1  # Computational Intelligence and Neuroscience

2  # The New High-Performance Face Tracking System based on

3  # Detection-Tracking and Tracklet-Tracklet Association in Semi-

4  # Online mode

5  NGOC Q. LY,[1] Tan T. Nguyen,[1,2] Tai C. Vong,[1] and Cuong V. Than [2]

6  [1] Faculty of Information Technology, VNUHCM-University of Science, 700000, Vietnam.
7  [2] AI Department, Axon, Ho Chi Minh, 700000, Vietnam.

8  Correspondence should be addressed to Tan T. Nguyen; ttnguyen@axon.com

9  ## Abstract

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

26  ## Introduction

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

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

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

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

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

43  These approaches tend to use commonly known detectors to generate all detection boxes
44  (stage 1). However, these methods are different from each other in defining relations between
45  nodes (stage 2) and objective functions (stage 3). [1] proposes to model all potential locations
46  over time, find trajectories that produce the minimum cost and track interacting objects
47  simultaneously by using intertwined flow and imposing linear flow constraints. [2] employs
48  an energy function that considers physical constraints such as target dynamics, mutual
49  exclusion, and track persistence. [4] proposes to jointly cluster detections over space and
50  time. The optimal number of people as well as the cluster of each person are obtained by
51  partitioning the graph with attractive and repulsive terms. [6] introduces two types of edges
52  (regular and lifted edges) for the tracking graph: the regular edges define the set of feasible
53  solutions in the graph and the lifted edges add additional long-range information to the
54  objective. The authors of [6] also employ human pose features extracted from a deep network
55  for the detection-detection association. Solving the problem with no constraints of speed
56  while having all the information beforehand, offline approaches often produce higher
57  accuracy than online approaches summarized as follows.

58  Online approaches mainly focus on tracking by detection [7]–[15]. Basically, they employ
59  three models: a state-of-the-art detection model to produce face detection bounding boxes, a
60  standalone tracker [16]–[19] to produce face track bounding boxes, and a deep feature model
61  [20]–[26] to extract representative features for matching. Combining detection and tracking
62  methods help alleviate challenges when using stand-alone trackers such as sudden
63  movements, blurring, pose variation. Specifically, because detection boxes are often neater
64  (close to faces) than track boxes while track boxes contain spatial and motion information of
65  the objects, fusing detection boxes and track boxes help mitigate the problem of accumulated
66  error in trackers. By adopting the detection-tracking framework, the problem of face tracking
67  is then reduced to data association [27], [28] problem, that is to assign detection boxes to
68  track boxes. Data association [27], [28] between detection boxes and track boxes is then can
69  be reduced to the bipartite matching problem (we assume no two detection boxes in one
70  frame belong to one individual, and so for track boxes) and can be efficiently solved by
71  Hungarian algorithm [29]. Because bipartite matching algorithms find 1-1 matches, it is
72  crucial to define a meaningful affinity metric, representing the relationship between two
73  nodes, for good performance.

74  These online approaches can be simplified as follows:

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

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

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

86  **Step 4**: After three steps above, we now have an affinity matrix (N detections x M
87  predictions). We apply a bipartite matching algorithm to associate new detections with
88  predictions. Unassigned detections are treated as new individuals while assigned detections
89  are used to update tracklets.

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

91  There are some disadvantages to these online approaches.

92  **Disadvantage 1**: At the ith frame we must assign identities to new detections at that frame.
93  This means we cannot take advantage of the information in the future.

94  **Disadvantage 2**: To decide whether a new detection belongs to a known identity or is a new
95  identity, we rely on the similarity matrix (computed by motion and appearance affinity)
96  thresholded by a scalar value. To have as few the number of tracklets for one individual as
97  possible, we must lower the threshold. However, doing that way, the possibility of one track
98  containing many individuals is high.

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

104  **Disadvantage 4**: Because these approaches compare detections with predictions, they fail to
105  employ very potential information that we can take advantage when we compare tracks with
106  tracks. That is the fact that two temporal-overlapped tracks cannot belong to the same
107  individual.

108  To resolve the issues stated above, we propose a semi-online framework for the multi-face
109  tracking problem. The framework consists of two stages: detection-tracking stage and
110  tracklet-tracklet association stage. For the detection-tracking stage, we employ the same
111  principle as in online approaches with a modification: we use two complementary trackers
112  (Kalman filter as a motion tracker and KCF (Kernelized Correlation Filter) as a visual
113  tracker) to improve accuracy. For the tracklet-tracklet association, inspired by offline
114  approaches, we treat each tracklet as a node of a graph and optimize the problem of assigning

115 identifications globally. In this stage, we also introduce an efficient metric to compare two
116 tracklets so that the framework can run with high speed.

