

# The New High-Performance Face Tracking System based on Detection-Tracking and Tracklet-Tracklet Association in Semi-Online Mode

Ngoc Q. Ly<sup>1</sup>, Tan T. Nguyen<sup>2</sup>, Tai C. Vong<sup>3</sup>  
Faculty of Information Technology  
VNUHCM-University of Science  
Ho Chi Minh City, Viet Nam

Cuong V. Than<sup>4</sup>  
AI Department  
Axon Company  
Seattle, USA

**Abstract**—Despite recent advances in multiple object tracking and pedestrian tracking, multiple-face tracking remains a challenging problem. In this work, the authors propose a framework to solve the problem in semi-online manner (the framework runs in real-time speed with two-second delay). The proposed framework consists of two stages: detection-tracking and tracklet-tracklet association. Detection-tracking stage is for creating short tracklets. Tracklet-tracklet association is for merging and assigning identifications to those tracklets. To the best of the authors' knowledge, the authors make contributions in three aspects: 1) the authors adopt a principle often used in online approaches as a part of the framework and introduce a tracklet-tracklet association stage to leverage future information; 2) the authors propose a motion affinity metric to compare trajectories of two tracklets; 3) the authors propose an efficient way to employ deep features in comparing tracklets of faces. The authors achieved 78.7% precision plot AUC, 68.1% success plot AUC on MobiFace dataset (test set). On OTB dataset, the authors achieved 78.2% and 72.5% precision plot AUC, 51.9% and 43.9% success plot AUC on normal and difficult face subsets, respectively. The average speed was maintained at around 44 FPS. In comparison to the state-of-the-art methods, the proposed framework's performance maintains high rankings in top 3 on two datasets while keeping the processing speed higher than the other methods in top 3.

**Keywords**—Face tracking; face re-identification; detection-tracking; tracklet-tracklet association

## I. INTRODUCTION

While multiple object tracking has been receiving much attention from researchers all over the world, multiple-face tracking has received much less attention due to two main reasons: face tracking is a sub-problem of object tracking thus many works focus on the general problem, and there is a lack of encompassing multiple-face tracking datasets. Therefore, multiple-face tracking remains a challenging problem. Recent advances in the field of multiple pedestrian tracking can be used to solve the problem of multiple-face tracking. There are two main research directions for the problem: online and offline.

Offline approaches [1]–[6] treat the problem as a global optimization one and solve it once having received all the information of all frames of a video. These approaches basically revolve in three stages:

Stage 1: Apply detection algorithms over all frames of the video to get detected bounding boxes of individuals, which are treated as nodes of a graph.

Stage 2: Define a meaningful metric to measure the relationship between two nodes of the graph by employing visual, spatial and temporal information.

Stage 3: Optimize an objective function globally to get clustered the bounding boxes of individuals.

These approaches tend to use commonly known detectors to generate all detection boxes (stage 1). However, these methods are different from each other in defining relations between nodes (stage 2) and objective functions (stage 3). Berclaz et al. [1] propose to model all potential locations over time, find trajectories that produce the minimum cost and track interacting objects simultaneously by using intertwined flow and imposing linear flow constraints. Milan et al. [2] employ an energy function that considers physical constraints such as target dynamics, mutual exclusion, and track persistence. Tang et al. [4] propose to jointly cluster detections over space and time by partitioning the graph with attractive and repulsive terms. Cruz et al. [6] introduce two lifted edges for the tracking graph that add additional long-range information to the objective. The authors of [6] also employ human pose features extracted from a deep network for the detection-detection association. Solving the problem with no constraints of speed while having all the information beforehand, offline approaches often produce higher accuracy than online approaches summarized as follows.

Online approaches mainly focus on tracking by detection [7]–[15]. Basically, they employ three models: a state-of-the-art detection model to produce face detection bounding boxes, a standalone tracker [16]–[19] to produce face track bounding boxes, and a deep feature model [20]–[26] to extract representative features for matching. Combining detection and tracking methods help alleviate challenges when using standalone trackers such as sudden movements, blurring, pose variation. By adopting the detection-tracking framework, the problem of face tracking is then reduced to data association [27], [28] problem, that is to assign detection boxes to track boxes. Data association [27], [28] between detection boxes and track boxes then can be reduced to the bipartite matching problem (assume no two detection boxes in one frame belong

to one individual, and so for track boxes) and can be efficiently solved by Hungarian algorithm [29]. Because bipartite matching algorithms find 1-1 matches, it is crucial to define a meaningful affinity metric, representing the relationship between two nodes, for good performance.

These online approaches can be simplified as follows:

Step 1: For each frame, run a detection model to get possible positions of faces in that frame (these results will be referred as detections). Then apply a deep feature model to extract features of these detections.

Step 2: Also, for that frame, run a tracker for each tracklet to get new possible positions from the previous position of each tracklet (these results will be referred as predictions). Then apply a deep feature model to extract features of these predictions.

Step 3: A defined metric is employed to relate detections with predictions. The metric consists of two parts: motion affinity and appearance affinity. Motion affinity is measured by the intersection over union (or Mahalanobis distance) of detections and predictions. Appearance affinity is measured by Euclidean (or cosine) distance between features of detections and features of predictions (or possibly of tracklets).

Step 4: After three steps above, the result is an affinity matrix (N detections x M predictions). Apply a bipartite matching algorithm to associate new detections with predictions. Unassigned detections are treated as new individuals while assigned detections are used to update tracklets.

Step 5: Repeat steps 1-4 consecutively for frames of a video.

There are some disadvantages to these online approaches.

Disadvantage 1: At the  $i$ -th frame, new detections must be assigned identifications at that frame. This means the information in the future cannot taken advantage of.

Disadvantage 2: To decide whether a new detection belongs to a known identity or is a new identity, the similarity matrix (computed by motion and appearance affinity) is used. To have the number of tracklets for one individual as low as possible, the threshold must be lowered. However, doing that way, the possibility of one track containing many individuals is high.

Disadvantage 3: Because detection-tracking method must run detection model and tracking algorithm for each frame to get new detections and new predictions, then run deep feature model (models used for feature extraction are computationally expensive) for new detections and new predictions, these models must be lightweight to run in real-time. This can lead to low accuracy in these models and causes errors for the whole framework.

Disadvantage 4: Because these approaches compare detections with predictions, they fail to employ very potential information that can be taken advantage of when comparing tracks to tracks. That is the fact that two temporal-overlapped tracks cannot belong to the same individual.

To resolve the issues stated above, the authors propose a semi-online framework for the multi-face tracking problem. The framework consists of two stages: detection-tracking stage and tracklet-tracklet association stage. For the detection-tracking stage, the authors employ the same principle as in online approaches with a modification: the authors use two complementary trackers (Kalman filter as a motion tracker and KCF (Kernelized Correlation Filter) as a visual tracker) to improve accuracy. For the tracklet-tracklet association, inspired by offline approaches, the authors treat each tracklet as a node of a graph and optimize the problem of assigning identifications globally. In this stage, the authors also introduce an efficient metric to compare two tracklets so that the framework can run with high speed.

The rest of this paper is organized as follows. In Related Works, the authors begin to cover current state-of-the-art methods for multiple-face tracking in two modes: offline and online. In Materials and Methods, the authors then turn to the proposed approach which is inspired by principles used in both offline and online multiple-face tracking. In this section, the authors illustrate the overview and detailed stages of the proposed framework. The authors conclude this section with contributions to literature. In Results and Discussions, the authors describe experiments and datasets, report experimental results, and discuss some implications. The final section concludes the proposed approach and considers ways to further improve multiple-face tracking.

## II. RELATED WORKS

### A. Offline Tracking

State-of-the-art methods for multi-face offline tracking are [30]–[32]. These approaches can be reduced to two main stages: tracklet creation (tracking-by-detection) and tracklet association. In [30], Zhang et al. first divide the video into many non-overlapping shots – music or film videos often contain many shots in different scenes. For each shot, the framework employs the tracking-by-detection paradigm to generate tracklets and merge those tracklets into groups by temporal, kinematic (motion, size) and appearance (deep feature) information. Then, Zhang et al. link tracklets across shots/scenes by treating each tracklet as a point, the appearance similarity between two tracklets as edge and applying the Hierarchical clustering algorithm to assign tracklets into groups. To increase the accuracy of the tracklet linking step, a discriminative feature extractor is needed. The authors of [30] introduce Learning Adaptive Discriminative Features whereby a deep extractor will be finetuned online based on samples from the video. Jin et al. [31] improve the performance of the mentioned method by using a more powerful detector (Faster R-CNN) in the tracking-by-detection stage and a more sophisticated tracklet association schedule. Lin et al. [32] push it further by applying body parts detector and introduce a co-occurrence model to generate longer tracklets when faces are out of camera (but body not) or detector cannot capture faces. Besides, the work also introduces a refinement scheme for tracklet association based on Gaussian Process.

