], show that the proposed approach obtains notably better accuracy and consistency. The proposed approach also has the benefit of not requiring a separate training phase.

Future research directions include extending the bag update process to follow a semi-supervised fashion, where the effectiveness of a new learned affine subspace is used to determine whether the subspace should be added to the bag. Furthermore, the bag size and update rate can be dynamic, possibly dependent on the degree of tracking difficulty in challenging scenarios.

#### ACKNOWLEDGEMENTS

The Australian Centre for Robotic Vision is supported by the Australian Research Council via the Centre of Excellence program. NICTA is funded by the Australian Government through the Department of Communications, and the Australian Research Council through the ICT Centre of Excellence program.

#### REFERENCES

- [1] P.-A. Absil, R. Mahony, and R. Sepulchre. *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, 2008.
- [2] A. Adam, E. Rivlin, and I. Shimshoni. Robust fragments-based tracking using the integral histogram. *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 798–805, 2006.
- [3] M. Arulampalam, S. Maskell, N. Gordon, and T. Clapp. A tutorial on particle filters for on-line nonlinear/non-Gaussian Bayesian tracking. *IEEE Transactions on Signal Processing*, 50(2):174–188, 2002.
- [4] B. Babenko, M. Yang, and S. Belongie. Robust object tracking with online multiple instance learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(8):1619–1632, 2011.
- [5] R. Basri, T. Hassner, and L. Zelnik-Manor. Approximate nearest subspace search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(2):266–278, 2011.
- [6] S. Birchfield. Elliptical head tracking using intensity gradients and color histograms. *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 232–237, 1998.
- [7] M. J. Black and A. D. Jepson. Eigentracking: Robust matching and tracking of articulated objects using a view-based representation. *Int. Journal of Computer Vision*, 26(1):63–84, 1998.
- [8] M. Brand. Fast low-rank modifications of the thin singular value decomposition. *Linear Algebra and its Applications*, 415(1):20–30, 2006.
- [9] D. Comaniciu, V. Ramesh, and P. Meer. Kernel-based object tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(5):564–577, 2003.
- [10] T. G. Dietterich, R. H. Lathrop, and T. Lozano-Pérez. Solving the multiple instance problem with axis-parallel rectangles. *Artificial Intelligence*, 89(1-2):31–71, 1997.
- [11] A. Edelman, T. Arias, and S. Smith. The geometry of algorithms with orthogonality constraints. *SIAM Journal on Matrix Analysis and Applications*, 20(2):303–353, 1998.
- [12] A. W. Fitzgibbon and A. Zisserman. Joint manifold distance: a new approach to appearance based clustering. *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 26–33, 2003.
- [13] M. Harandi, R. Hartley, C. Shen, B. Lovell, and C. Sanderson. Extrinsic methods for coding and dictionary learning on Grassmann manifolds. *International Journal of Computer Vision*, 114(2):113–136, 2015. <http://dx.doi.org/10.1007/s11263-015-0833-x>
- [14] M. Harandi, C. Sanderson, S. Shirazi, and B. C. Lovell. Kernel analysis on Grassmann manifolds for action recognition. *Pattern Recognition Letters*, 34(15):1906–1915, 2013. <http://dx.doi.org/10.1016/j.patrec.2013.01.008>
- [15] J. Ho, K. Lee, M. Yang, and D. Kriegman. Visual tracking using learned linear subspaces. *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 782–789, 2004.
- [16] W. Hu, X. Li, W. Luo, X. Zhang, S. Maybank, and Z. Zhang. Single and multiple object tracking using log-Euclidean Riemannian subspace and block-division appearance model. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(12):2420–2440, 2012.
- [17] R. A. Johnson, I. Miller, and J. Freund. *Probability and Statistics for Engineers*. Pearson, 8th edition, 2010.
- [18] T. Kailath. The divergence and Bhattacharyya distance measures in signal selection. *IEEE Transactions on Communication Technology*, 15(1):52–60, 1967.
- [19] Z. Kalal, K. Mikolajczyk, and J. Matas. Tracking-learning-detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(7):1409–1422, 2012.