## Materials and Methods

### *Related works*

### Offline tracking

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

### Online tracking

### *Hand-crafted features*

139 One of the attempts to solve the multi-face online tracking problem that yield good results is
140 [33]. In this work, the authors adopt the tracking-by-detection mechanism for the pipeline
141 (Figure 1). Because of the frontal characteristics of the dataset being used, the authors
142 employ a Haar-like cascade face detector [34] to attain computational efficiency. In any
143 tracking problem, the ability to learn appearance change and predict future states of objects is
144 crucial for the model. Thus, the authors introduce a structured SVM tracker that store
145 previous patterns and positions of an object and can predict the new state of an object based
146 on current spatial and visual information. The tracker is updated online based on both track
147 prediction and detection. In the data association step, this work applies Hungarian algorithm
148 for the cost matrix computed by the intersection over union of detection boxes and track
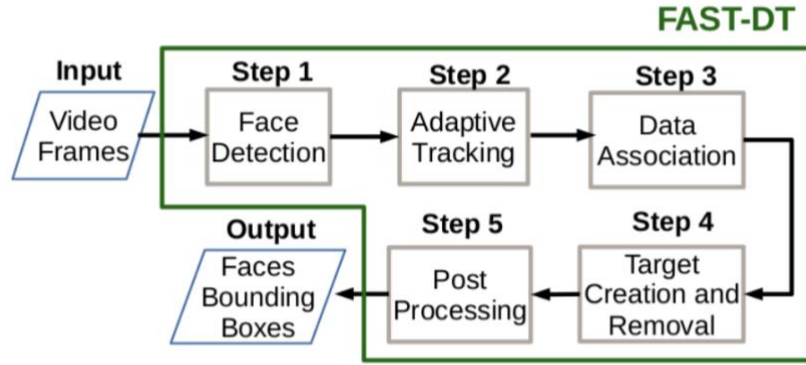149 boxes.

Figure 1. Multi-face detection and tracking framework [33]

Similar to the above work, [35] also adopt tracking-by-detection mechanism but with a more sophisticated tracker update routine. [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, the authors adopt an advancement of [34] and a color-assisted tracker as detect and track components respectively (Figure 2). The novelty of this work lies in the combined framework. Instead of running a detector for every frame like previous work, the authors 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.
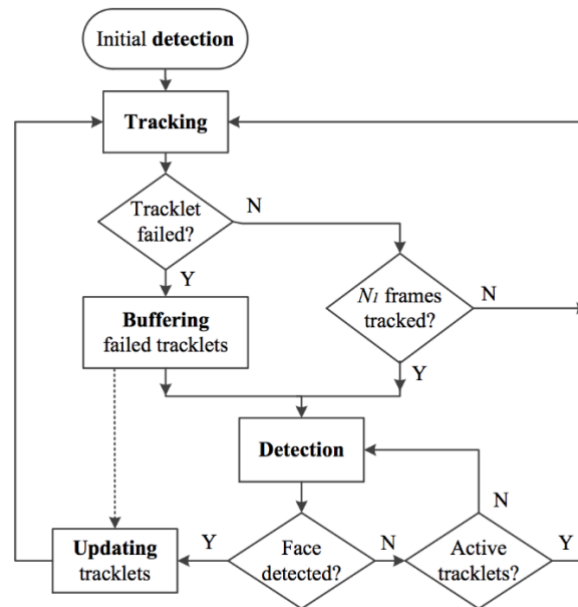


Figure 2. Multi-face tracking detection and tracking flow[36]