B. Online Tracking

1) *Hand-crafted features*: One of the attempts to solve the multi-face online tracking problem that yield good results is [33]. In this work, Comaschi et al. adopt the tracking-by-detection mechanism for the pipeline (Fig. 1). Because of the frontal characteristics of the dataset being used, the work employs a Haar-like cascade face detector [34] to attain computational efficiency. In any tracking problem, the ability to learn appearance change and predict future states of objects is crucial for the model. Thus, the work introduces a structured SVM tracker that stores previous patterns and positions of an object and can predict the new state of an object based on current spatial and visual information. The tracker is updated online based on both track prediction and detection. In the data association step, this work applies Hungarian algorithm for the cost matrix computed by the intersection over union of detection boxes and track boxes.

Similar to the above work, Lan et al. [35] also adopt tracking-by-detection mechanism but with a more sophisticated tracker update routine. Naiel et al. [36] try to decrease the false negative rate (miss detection caused by a simple detector) of the previous pipeline without reducing speed. In this work, Naiel et al. adopt an advancement of [34] and a color-assisted tracker as detect and track components respectively (Fig. 2). The novelty of this work lies in the combined framework. Instead of running a detector for every frame like previous work, Naiel et al. propose a trigger mechanism so that the detector only need to run on some specific frames. Specifically, the detector is only triggered after a fixed interval (N frames) or earlier, when there is any tracking fail. The authors compare the histogram of the new track box with histograms of previous track boxes. If there is any large discrepancy, the track fail will trigger detection.

Similarly, the authors of [37] adopt the idea of sparse detection, modifies Viola-Jones detector in conjunction with a variant of optical flow to create a combined detection-tracking model.

2) *Deep features*: Recently, many works [38]–[42] integrate deep feature extractors into the tracking framework. Of those works, Chen et al. [38] adopt the sparse detection mechanism as described above and use KLT tracker [43] for the tracking-by-detection stage. In the data association step between detection boxes and track boxes, deep feature vectors are used as visual information in addition to spatial information.

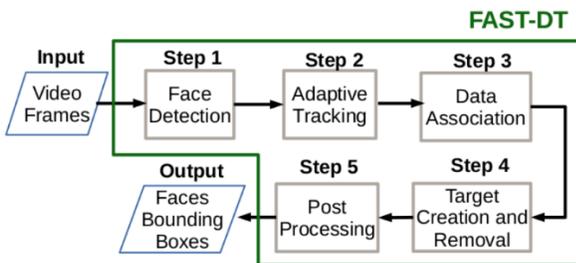


Fig. 1. Multi-Face Detection and Tracking Framework [33].

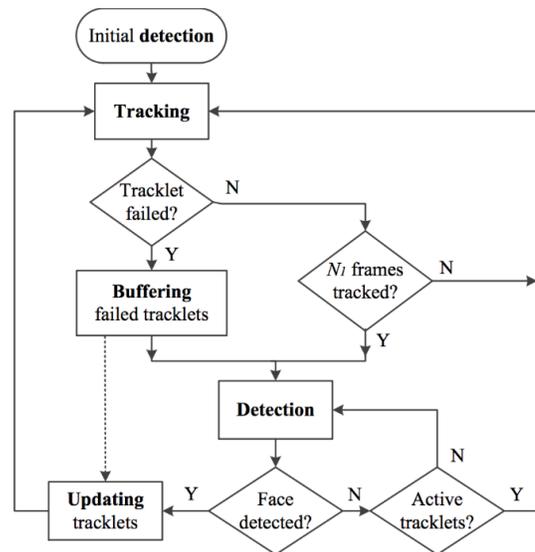


Fig. 2. Multi-Face Tracking Detection and Tracking Flow [36].

III. METHOD

A. Overview

1) *Semi-online tracking*: Aiming for practical usage and from the analysis of the online detection-tracking approaches, the authors propose a new approach in semi-online manner by introducing the tracklet-tracklet association stage (Fig. 3).

After getting the detections of a frame, the authors should match it with tracklets up until the previous frame to determine identifications for new detections. To achieve this criterion, using a deep feature extractor is a heavy waste. The authors propose a way to lighten the process while keeping the accuracy as high as possible. First, the authors use a light feature LBPH (Local Binary Pattern Histogram) extractor in the detection-tracking stage (Fig. 5) for efficient computation and combine it with information from a tracking method (Kalman filter) to reduce the errors as much as possible in creating short tracklets (the authors have not yet assigned identifications for those tracklets). Then the authors observe that consecutive face boxes of one tracklet are nearly the same, thus in the tracklet-tracklet association stage (Fig. 7), the authors introduce a compression method to get representatives of a tracklet and apply a deep feature extractor on these representatives instead of all boxes. The authors then link short tracklets into long tracklets by using those features as appearance information. In the linking step, the authors also introduce a new method for motion similarity between two tracklets. The tracklet-tracklet association stage resolved much problems stated above: the future information of frames sequences is well manipulated; the computational complexity is cut off from deep feature comparison by applying the new compression method.

Detection-Tracking stage: The main role of this stage is to extract the track information of targets in a frame using detecting and tracking methods. Technically, the detection-tracking stage processes frame-by-frame for every mini-batch interval (64 frames) and yields a list of tracklets. The process is illustrated in Fig. 4.

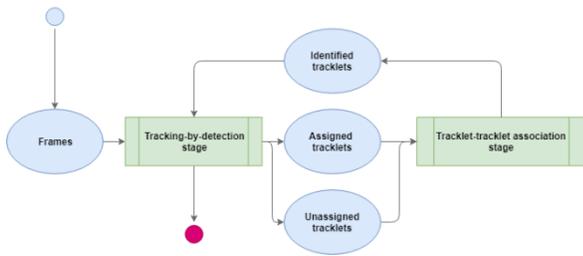


Fig. 3. Our Proposed Method. The Extra Tracklet-Tracklet Association is Introduced to Improve Accuracy by using more Information and Lighten the Process before.

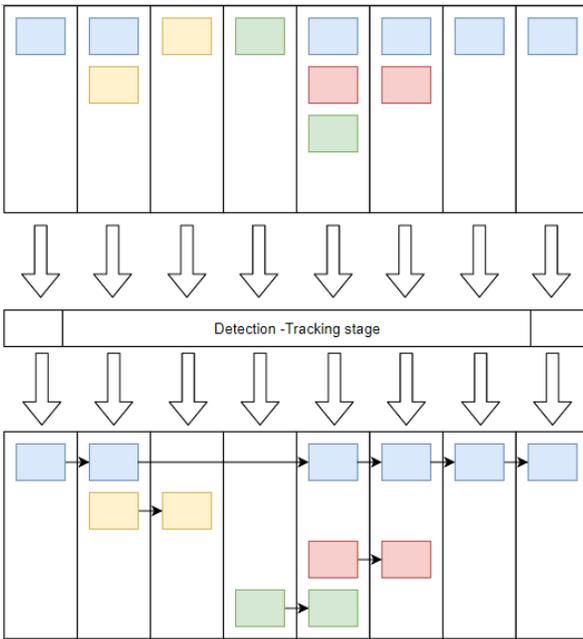


Fig. 4. Detection-Tracking Stage (Frame by Frame). Columns are Consecutive Frames; each Box is a Tracked Box in each Frame; the Arrows show how a Tracklet is Formed; Each identity is Marked by different Colors in Each Box.

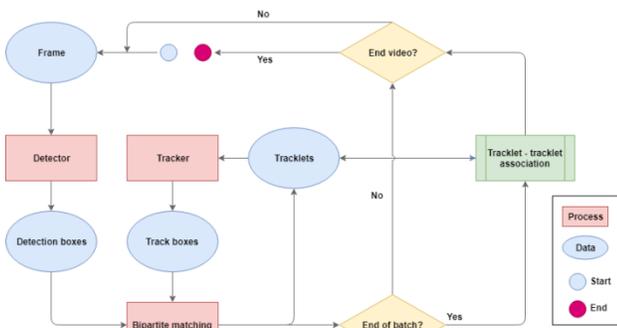


Fig. 5. Our Detection – Tracking flow Diagram.