- [20] T. Kim, J. Kittler, and R. Cipolla. Discriminative learning and recognition of image set classes using canonical correlations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(6):1005–1018, 2007.
- [21] J. Kittler, M. Hatef, R. Duin, and J. Matas. On combining classifiers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(3):226–239, 1998.
- [22] M. La Cascia, S. Sclaroff, and V. Athitsos. Fast, reliable head tracking under varying illumination: An approach based on registration of texture-mapped 3d models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(4):322–336, 2000.
- [23] G. Li, Q. Huang, L. Qin, and S. Jiang. SSOCBT: A robust semisupervised online covboost tracker that uses samples differently. *IEEE Transactions on Circuits and Systems for Video Technology*, 23(4):695–709, 2013.
- [24] X. Li, W. Hu, Z. Zhang, X. Zhang, and G. Luo. Robust visual tracking based on incremental tensor subspace learning. *Int. Conference on Computer Vision (ICCV)*, pages 1–8, 2007.
- [25] Y. Li. On incremental and robust subspace learning. *Pattern Recognition*, 37(7):1509–1518, 2004.
- [26] J. Lim, D. A. Ross, R.-S. Lin, and M.-H. Yang. Incremental learning for visual tracking. *Advances in Neural Information Processing Systems*, pages 793–800, 2004.
- [27] H. Liu, S. Chen, and N. Kubota. Intelligent video systems and analytics: A survey. *IEEE Trans. on Industrial Informatics*, 9(3):1222–1233, 2013.
- [28] H. Lu, S. Lu, D. Wang, S. Wang, and H. Leung. Pixel-wise spatial pyramid-based hybrid tracking. *IEEE Transactions on Circuits and Systems for Video Technology*, 22(9):1365–1376, 2012.
- [29] Y. M. Lui. Advances in matrix manifolds for computer vision. *Image and Vision Computing*, 30(6-7):380–388, 2012.
- [30] D. A. Ross, J. Lim, R.-S. Lin, and M.-H. Yang. Incremental learning for robust visual tracking. *Int. Journal of Computer Vision*, 77(1-3):125–141, 2008.
- [31] C. Sanderson, M. Harandi, Y. Wong, and B. C. Lovell. Combined learning of salient local descriptors and distance metrics for image set face verification. *IEEE International Conference on Advanced Video and Signal-Based Surveillance (AVSS)*, pages 294–299, 2012. <http://dx.doi.org/10.1109/AVSS.2012.23>
- [32] C. Sanderson and K. K. Paliwal. Identity verification using speech and face information. *Digital Signal Processing*, 14(5):449–480, 2004.
- [33] S. Shirazi, M. T. Harandi, B. C. Lovell, and C. Sanderson. Object tracking via non-Euclidean geometry: A Grassmann approach. *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 901–908, 2014.
- [34] D. Skocaj and A. Leonardis. Weighted and robust incremental method for subspace learning. *Int. Conference on Computer Vision (ICCV)*, pages 1494–1501, 2003.
- [35] D. Wang, H. Lu, and M.-H. Yang. Least soft-threshold squares tracking. *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 2371–2378, 2013.
- [36] R. Wang, S. Shan, X. Chen, and W. Gao. Manifold-manifold distance with application to face recognition based on image set. *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 1–8, 2008.
- [37] S. Wang, H. Lu, F. Yang, and M.-H. Yang. Superpixel tracking. *Int. Conference on Computer Vision (ICCV)*, pages 1323–1330, 2011.
- [38] T. Wang, A. Backhouse, and I. Gu. Online subspace learning on Grassmann manifold for moving object tracking in video. *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 969–972, 2008.
- [39] C. Yang, R. Duraiswami, and L. Davis. Efficient mean-shift tracking via a new similarity measure. *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 176–183, 2005.
- [40] M. Yang, Z. Fan, J. Fan, and Y. Wu. Tracking nonstationary visual appearances by data-driven adaptation. *IEEE Transactions on Image Processing*, 18(7):1633–1644, 2009.
- [41] M.-H. Yang and J. Ho. Toward robust online visual tracking. *Distributed Video Sensor Networks*, pages 119–136. Springer, 2011.
- [42] A. Yilmaz, O. Javed, and M. Shah. Object tracking: A survey. *ACM Computing Surveys*, 38(4):1–45, 2006.
- [43] K. Zhang, L. Zhang, and M. Yang. Real-time object tracking via online discriminative feature selection. *IEEE Transactions on Image Processing*, 22(12):4664–4677, 2013.
- [44] W. Zhong, H. Lu, and M.-H. Yang. Robust object tracking via sparsity-based collaborative model. *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 1838–1845, 2012.