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

*Deep features*

168 Recently, many works [38]–[42] integrate deep feature extractors into the tracking
169 framework. Of those works, [38] adopts the sparse detection mechanism as described above
170 and use KLT tracker [43] for the tracking-by-detection stage. In the data association step
171 between detection boxes and track boxes, deep feature vectors are used as visual information
172 in addition to spatial information.
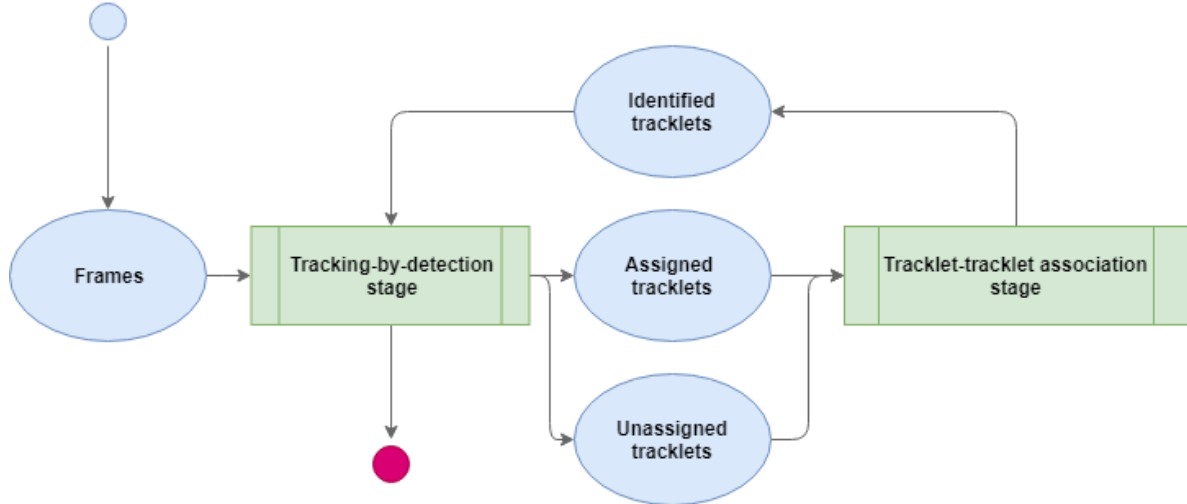
173 *Our approach overviews*



174

175 Figure 3. Our proposed method. The extra tracklet-tracklet association is introduced to improve accuracy by
176 using more information and lighten the process before.

177 **Semi-online tracking**

178 Aiming for practical usage and from the analysis of the online detection-tracking approaches,
179 we propose a new approach in semi-online manner by introducing the tracklet-tracklet
180 association stage (Figure 3).

181 After getting the detections of a frame, we should match it with tracklets up until the previous
182 frame to determine the identification of new detections. To achieve this criterion, using a
183 deep feature extractor is a heavy waste. We propose a way to lighten the process while
184 keeping the accuracy as high as possible. First, we use a light feature LBPH (Local Binary
185 Pattern Histogram) extractor in **the detection-tracking stage** (Figure 5) for efficient
186 computation and combine it with information from a tracking method (Kalman filter) to
187 reduce the errors as much as possible in creating short tracklets (we have not yet assigned
188 identifications for those tracklets). Then we observe that consecutive face boxes of one
189 tracklet are nearly the same, thus in the **tracklet-tracklet association stage** (Figure 8) **,** we
190 introduce a compression method to get representatives of a tracklet and apply a deep feature
191 extractor on these representatives instead of all boxes. We then link short tracklets into long
192 tracklets by using those features as appearance information. In the linking step, we also
193 introduce a new method for motion similarity between two tracklets. The tracklet-tracklet
194 association stage resolved much problems stated above: the future information of frames
195 sequences is well manipulated; the computational complexity is cut off from deep feature
196 comparison by applying the new compression method.

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

6

198 **Detection-Tracking stage**: The main role of this stage is to extract the track information of
199 targets in a frame using detecting and tracking methods. Technically, the detection-tracking
200 stage processes frame-by-frame for every mini-batch interval (i.e. 60 frames) and yield a list
201 of tracklets. The process is illustrated in Figure 4.

202 **Tracklet-tracklet association stage**: At the end of each mini-batch process, the list of
203 tracklets is passed to this stage. The main role of this stage is to correct false positives of the
204 previous stage and connect related tracklet to create long tracklets and then assign
205 identifications to these new tracklets. The process is showed in Figure 7.

206 *Computational complexity*

207 Our framework can process video streaming in real-time. The speed can reach around 60fps,
208 which is greater or equal the frequency of common videos (from 30 to 60fps).

209 *Detection-Tracking stage*

210 We leverage known detection-tracking approaches with some modification to speed up the
211 stage without sacrificing much performance and introduce a new stage to improve the
212 performance. We also implemented a framework: the detection-tracking stage combining
213 S3FD face detector to produce detection boxes, LHBPs feature extractor to extract the global
214 features,

215 Kalman Filter tracker to produce tracking boxes, then Hungarian algorithms for matching the
216 corresponding boxes to create tracklets.