The end-to-end framework consists of two stages:

**Tracklet-tracklet association stage:** At the end of each mini-batch process, the list of tracklets is passed to this stage. The main role of this stage is to correct false positives of the previous stage and connect related tracklet to create long tracklets and then assign identifications to these new tracklets. The process is shown in Fig. 6.

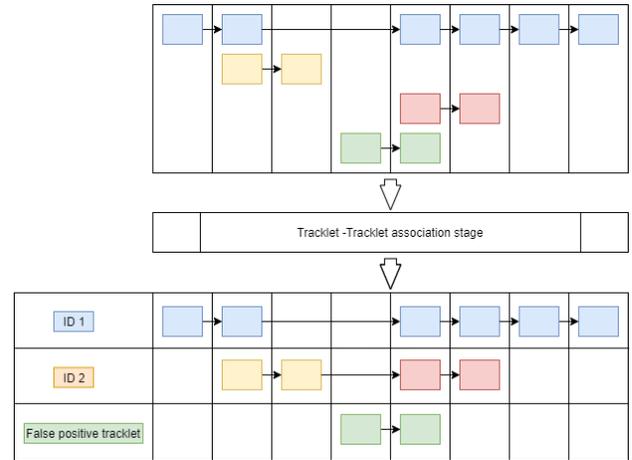


Fig. 6. Tracklet-Tracklet Association Stage. from Tracklets Formed before, the Identities will be Determined in this Stage.

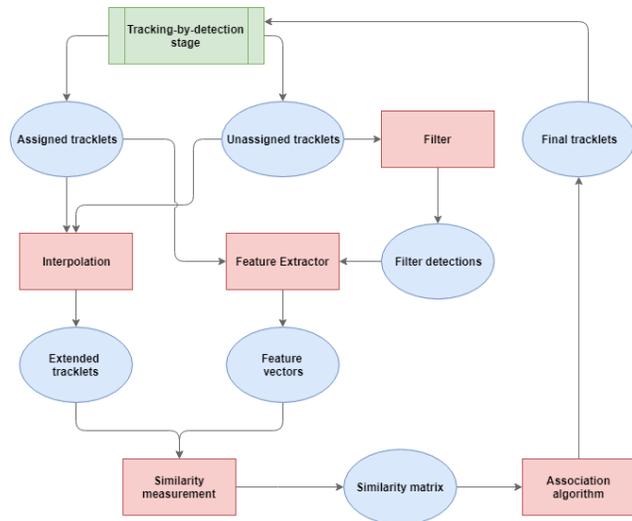


Fig. 7. Our Tracklet-Tracklet Association flow Diagram.

The proposed framework returns results after the tracklet-tracklet association stage. For instance, it returns results of frames 1-st to 64-th after seeing the information of frame 64-th. This induces a delay of over 2 seconds (64 frames ~ 2 seconds in normal 30fps videos). The details of the proposed framework are explained follow.

2) **Computational complexity:** The proposed framework can process video streaming in real-time. The speed can reach around 60fps, which is greater or equal the frequency of common videos (from 30 to 60fps).

3) **Detection-tracking stage:** The authors leverage known detection-tracking approaches with some modifications to speed up the stage without sacrificing much performance and introduce a new stage to improve the performance. The authors also implemented a framework: the detection-tracking stage combining S3FD face detector to produce detection boxes, LBPHs feature extractor to extract the global features, Kalman Filter tracker to produce tracking boxes, then Hungarian algorithms for matching the corresponding boxes to create tracklets.

4) *Tracklet-tracklet association stage*: The tracklet - tracklet association stage uses the motion information simulated by the spline interpolation and appearance information from FaceNet deep feature extractor to drop the false positives and match the suitable tracklets to accurately assign the ids for targets.

### B. Detection – Tracking Stage

1) *Goal*: In this stage, all the detection boxes of all frames in a batch will be grouped into short tracklets with the help of a single object tracking method.

2) *Principle*: Combining a single tracker and a detector helps a lot in overcoming the limitation of each single method. Using single trackers [16]–[19] to track faces in the wild situation is hard due to occlusion, illumination change, pose variation, sudden movement, etc. These issues can lead to track losses, inaccurate boxes (boxes that capture part of the face), incorrect boxes (boxes that capture the face of another individual). Moreover, using only a detector faces the appearance feature confusion if there are faces of different individuals with high appearance similarity.

The authors observe that detection models yield neater boxes than single trackers so using detection boxes as new information for updating single trackers is reasonable.

3) *Method*: In this stage, a detection model is used to generate possible bounding boxes of faces in a frame. During that time, a tracker is also used to predict a new possible bounding boxes positions from previous frames. Our detection-tracking algorithm will try to fuse these detection results with track results in order to better enhance the output, create more reliable tracklets.

At each frame, after running the detection and tracking process, the authors get a list of (N) detection boxes and (M) track boxes. The track boxes are the spatial predictions of bounding boxes from previous tracklets, while the detection boxes are the bounding boxes of faces that existed in that frame. Those faces may be the old faces from the previous frames, but they may also be the new faces that only exist from that frame. The main purpose of the detection-tracking algorithm is to define a meaningful affinity matrix (N x M) so that it can reflect the relationships between those detection boxes and track boxes.

Two features that are commonly leveraged are motion and appearance:

Motion affinity between a detection box and a track box is defined by the intersection over union (IoU) of them.

Appearance affinity between a detection box and a track box is defined by cosine affinity between LBPH features of them.

Those two features are used because for a pair of detection box and track box to be matched, two boxes should be close to each other with similar size and visual feature.

The authors define a gating unit for each affinity in order to filter out less likely matches. Because of our intention that if a

detection box and a track box are considered a possible match, they must satisfy motion affinity alone and appearance affinity alone first.

As explained, the authors want both metrics to be high to treat a pair of detection box and track box a likely match; thus, if both affinity metrics pass the threshold then the final affinity is the multiplicative result of motion and appearance affinity, otherwise is zero.

$$Match(i, j) = \begin{cases} s_m(i, j) \cdot s_a(i, j) & \text{if } s_m(i, j) > \gamma_M \text{ and } s_a(i, j) > \gamma_A \\ 0 & \text{else} \end{cases} \quad (1)$$

where,

$s_a(i, j)$  describes the appearance similarity distance between bounding boxes  $i$  and  $j$ , its range is from 0 to 1.

$s_m(i, j)$  describes the space similarity distance between bounding boxes  $i$  and  $j$ , its range is from 0 to 1.

$\gamma_M$  is the threshold for space similarity distance determined by heuristic (the authors reason that detection box and track box should be near to be of one individual, so the authors set this value to 0.3).

$\gamma_A$  is the threshold for appearance similarity distance determined by heuristic (the purpose of this stage is to create short tracklets, the authors use a high threshold to prevent wrong matches, specifically 0.9).

$Match(i, j)$  will be used to determine if a detection box and a track box is a possible match. It only has value if both motion and appearance metrics are over their thresholds. If one of the metrics is lower than its respective threshold,  $Match(i, j)$  is set to 0. The thresholds for  $Match(i, j)$  are determined through experiments (value search).

### C. Tracklet-Tracklet Association

1) *Goal*: Short tracklets from the detection-tracking stage are passed to this stage. The authors will group short tracklets into long tracklets and assign identifications for them. After this stage, the boxes in each frame will be marked with identifications and ready to deliver to the result stream.

2) *Principle*: The objective of face tracking is that for everyone existed in a video, the framework should output as few as possible the number of tracklets for that individual without wrongly including other faces of other individuals. This leads to the tradeoff mentioned in Section I. The authors tackle this with two principles:

Make sure the possibility of wrongly matching is as low as possible by using tight constraints (high affinity thresholds).

Adopt efficient motion and appearance affinity metrics between tracklets (different from track-detection) to group tracklets into identities based on a community discovery algorithm in this stage.

3) *Method*: After each batch processing the detection-tracking stage, the authors have a list of unknown-id tracklets that are needed to be assigned identifications in this stage. The authors also have a list of known-id tracklets in the past

(previous batches). Our job is now trying to assign identifications to unknown-id tracklets.