217 After getting the detection information of a frame, we should match it with detection
218 information of the previous frame to determine if they are the same identification. To achieve
219 this criterion, using a deep feature extractor is a heavy waste. We proposed a way to lighten
220 the process while keeping the accuracy as high as possible. First, we use a light feature
221 extractor for efficient computation and combine it with information from another tracking
222 method like Kalman to reduce the errors as much as possible. Then we can extract
223 information from a mini-batch of frames to correct the process later while gaining more
224 useful information, trade of some delays.

225 *Tracklet-Tracklet association stage*

226 The tracklet - tracklet association stage uses the motion information simulated by the spline
227 interpolation and appearance information from FaceNet deep feature extractor to drop the
228 false positives and match the suitable tracklets to accurately assign the ids for targets.
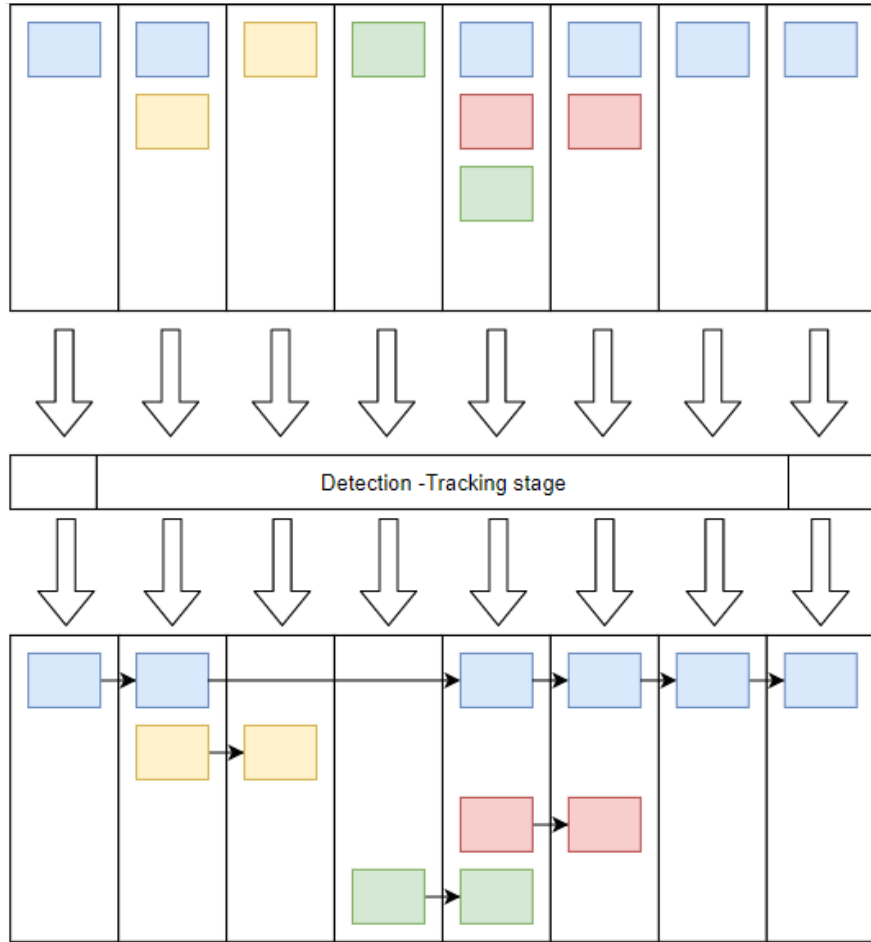
229 *Detection – tracking stage*

230

Figure 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.

**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.

**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.

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

**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

250 positions from previous frames. Our detection-tracking algorithm will try to fuse these
251 detection results with track results in order to better enhance the output, create more reliable
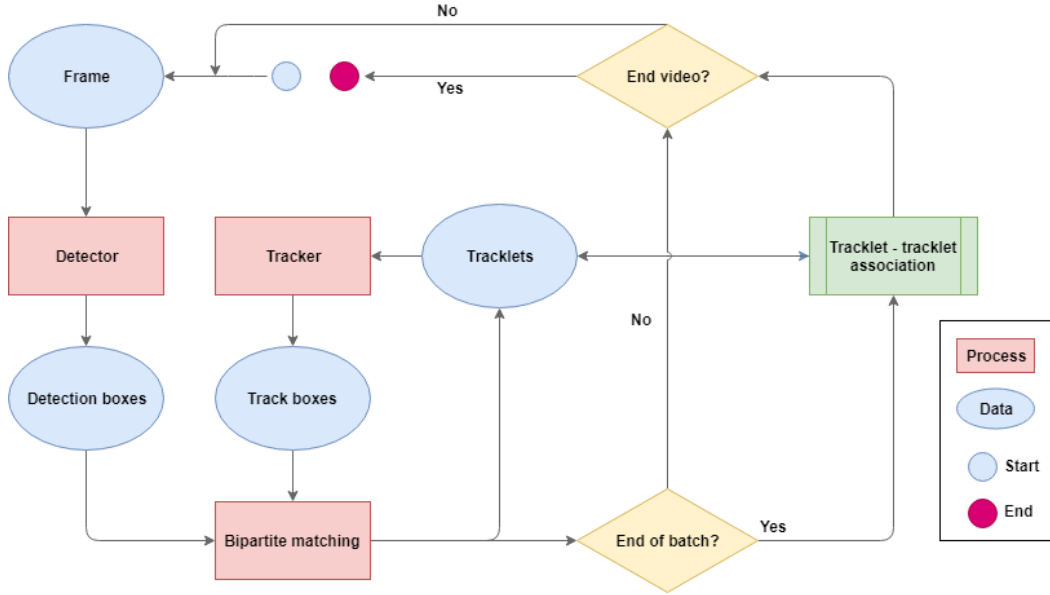252 tracklets



253

Figure 5. Our detection – tracking flow diagram.

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

262 Two features that are commonly leveraged are motion and appearance:

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

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

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

269 We define a gating unit for each affinity in order to filter out less likely matches. Because of
270 our intention that if a detection box and a track box are considered a possible match, they
271 must satisfy motion affinity alone and appearance affinity alone first.

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

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

9

275  where,

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

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

280  $\gamma_M$ is the threshold for space similarity distance determined by heuristic (we reason that
281  detection box and track box should be near to be of one individual, so we set this value to
282  0.3).

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

286  $Match(i,j)$ will be used to determine if a detection box and a track box is a possible match.
287  It only has value if both motion and appearance metrics are over their thresholds. Otherwise,
288  its value is zero, its range is from 0 to 1. The thresholds for $Match(i,j)$ are determined
289  through experiments (value search).

290  **Algorithm Tracklet-Tracklet Association**

291  We present the algorithm as pseudo-code in Figure 6.

292  There are some variables we must notice:

293  - *F(i)*: the i-th frame.

294  - *affinity_matrix*: an affinity matrix between detection boxes and track boxes.

295  - *not_match_threshold*: number of unassigned times for a tracklet to be considered
296    inactive.

297  - *motion_threshold*: minimum IoU to consider a pair of detection box and track box a
298    possible match, stated as $\gamma^M$ in equation (1).

299  - *appearance_threshold*: minimum cosine affinity threshold to consider a pair of
300    detection box and track box a possible match, stated as $\gamma^A$ in equation (1).

301  - *affinity_threshold*: minimum threshold to consider a pair of detection box and track
302    box a match.

---

**Algorithm 1:** Detection-tracking tracklet assignment algorithm

---

**Function** Detection-tracking tracklet assignment

**Input**:

**Ta**: a list of active tracklets (a tracklet is considered active if temporal difference between
the current frame and the last frame it has matched a detection box is less than
*not_match_threshold*)

---

**Ti**: a list of inactive tracklets (a tracklet is considered inactive if temporal difference between the current frame and the last frame it has matched a detection box is greater than or equal *not_match_threshold*)

**Output**:

**Ta**: a list of active tracklets (updated)

**Ti**: a list of inactive tracklets (updated)

1. **Begin**

2. **for** frame **in** frames_of_batch **do**

3.      *// Run detector on frame F(i)*

4.      detections = detect_func(frame)

5.      track_predictions = []

6.      just_stopped_boxes = []

7.