The authors formulate the assignment puzzle as an optimization problem. Each tracklet is treated as a node of a graph. The edge of two nodes indicates the affinity between the two. The authors then apply a clustering algorithm, in this situation, Leiden algorithm [28] on this graph in order to partition it into subgraphs – groups, each containing tracklets – nodes of the same individual. The authors put constraints so that each subgraph will not contain two known-id tracklets or two temporally overlapped tracklets. One of the essential parts of this stage is defining a meaningful metric representing the edge of two nodes. To do that, the authors adopt the complementary nature of motion and appearance.

a) Motion distance: For motion, the authors introduce a trajectory difference metric. Given two tracklets (t(i), t(j)), it is safe to assume that t(i) predate t(j) and there is no temporal overlap between two tracklets. From the boxes of t(i), the authors extrapolate forward to get the possible boxes in the future relative to t(i). From the boxes of t(j), the authors extrapolate backward to get the possible boxes in the past relative to t(j). For extrapolation, the authors assume that face movement can be modeled as a polynomial function and apply spline extrapolation. The authors ran model selection to determine the degree of movement and found that 1-degree spline performs best. Now the extrapolated parts of the two overlap temporally, the authors have a pair of overlapped extrapolated boxes in the same frame f(k). The authors now calculate a spatial distance between two boxes using two centers and a diagonal distance between two boxes according to their diagonals. The authors introduce a weight parameter to fuse the two distances into one unified box-box distance.

The box-box distance at frame k can be formulated in the following equation:

$$d_{M,k} = \lambda \cdot d_{S,k} + (1 - \lambda) \cdot d_{D,k} \quad (2)$$

In that,

$d_{S,k}$  is the Euclidean distance between two centers of two boxes.

$d_{D,k}$  is the diagonal distance between two boxes calculated by the difference in length between two diagonals.

$\lambda$  is the weight parameter to fuse above distances into one unified distance (the authors search from 0 to 1 with 0.1 interval and choose 0.4 to maximize area under the curve of success plot).

$d_{M,k}$  is the box-box distance at frame k the authors are going to obtain.

Then the trajectory distance is the average of pair distances:

$$d_M = \frac{1}{n-m+1} \sum_{k=m}^n d_{M,k} \quad (3)$$

where,

$k = m \rightarrow n$  are overlapped frame indices.

$d_{M,k}$  is the box-box distance at frame k.

$d_M$  is the trajectory distance, the average box-box distance over  $m - n + 1$  frames.

b) Appearance distance: For appearance, the authors use average Euclidean distance between two feature sets of two tracklets. For each box of a tracklet, the authors have a respective LBPHs feature (referred to as light feature) extracted from the detection-tracking stage. Assume t(i) have N light feature vectors and t(j) have M light feature vectors, one straightforward method is to compute N\*M Euclidean distances and use the average as the distance between two tracklets. However, the task is to distinguish between human faces, LBPHs feature is not discriminative enough for this task that requires fine-grained features. Besides, deep neural networks have outperformed hand-crafted methods on many visual tasks that require fine-grained features. Thus, the authors employ a deep feature extractor (Facenet) [20] for this task. Specifically, the authors deploy the pretrained model and feedforward to extract features.

However, deep feature extractors are computationally expensive and if the authors compute deep features for all boxes of a tracklet the framework would not run in real-time. Moreover, temporally adjacent boxes often contain similar information, so it would be redundant to compute all the deep features. The authors introduce our compression method to lower the number of boxes needed to be passed through a deep feature extractor using already computed light features.

Given a list of light feature vectors of a tracklet, the authors apply a clustering algorithm on these light feature vectors and pick out centroids, i.e.  $N_{compressed}$  boxes, for deep feature extraction. Only centroids are then passed to the deep feature extractor to extract 128-dimensional vectors. This way the authors save a lot of time computing deep features while keeping the diversity of a tracklet. The authors then use average Euclidean distance between two deep feature sets of two tracklets as tracklet - tracklet appearance distance:

$$d_A = \frac{1}{N_{compressed}} \cdot \frac{1}{M_{compressed}} \times \sum_{n1}^{N_{compressed}} \sum_{m1}^{M_{compressed}} Euclid(f(n), f(m)) \quad (4)$$

In that,

$M_{compressed}$  is the number of filtered boxes of the first track for deep feature extraction.

$N_{compressed}$  is the number of filtered boxes of the second track for deep feature extraction.

$d_A$  is our tracklet – tracklet appearance distance, calculated as the average Euclidean distance between two deep feature sets of two tracklets.

$f(n)$  is the feature extracted from the n-th box of  $N_{compressed}$  boxes.

$f(m)$  is the feature extracted from the m-th box of  $M_{compressed}$  boxes.

c) Fusing results: A weighted sum of appearance and motion affinities is the affinity between two tracklets (used as

the weight of the edge between two nodes). The authors fuse two affinities by taking the addition rather than multiplication as used in the detection-tracking stage because motion affinity is not reliable enough in case of long-term occlusion or camera shake. Thus, the authors set the weight for motion affinity low so that it plays as extra information.

$$d_{AM}(i, j) = \lambda \cdot d_M(i, j) + (1 - \lambda) \cdot d_A(i, j) \quad (5)$$

Where

$d_M(i, j)$  is the motion dissimilarity distance, calculated as explained.

$d_A(i, j)$  is the appearance dissimilarity distance, calculated as explained.

$\lambda$  is the weight parameter to adjust the importance of each distance. This value is determined through experiments (the authors search from 0 to 1 with 0.1 interval and choose 0.3 to maximize area under the curve for success plot).

$d_{AM}(i, j)$  is the dissimilarity distance of tracklet  $i$  and  $j$ .

#### D. Contributions

This proposed approach tackles challenges related to online approach above:

- Instead of computing deep features for all faces of one tracklet as online approaches do, the authors leverage light features (LBPHs) in the context of tracklet to efficiently compute deep features (extracted by deep network) without compromising representative power. In fact, the compressing method produces a more accurate representation for a tracklet thanks to diversity and high detection quality (high-score detected boxes).
- Using this framework, the authors can tighten the constraints in the tracking-by-detection stage so that the possibility of wrongly matching is low. Though having many tracklets after the tracking-by-detection stage, these tracklets will be grouped in the tracklet-tracklet association stage.
- The authors do not have to assign identifications to new detections right away in the detection-tracking stage but leave it to the tracklet-tracklet association stage. This way the authors can filter out false positives efficiently in the pre-processing step.
- The identification assignment step is tracklet-based; thus, the authors can take advantage of temporal information of tracklets (co-extant tracklets belong to different individuals).
- The authors also propose the trajectory difference metric to account for motion in tracklet-tracklet comparison.

In application, dataset is limited so using a pre-trained model and finetuning on small dataset is a reasonable choice. In this work, the authors show that simply adopting deep features (extracted by Facenet) and employ Euclidean (or cosine) metric is not discriminative enough in reference to real-life data. Therefore, the authors propose to apply Logistic

discriminant metric learning so that the new embedding space for real-life data is more discriminative.

The authors speculate that other regions of person, besides the face, also contain discriminating features. The authors tried to employ some color-based feature (color name) and texture-based feature (LOMO) but the results were not comparable, thus leaving this part for future work.

## IV. RESULTS AND DISCUSSIONS

Our experiments are conducted by python on the hardware GTX 1080 GPU, Intel(R) Xeon(R) CPU E5-2620 v4 @ 2.10GHz, 16GB RAM, while the MobiFace paper [44] used a desktop machine with Intel i9-7900X CPU (3.30GHz) and one GTX 1080 Ti GPU. Therefore, it's fair to compare the speed of our method versus other methods on MobiFace. For OTB dataset [45], RFTD method [46] used a setup with Intel Core i7 with 3.07GHz clock with no GPU and CXT and SCM used similar computational power, so the authors only compare the performance of our method versus other methods in terms of accuracy.

### A. The Purpose of Experiments on MobiFace and OTB Datasets

In order to prove the efficiency of our tracking framework, the authors conducted two comparisons:

Comparing single trackers with tracking-by-detection approaches through results from MobiFace Dataset. The purpose is to prove that integrate the detection method will enhance the result more than using a single tracker.

Comparing tracking-by-detection approaches with our approach through results from OTB Dataset. The purpose is to prove that using the light feature to process in the tracking-by-detection stage and using the deep feature in the tracklet - tracklet association stage in conjunction with motion affinity is a significant improvement.

#### 1) Experiments on MobiFace dataset

a) *About the dataset:* MobiFace dataset [44] is the first dataset for single face tracking in mobile situations. Due to the lack of engrossing face tracking datasets before MobiFace, the performance of pioneer face trackers was reported on a few videos or on small subsets of the OTB dataset, and the comparison between approaches was limited. The introduced dataset provides a unified benchmark with different attributes for future development in this field. Some samples of the dataset are illustrated in Fig. 8.

The authors collected 80 unedited live-streaming mobile videos captured by 70 different smartphone users in fully unconstrained environments and manually labeled over 95.000 bounding boxes on all frames. In order to cover typical usage of mobile device camera, the authors fetched videos from YouTube mobile live-streaming channels. Most of the videos are captured and uploaded under fully unconstrained environments without any extra video editing or visual effects. 6021 videos were collected and discarded under strict criteria that the target faces should appear at least in 10% of the video frames, and the target faces should not always stay still to serve the purpose of visual tracking. Besides the common 8 attributes

in object tracking datasets, the authors proposed six additional attributes commonly seen in mobile situations.

The authors also fine-tuned and improved a handful of state-of-the-art trackers and perform evaluations on the dataset. Through comparing with those results, the authors can evaluate the efficiency of our method.

*b) Setup the experiments:* Note that MobiFace dataset is designed for supervised trackers - an initial box of a targeted face is specified in the first frame. However, our method is designed to work in an unsupervised way (the authors do not need initial boxes) and can track multiple targets at a time. In order to adapt to the dataset, the authors must reduce the system to fit with the protocol of the dataset. Specifically, in the first frame of each video, the authors compare the detected result of our system with the initial box provided by the dataset to specify the targeted face and then return track results of that target only.

The video is only stored in YouTube so from the time the authors access it, the authors are unable to collect all videos from the dataset because some has been deleted by the owners.

The authors consider the three metrics proposed in the dataset: normalized precision, success rate, frames per second. As most of the metrics are in plot form, the authors will explain the way to extract an important metric from the plot, the area under the curve (AUC). With  $N$  is the number of thresholds used to draw the plot and  $n = 1, 2, 3, \dots, N$ . The curve was drawn from points with coordinate  $(t_n, f_n)$ ,  $t_n$  is the threshold value at that point and  $f_n$  is the evaluated value of our algorithm at that threshold, i.e. location error of precision plot, overlap score of success plot. The AUC is then calculated by

$$AUC = \sum_n (t_n - t_{n-1}) f_n \quad (6)$$

Normalised precision plot: Precision plot is a widely used evaluation metric for the tracking field. The precision is described as the location error, which is the Euclidean distance between the center location of the tracked face and the ground truth bounding box. This metric reflects how far the tracker has drifted from the targeted face. However, as the videos differ greatly in resolution, the authors of [44] adopt the recently proposed normalised precision value. The size of the frame is used for the normalisation, and the authors of [44] rank the trackers based on the area under the curve (AUC) for normalised precision value between 0 and 0.5.

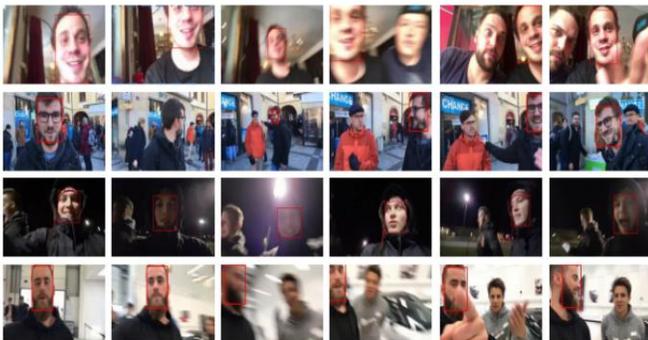


Fig. 8. Some Example Frame from the MobiFace Dataset [44]. Red ground Truth Bounding Boxes are Annotated by the Authors.

Success plot: Overlap score is also another commonly used metric in the tracking field. Given a ground truth bounding box  $r_{gt}$  of the target, the predicted bounding box of our algorithm is  $r_p$ . Then the overlap score can be computed by the intersection over union (IoU) of those two boxes as  $S = \frac{r_{gt} \cap r_p}{r_{gt} \cup r_p}$ , where the  $\cap$  and  $\cup$  represent the intersection and union of two rectangles, respectively. The success plot reflects the percentage of frames in which the intersection over union (IoU) of the predicted and ground truth bounding box is greater than a given threshold. Usually, the average success rate at 0.5 threshold is enough for evaluation. In addition, the area under the curve (AUC), which is the accumulated success rate can also be used for measurement. The authors can use those metrics interchangeably to summarize the performance.

Frames Per Second (FPS): the average speed of the evaluated tracker running across all the sequences. The initialization time is not considered. Because of the applicability concern, a mobile face tracker must be able to run at high speed (either on CPU or GPU) to allow maximum potential migration to actual mobile devices. Due to the lack of implementation of competitive trackers on mobile platforms, the authors can only use the FPS measured on the desktop environment, which indicate the relative efficiency of the trackers for evaluating and comparing.

*c) Experiment results:* Evaluation metrics of our method and state-of-the-art methods are illustrated in Fig. 9 and a detailed comparison is shown in Table I.

TABLE I. A DETAILED COMPARISON BETWEEN OUR METHOD AND MOBIFACE EVALUATED RESULTS

Tracker	Normalised Precision plot (AUC)	Success plot (AUC)	FPS
MDNet-MBF+R	<b>0.800</b>	<b>0.601</b>	1.79
MetaMDNet-MBF+R	0.767	<b>0.571</b>	1.03
MetaMDNet-YTF+R	0.744	0.566	1.06
MDNet-MBF	<b>0.772</b>	0.549	1.58
SiamFC-MBF+R	0.758	0.526	<b>53.14</b>
SiamFC-MBF	0.750	0.521	<b>81.54</b>
Proposed framework	<b>0.787</b>	<b>0.681</b>	<b>44.38<sup>a</sup></b>

<sup>a</sup> The authors profile the program and exclude reading image from disk time and writing image to disk time before calculating speed (details are in test.profile file in our source code).

*d) Discussion:* Because our approach is targeted for the multi-face tracking field. In order to make it work with the dataset, the authors run the framework over the dataset and get all tracks of targets in the video, then according to the initialized ground truth box, the authors define the target and return the target track results only. Because the dataset is from unconstrained environments with many existing faces, it is a noticeable effort of our tracker to avoid mistakes between tracklets and output the correct results.

As shown in the above plots, our method has an advantage in the success plot, but not the precision plot. Precision is affected by the Euclidean distance between the center of a ground truth bounding box and the center of a tracked box. Because high normalised error still treats a tracked box that

drifts out of a face (high Euclidean distance between two centers) as a true prediction, trackers that still maintain a track when the box drifts out of a face perform better with high normalized error. In the proposed framework, when the tracked box drifts out of a face, the algorithm terminates the tracklet instantly; therefore, with high normalised error, our tracker performs the same as with low normalised error while other trackers yield noticeably different results with different normalised errors.

The success plot might be more practical for applications that require high IoU between prediction boxes and ground truth boxes. The success plots of trackers evaluated in MobiFace dataset start very high, but the slope is very steep. Starting from above 0.8 success rate for threshold 0, to threshold 0.5, they drop to below 0.7 success rate. The steep slope indicates predicted boxes of those trackers are not always aligned with ground truth boxes. Our starting point is somewhere below 0.8 success rate but maintains the success rate over the overlap threshold change. At threshold 0.5, our approach still has a high success rate, above 0.7, indicating our boxes is closely aligned with ground truth boxes. At 0.5 threshold, the predicted boxes cover most of the track target and can be well used in application. Besides, as the main target of ours is for practical usages, a good success plot and success rate at 0.5 threshold - while keeping the speed - are acceptable.

## 2) Experiments on OTB (Object Tracking Benchmark) dataset

*a) About the dataset:* OTB Dataset [45] is one of the most famous datasets specifically used for benchmarking the object trackers since its appearance. The authors worked to collect and annotate most of the common tracking sequences from different datasets. They also classified those sequences into multiple categories by challenges as in Table II and selected 50 difficult and representative ones in the TB-50 dataset for an in-depth analysis. The full dataset contains more sequences of human (36 body and 26 face/head videos) than other categories because human target objects have the most practical usages, some samples of the dataset is illustrated in Fig. 10.