8.      *// Run tracker for each active tracklet on frame F(i)*

9.      **for** tracklet **in** Ta **do**

10.        *// run tracker on frame i-th to find track box*

11.        found, prediction = track_func(tracklet, F(i))

12.        *// if found prediction, add to a list to match with detections*

13.        **if** found:

14.          track_predictions.append(prediction)

15.          *// if not found, lower its priority in association step*

16.        **else**:

17.          just_stopped_boxes.append(tracklet.get_box_at_previous_frame)

18.        **end if**

19.        Ta.remove(tracklet)

20.      **end** **//**for tracklet

21.      *// Construct distance matrix for association*

22.      affinity_matrix = matrix(len(track_predictions), len(detections))

23.      **for** prediction **in** track_predictions **do**

24.        **for** detection **in** detections **do**

**25.**          affinity_matrix[prediction][detection] = get_affinity(detection, prediction)

**26.**     **end** **//**for detection

**27.**   **end** **//**for prediction

**28.**   *// Hungarian function return:*

**29.**   *// a list of unassigned tracklets*

**30.**   *// a list of unassigned detections*

**31.**   *// a list of pairs of detection-tracklet*

**32.**   unassigned_tracklets, unassigned_detections, matches = Hungarian(affinity_matrix)

**33.**   **for** detection, tracklet **in** matches **do**

**34.**      tracklet.update(detection)

**35.**   **end** **//**for detection, tracklet

**36.**   **for** tracklet **in** unassigned_tracklets **do**

**37.**      tracklet.non_detection_streak += 1

**38.**      *// if tracklet has not been assigned for any detection in streak_threshold times*

**39.**      *// assign tracklet to stopped state*

**40.**      **if** tracklet.non_detection_streak > streak_threshold:

**41.**         tracklet.state = inactive

**42.**         Ta.remove(tracklet)

**43.**         Ti.append(tracklet)

**44.**      **endif**

**45.**   **end** **//for** tracklet

**46.**   affinity_matrix = matrix(len(just_stopped_tracklets), len(unassigned_detections))

**47.**   **for** prediction **in** just_stopped_boxes **do**

**48.**      **for** detection **in** unassigned_detections **do**

**49.**         affinity_matrix[prediction][detection] = get_affinity(detection, prediction)

**50.**      **end** **//for** prediction

**51.**   **end** **//for** detection

**52.**   unassigned_tracklets, unassigned_detections, matches = Hungarian(affinity_matrix)

**53.**   **for** tracklet **in** unassigned_tracklets **do**

**54.**      Ta.remove(tracklet)

**55.**      Ti.append(tracklet)

**56.**    **end //for** tracklet

**57.**    *// initialize new tracklet by unassigned detections*

**58.**    **for** detection **in** unassigned_detections **do**

**59.**      Ta.append(init_tracklet(detection))

**60.**    **end //for** detection

**61. end //for** frame

**62. End //Function** Detection-tracking tracklet assignment

303                     Figure 6. Detection – tracklet assignment

304 Tracklet stores an internal tracker, a list of boxes (these boxes can be detection boxes or track

305 boxes) capturing one individual across frames and other information used in the algorithm

306 (state, non_detection_streak, etc.). The tracklet assignment process can be summarized as

307 follow:

308      Step 1.      For frame i-th, we run a face detector on that frame to get detection boxes

309              (referred in the pseudo-code as detections).

310      Step 2.      For each tracklet that is active, we run an internal tracker of that tracklet on

311              frame i-th to get a track box, if the tracker of the tracklet die (i.e. cannot find a box) in

312              this frame, we put the last box (either detection or track) of that tracklet in

313              *just_stopped_boxes*, otherwise we put the new track box in *track_predictions*.

314      Step 3.      From the above lists of detection boxes and tracking boxes, we construct an

315              affinity matrix between track predictions and detections. We then apply Hungarian

316              algorithm to assign detections to *track_predictions* based on this *affinity_matrix*. The

317              result of this algorithm is a list of pairs of detection boxes and tracklets being

318              assigned, a list of unassigned detection boxes, a list of unassigned tracklets.

319      Step 4.      Based on the results of Hungarian algorithm, we update matched tracklets and

320              unmatched tracklets.

321      Step 5.      Repeat step 3 and step 4 for *just_stopped_boxes* and *unassigned_detections*.

322      Step 6.      Repeat the above steps for the whole batch (i. e. 60 frames).

323 We then use the output of this algorithm for the tracklet-tracklet association stage. All the

324 tracklets taken from this stage is regarded as unknown id tracklets in the next stage.