Before the introduction of MobiFace dataset, face tracking methods could only be evaluated on small self-collected datasets or a subset of OTB dataset. The whole dataset is designed for the object tracking algorithms, but the authors selectively pick out the sequences with faces to conduct experiments and compare with those methods mentioned before. The chosen face subset is described in Table III, the top 10 sequences are referred to as the difficult set and top 15 is the normal set [46].

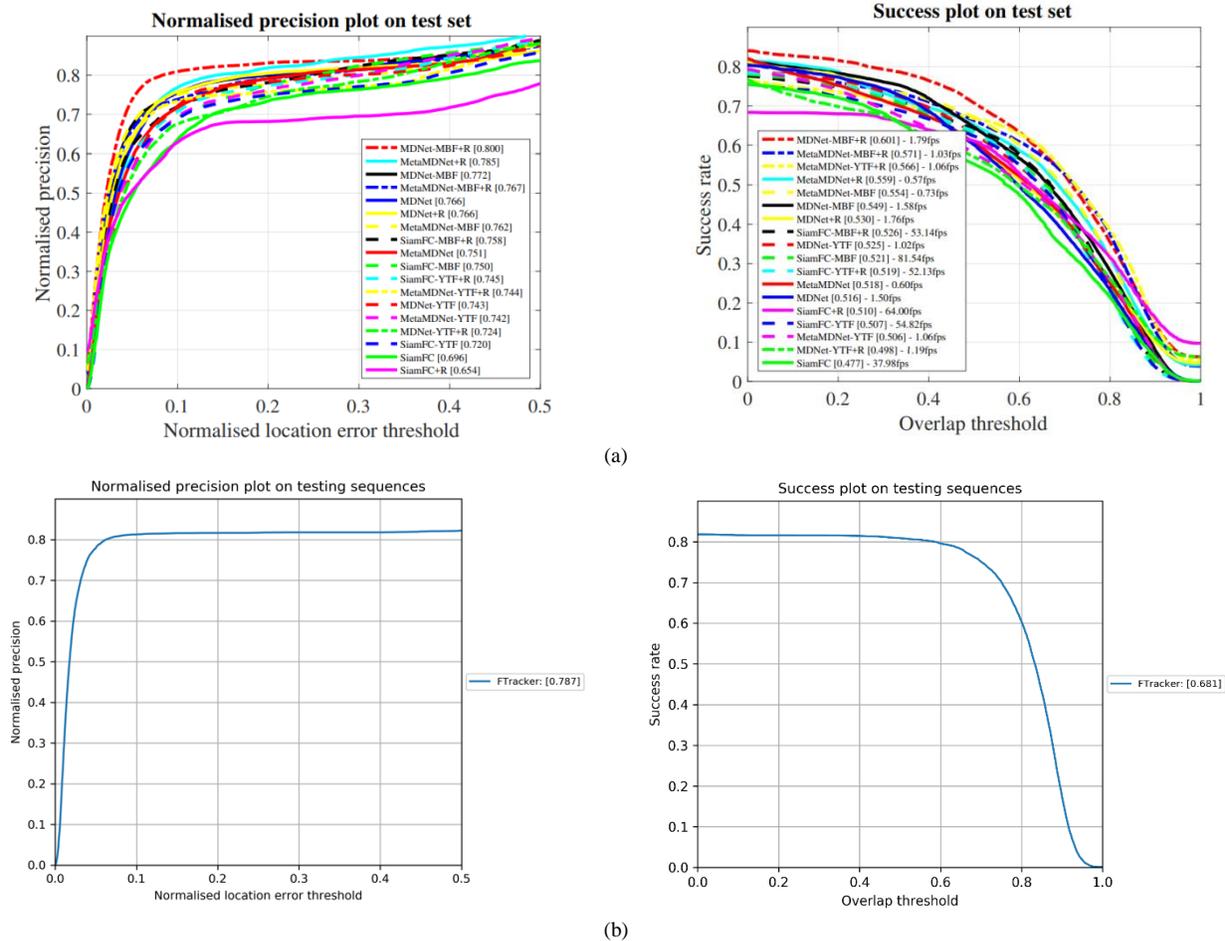


Fig. 9. Evaluation Results of Trackers on MobiFace Test Set: (a) Results from MobiFace Paper [44], (b) Results on our Method.

TABLE. II. ANNOTATED SEQUENCE ATTRIBUTES WITH THE THRESHOLD VALUES IN THE PERFORMANCE EVALUATION FROM OTB DATASET [45]

Attribute	Description
IV	Illumination Variation - The illumination in the target region is significantly changed
SV	Scale Variation - The ratio of the bounding boxes of the first frame and the current frame is out of range. $\left[\frac{1}{t_s}, t_s\right], t_s > 1 (t_s = 2)$
OCC	Occlusion - The target is partially or fully occluded.
DEF	Deformation - Non-rigid object deformation.
MB	Motion Blur - The target region is blurred due to the motion of the target or the camera.
FM	Fast Motion - The motion of the ground truth is larger than $t_m$ pixels ( $t_m = 20$ )
IPR	In-Plane Rotation - The target rotates in the image plane.
OPR	Out-of-Plane Rotation - The target rotates out of the image plane
OV	Out-of-View - Some portion of the target leaves the view
BC	Background Clutters - The background near the target has similar color or texture as the target
LR	Low Resolution - The number of pixels inside the ground-truth bounding box is less than $t_r$ ( $t_r = 400$ )

TABLE. III. ANNOTATED SEQUENCE ATTRIBUTES WITH THE THRESHOLD VALUES IN THE PERFORMANCE EVALUATION FROM OTB DATASET [45]

#	Sequence	Challenge
1	Soccer	IV, SV, OCC, MB, FM, IPR, OPR, BC
2	Freeman4	SV, OCC, IPR, OPR
3	Freeman1	SV, IPR, OPR
4	FleetFace	SV, DEF, MB, FM, IPR, OPR
5	Freeman3	SV, IPR, OPR
6	Girl	SV, OCC, IPR, OPR
7	Jumping	MB, FM
8	Trellis	IV, SV, IPR, OPR, BC
9	David	IV, SV, OCC, DEF, MB, IPR, OPR
10	Boy	SV, MB, FM, IPR, OPR
11	FaceOcc2	IV, OCC, IPR, OPR
12	Dudek	SV, OCC, DEF, FM, IPR, OPR, OV, BC
13	David2	IPR, OPR
14	Mhyang	IV, DEF, OPR, BC
15	FaceOcc1	OCC

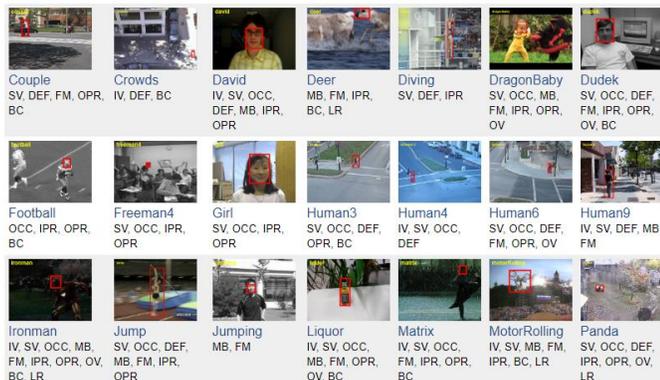


Fig. 10. Some Example Sequences from the OTB Dataset [45].

However, the dataset is also designed for the single object tracker. So, evaluation on this dataset also cannot reflect all the potential power of our system, but the authors can use that result to relatively compare with previous trackers in order to verify the power of the proposed framework.

b) *Set up the experiments:* Because the authors of MobiFace dataset inherit a lot of legacy from OTB dataset, in general, the setup stage and evaluation stage for OTB Dataset are the same as the MobiFace dataset.

c) *Experimental results:* Evaluation metrics of our method and state-of-the-art methods are illustrated in Fig. 11, Fig. 12, and a detailed comparison is shown in Table IV and Table V.

d) *Discussion:* The precision plots in Fig. 11 are good. The overall results are quite good, and the slope is shallow as predicted after witnessing above experiments. However, the authors have no data from other works to have an in-depth comparison.