325

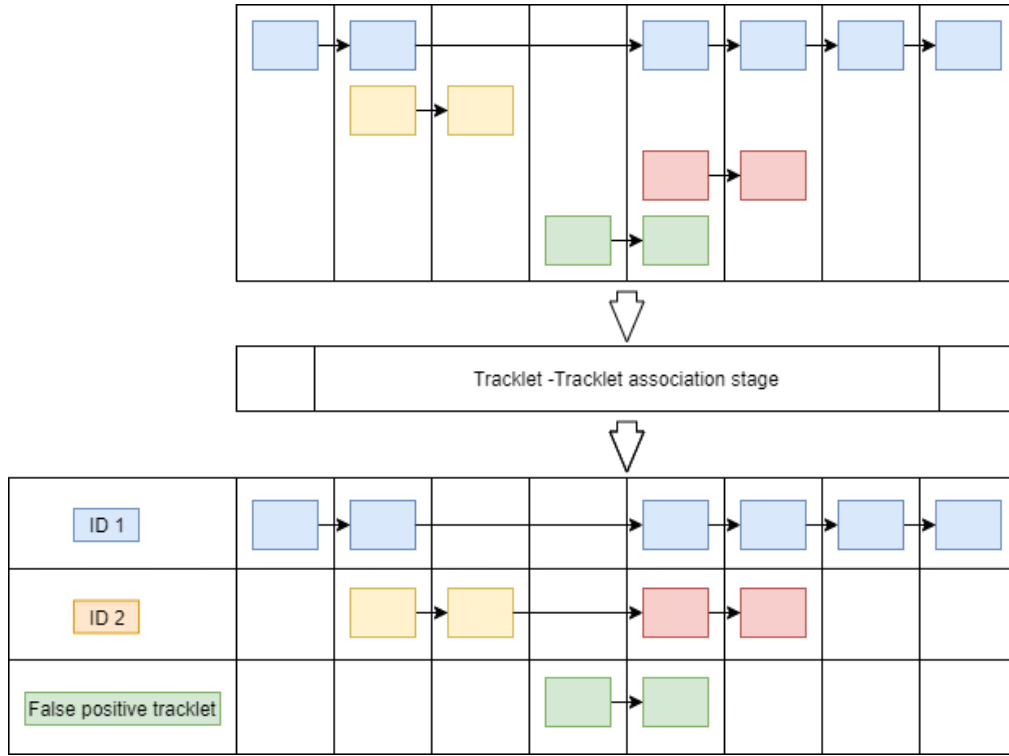326 *Tracklet-tracklet association*



327

Figure 7. Tracklet-tracklet association stage. From tracklets formed before, the identities will be determined in
this stage.

**Goal**

Short tracklets from the detection-tracking stage are passed to this stage. We 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.

**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 1.
We 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.

**Method**

After each batch processing the detection-tracking stage, we have a list of unknown-id
tracklets that are needed to be assigned identifications in this stage. We also have a list of
known-id tracklets in the past (previous batches). Our job is now trying to assign
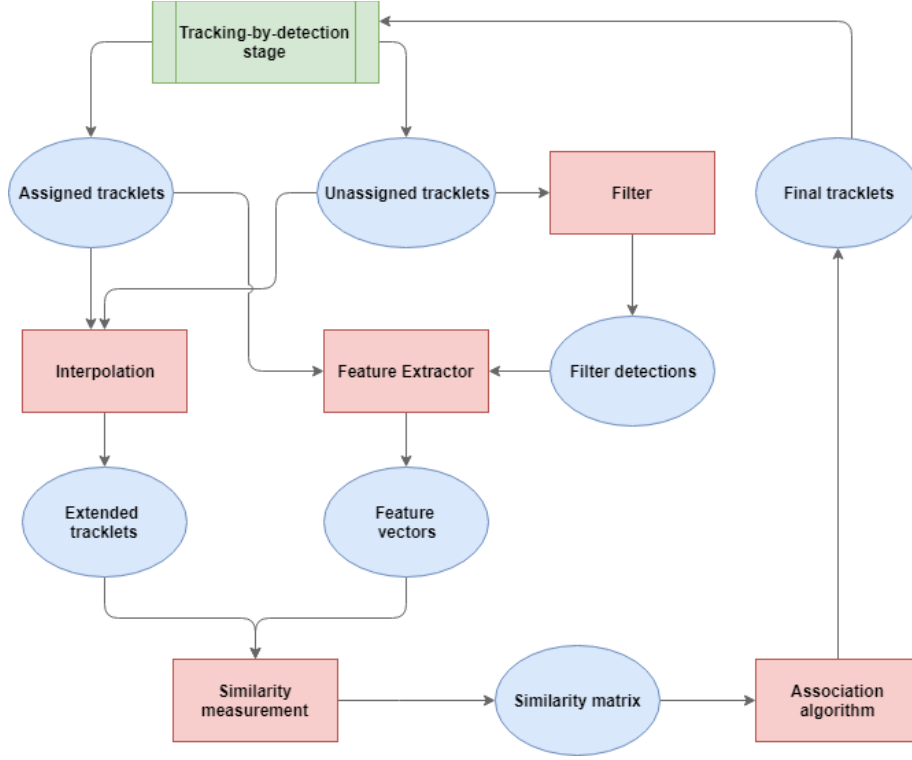identifications to unknown-id tracklets.