TABLE. IV. TOP TRACKER COMPARISON ON OTB DATASET FACE SUBSET (NORMAL SET). EVALUATED RESULTS ARE FROM RFTD PAPER [46]

Face Tracker	Success Plot AUC	Success plot Threshold (0.5)
RFTD	55.2	<b>71.3</b>
Struck	<b>55.9</b>	67.6
SCM	<b>58.3</b>	<b>72.6</b>
ASLA	53.8	62.9
CSK	48.0	56.8
LIAPG	50.7	59.7
OAB	42.6	48.9
TLD	51.8	67.3
CXT	<b>57.3</b>	65.7
BSBT	40.6	47.0
<b>Our framework</b>	51.9	<b>68.3</b>

TABLE. V. TOP TRACKER COMPARISON ON OTB DATASET FACE SUBSET (DIFFICULT SET). EVALUATED RESULTS ARE FROM RFTD PAPER [46]

Face Tracker	Success Plot AUC	Success plot Threshold (0.5)
RFTD	<b>49.7</b>	<b>62.0</b>
Struck	45.2	51.7
SCM	<b>49.7</b>	<b>61.3</b>
ASLA	46.1	54.7
CSK	33.5	52.2
LIAPG	38.5	43.9
OAB	34.4	36.6
TLD	46.3	57.4
CXT	<b>48.2</b>	52.2
BSBT	29.0	29.7
<b>Proposed framework</b>	43.9	<b>59.7</b>

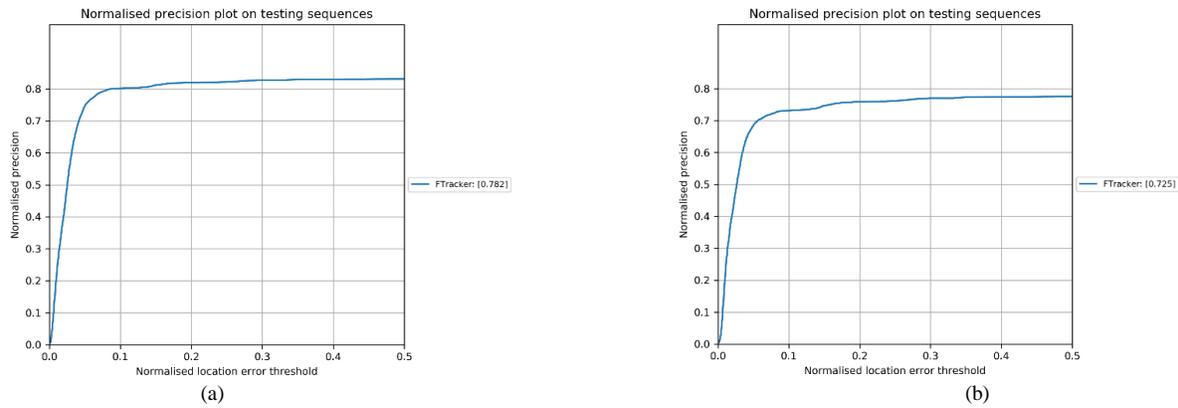


Fig. 11. Our Normalised Precision Plot on OTB Dataset Face Subsets (a) Normal Set (b) Difficult Set.

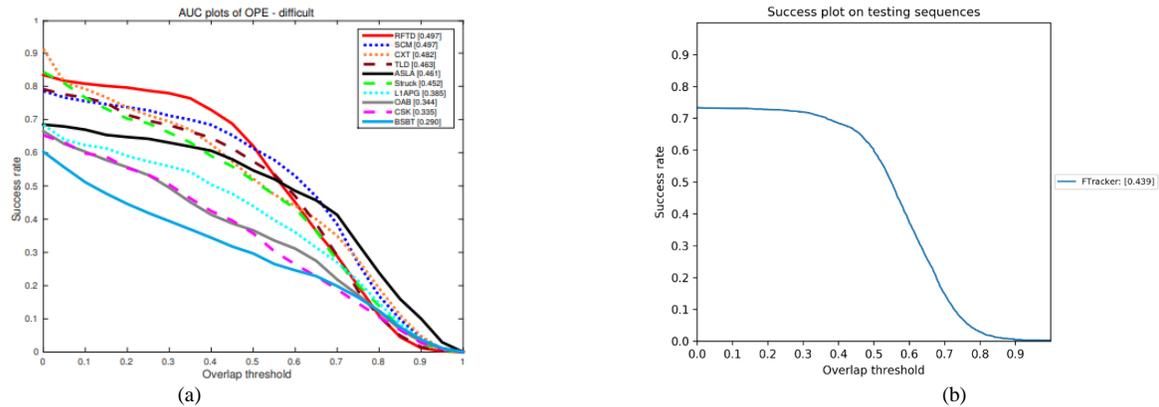


Fig. 12. Success Plots of Trackers on OTB Dataset Face Subset (difficult set): (a) Results from RFTD Paper[46] (b) Results on our Method.

As first sight from the metric Table IV and Table V, the proposed framework has average AUC while the slope of the proposed framework is also shallow as predicted. The main reason here is because when the predicted box is drifted from the face, the algorithm terminates the tracklet instantly; therefore, with high normalised error, our tracker performs the same as with low normalised error while other trackers yield noticeably different results with different normalised errors. The initial modest success rate leads to a modest average value. The success rate at threshold 0.5 is still good, ranking third in that section in both subsets.

## V. CONCLUSIONS

In this work, the authors proposed a method for face tracking problem in semi-online manner - the online process with some minor delay. The comparing experiments are conducted on two datasets: MobiFace dataset and OTB dataset with many state-of-the-arts works in the field. The results show that our method can produce robust accuracy while keeping a good speed. With that, the effectiveness of adding the tracklet-tracklet association stage after detection stage in semi-online manner is proven. The manipulation of appearance affinity and motion affinity have brought us the accuracy of the framework, while the workload division and information sharing of the two main stages make our process lighter and achieve better speed. With the improvements, all the disadvantages pointed out in Section I are solved.

The demonstrated framework has many advantages that can be applied to the production environment. First, the process as a whole was cut off to achieve the speed which is suitable for continuous streaming with a little delay. Second, the accuracy maintains at an acceptable value, which makes the proposed framework robust in many unconstraint environments. Finally, the framework can work without supervision, and is a high-performance multi-face tracking system.

There are many ways to develop from this work. First, because the framework consists of many components, researchers can try other combinations of related techniques (detector, tracker, feature extractor) to achieve better results. Second, the concept of semi-online tracking (use some delay for better results) can be applied to current work on face tracking.

## ACKNOWLEDGMENT

This research is funded by Viet Nam National University Ho Chi Minh City (VNUHCM) under grant no. B2018-18-01.

Thank to Axon company for the valuable support on cooperation.

## REFERENCES

- [1] J. Berclaz, F. Fleuret, E. Türetken, and P. Fua, "Multiple Object Tracking Using K-Shortest Paths Optimization," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, pp. 1806–1819, 2011.
- [2] A. Milan, S. Roth, and K. Schindler, "Continuous Energy Minimization for Multitarget Tracking," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 36, no. 1, pp. 58–72, Jan. 2014.

- [3] C. Kim, F. Li, A. Ciptadi, and J. M. Reh, "Multiple Hypothesis Tracking Revisited," 2015 IEEE Int. Conf. Comput. Vis. ICCV, pp. 4696–4704, 2015.
- [4] S. Tang, B. Andres, M. Andriluka, and B. Schiele, "Multi-Person Tracking by Multicut and Deep Matching," ArXiv E-Prints, p. arXiv:1608.05404, Aug. 2016.
- [5] L. Leal-Taixé, C. Canton Ferrer, and K. Schindler, "Learning by tracking: Siamese CNN for robust target association," ArXiv E-Prints, p. arXiv:1604.07866, Apr. 2016.
- [6] C. Cruz, L. Sucar, and E. Morales, "Real-Time face recognition for human-robot interaction," in Proceedings of the 8th IEEE International Conference on Automatic Face and Gesture Recognition, 2008, pp. 1–6.
- [7] A. V. Segal and I. D. Reid, "Latent Data Association: Bayesian Model Selection for Multi-target Tracking," 2013 IEEE Int. Conf. Comput. Vis., pp. 2904–2911, 2013.
- [8] A. Sadeghian, A. Alahi, and S. Savarese, "Tracking The Untrackable: Learning To Track Multiple Cues with Long-Term Dependencies," ArXiv E-Prints, p. arXiv:1701.01909, Jan. 2017.
- [9] Q. Chu, W. Ouyang, H. Li, X. Wang, B. Liu, and N. Yu, "Online Multi-Object Tracking Using CNN-based Single Object Tracker with Spatial-Temporal Attention Mechanism," ArXiv E-Prints, p. arXiv:1708.02843, Aug. 2017.
- [10] L. Chen, H. Ai, Z. Zhuang, and C. Shang, "Real-time Multiple People Tracking with Deeply Learned Candidate Selection and Person Re-Identification," ArXiv E-Prints, p. arXiv:1809.04427, Sep. 2018.
- [11] C. Kim, F. Li, and J. M. Reh, "Multi-object Tracking with Neural Gating Using Bilinear LSTM," in ECCV, 2018.
- [12] M. Thoreau and N. Kottege, "Improving Online Multiple Object tracking with Deep Metric Learning," ArXiv E-Prints, p. arXiv:1806.07592, Jun. 2018.
- [13] Y. Yoon, A. Boragule, Y. Song, K. Yoon, and M. Jeon, "Online Multi-Object Tracking with Historical Appearance Matching and Scene Adaptive Detection Filtering," ArXiv E-Prints, p. arXiv:1805.10916, May 2018.
- [14] N. Narayan, N. Sankaran, S. Setlur, and V. Govindaraju, "Re-identification for Online Person Tracking by Modeling Space-Time Continuum," 2018 IEEE CVF Conf. Comput. Vis. Pattern Recognit. Workshop CVPRW, pp. 1519–1519, 2018.
- [15] S. Zhang et al., "Improved Selective Refinement Network for Face Detection," ArXiv E-Prints, p. arXiv:1901.06651, Jan. 2019.
- [16] R. Kalman, "A New Approach to Linear Filtering and Prediction Problems," Trans. ASME - J. Basic Eng., vol. 82, pp. 35–45, 1960.
- [17] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista, "High-Speed Tracking with Kernelized Correlation Filters," ArXiv E-Prints, p. arXiv:1404.7584, Apr. 2014.
- [18] D. Held, S. Thrun, and S. Savarese, "Learning to Track at 100 FPS with Deep Regression Networks," in Unknown, 2016, vol. 9905, pp. 749–765.
- [19] L. Bertinetto, J. Valmadre, J. F. Henriques, A. Vedaldi, and P. H. S. Torr, "Fully-Convolutional Siamese Networks for Object Tracking," ArXiv E-Prints, p. arXiv:1606.09549, Jun. 2016.
- [20] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embedding for face recognition and clustering," in 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2015, pp. 815–823.
- [21] Z. Zhong, L. Zheng, Z. Zheng, S. Li, and Y. Yang, "Camera Style Adaptation for Person Re-identification," ArXiv E-Prints, p. arXiv:1711.10295, Nov. 2017.
- [22] J. Zhuo, Z. Chen, J. Lai, and G. Wang, "Occluded Person Re-identification," ArXiv E-Prints, p. arXiv:1804.02792, Apr. 2018.
- [23] Y. Suh, J. Wang, S. Tang, T. Mei, and K. M. Lee, "Part-Aligned Bilinear Representations for Person Re-identification," ArXiv E-Prints, p. arXiv:1804.07094, Apr. 2018.
- [24] M. M. Kalayeh, E. Basaran, M. Gokmen, M. E. Kamasak, and M. Shah, "Human Semantic Parsing for Person Re-identification," ArXiv E-Prints, p. arXiv:1804.00216, Mar. 2018.
- [25] J. Almazán, B. Gajic, N. Murray, and D. Larlus, "Re-ID done right: towards good practices for person re-identification.," CoRR, vol. abs/1801.05339, 2018.
- [26] H. Wang et al., "CosFace: Large Margin Cosine Loss for Deep Face Recognition," ArXiv E-Prints, p. arXiv:1801.09414, Jan. 2018.
- [27] M. Keuper, E. Levinkov, N. Bonneel, G. Lavoué, T. Brox, and B. Andres, "Efficient Decomposition of Image and Mesh Graphs by Lifted Multicuts," ArXiv E-Prints, p. arXiv:1505.06973, May 2015.
- [28] V. Traag, L. Waltman, and N. J. van Eck, "From Louvain to Leiden: guaranteeing well-connected communities," ArXiv E-Prints, p. arXiv:1810.08473, Oct. 2018.
- [29] H. W. Kuhn, "The Hungarian Method for the Assignment Problem," in 50 Years of Integer Programming, 2010.
- [30] S. Zhang et al., "Tracking Persons-of-Interest via Adaptive Discriminative Features," in Computer Vision – ECCV 2016, Cham, 2016, pp. 415–433.
- [31] S. Jin, H. Su, C. Stauffer, and E. Learned-Miller, "End-to-End Face Detection and Cast Grouping in Movies Using Erdős-Rényi Clustering," in ArXiv e-prints, 2017, pp. 5286–5295.
- [32] C. Lin and Y. Hung, "A Prior-Less Method for Multi-face Tracking in Unconstrained Videos," in 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp. 538–547.
- [33] F. Comaschi, S. Stuijk, T. Basten, and H. Corporaal, "Online multi-face detection and tracking using detector confidence and structured SVMs," in 2015 12th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS), 2015, pp. 1–6.
- [34] P. Viola and M. J. Jones, "Robust Real-Time Face Detection," Int. J. Comput. Vis., vol. 57, no. 2, pp. 137–154, May 2004.
- [35] M. Naiel, M. O. Ahmad, M. N. s Swamy, J. Lim, and M.-H. Yang, "Online Multi-Object Tracking via Robust Collaborative Model and Sample Selection," Comput. Vis. Image Underst., vol. 154, 2016.
- [36] X. Lan, Z. Xiong, W. Zhang, S. Li, H. Chang, and W. Zeng, "A super-fast online face tracking system for video surveillance," in 2016 IEEE International Symposium on Circuits and Systems (ISCAS), 2016, pp. 1998–2001.
- [37] A. Ranftl, F. Alonso-Fernandez, S. Karlsson, and J. Bigun, "Real-time AdaBoost cascade face tracker based on likelihood map and optical flow," IET Biom., vol. 6, no. 6, pp. 468–477, 2017.
- [38] J. Chen, R. Ranjan, A. Kumar, C. Chen, V. M. Patel, and R. Chellappa, "An End-to-End System for Unconstrained Face Verification with Deep Convolutional Neural Networks," in 2015 IEEE International Conference on Computer Vision Workshop (ICCVW), 2015, pp. 360–368.
- [39] M. Hayat, S. H. Khan, N. Werghi, and R. Goecke, "Joint Registration and Representation Learning for Unconstrained Face Identification," in 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017, pp. 1551–1560.
- [40] N. Crosswhite, J. Byrne, C. Stauffer, O. Parkhi, Q. Cao, and A. Zisserman, "Template Adaptation for Face Verification and Identification," in 2017 12th IEEE International Conference on Automatic Face Gesture Recognition (FG 2017), 2017, pp. 1–8.
- [41] R. Ranjan et al., "A Fast and Accurate System for Face Detection, Identification, and Verification," IEEE Trans. Biom. Behav. Identity Sci., vol. 1, pp. 82–96, 2018.
- [42] Y. Wang, J. Shen, S. Petridis, and M. Pantic, "A real-time and unsupervised face Re-Identification system for Human-Robot Interaction," Pattern Recognit. Lett., 2018.
- [43] Jianbo Shi and Tomasi, "Good features to track," in 1994 Proceedings of IEEE Conference on Computer Vision and Pattern Recognition, 1994, pp. 593–600.
- [44] Y. Lin, S. Cheng, J. Shen, and M. Pantic, "MobiFace: A Novel Dataset for Mobile Face Tracking in the Wild," ArXiv E-Prints, p. arXiv:1805.09749, May 2018.
- [45] Y. Wu, J. Lim, and M.-H. Yang, "Object Tracking Benchmark," IEEE Trans. Pattern Anal. Mach. Intell., vol. 37, pp. 1–1, 2015.
- [46] F. Comaschi, S. Stuijk, T. Basten, and H. Corporaal, "Robust online face tracking-by-detection," in 2016 IEEE International Conference on Multimedia and Expo (ICME), 2016, pp. 1–6.