14

Figure 8. Our tracklet-tracklet association flow diagram.

We 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. We 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. We 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, we adopt the complementary nature of motion and appearance.

**Motion distance**

For motion, we 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), we extrapolate forward to get the possible boxes in the future relative to t(i). From the boxes of t(j), we extrapolate backward to get the possible boxes in the past relative to (t(j). For extrapolation, we assume that face movement can be modeled as a polynomial function and apply spline extrapolation. We 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, we have a pair of overlapped extrapolated boxes in the same frame f(k). We now calculate a spatial distance between two boxes using two centers and a diagonal distance between two boxes according to their diagonals. We 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:

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

In that,

15

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

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

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

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

379     Then the trajectory distance is the average of pair distances:

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

380     where,

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

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

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

384     **Appearance distance**

385     For appearance, we use average Euclidean distance between two feature sets of two tracklets.
386     For each box of a tracklet, we have a respective LBPHs feature (referred to as light feature)
387     extracted from the detection-tracking stage. Assume t(i) have N light feature vectors and t(j)
388     have M light feature vectors, one straight forward method is to compute N*M Euclidean
389     distances and use the average as the distance between two tracklets. From experiments, we
390     observe that LBPHs feature is not representative enough for this task. Thus, we adopt a deep
391     feature extractor [20] for this task.

392     However, deep feature extractors are computationally expensive and if we compute deep
393     features for all boxes of a tracklet the framework would not run in real-time. Moreover,
394     temporally adjacent boxes often contain similar information, so it would be redundant to
395     compute all the deep features. We introduce a mechanism to lower the number of boxes
396     needed to be passed through a deep feature extractor using already computed light features.

397     Given a list of light feature vectors of a tracklet, we apply a clustering algorithm on these
398     light feature vectors and pick out centroids, i.e. N_compressed boxes, for deep feature extraction.
399     This way we save a lot of time computing deep features while keeping the diversity of a
400     tracklet. We then use average Euclidean distance between two deep feature sets of two
401     tracklets as tracklet - tracklet appearance distance:

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

402     In that,

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

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

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

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

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

409     **Fusing results**

410     A weighted sum of appearance and motion affinities is the affinity between two tracklets
411     (used as the weight of the edge between two nodes). We fuse two affinities by taking the
412     addition rather than multiplication as used in the detection-tracking stage because motion
413     affinity is not reliable enough in case of long-term occlusion or camera shake. Thus, we set
414     the weight for motion affinity low so that it plays as extra information.

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

415     Where

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

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

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

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

422     **Algorithm Tracklet-Tracklet Association**

423     We present the algorithm as pseudo-code, Figure 9 presents preprocess, Figure 11 presents
424     utility functions and Figure 10 presents tracklet-tracklet association.

425

**Preprocess before the tracklet-tracklet association**

**Function Preprocess before the Tracklet-Tracklet Association**

**Input**:

C: a list of continuous tracklets

H: a list of head-interrupted tracklets

T: a list of tail-interrupted tracklets

**Output**:

U: a list of unknown-id tracklets

1. **for** th **in** H **do**

2.     H.remove(th)

3.     **if** its tail-counterparted has an id:

4.       assign that id for th

5.     **else**:

6.       tracklet = merge(th, tail(th))

7.       H.append(tracklet)

8.     **end if**

9. **end** //for th

10.

11. **for** tt **in** T **do**

12.     **if** is_face_test(tt) and num_detections(tt) < 5:

13.       T.remove(tt)

14.       reserve tt for next mini-batch

15.     **end if**

16. **end** //for tt

17.

18. U = H + T + C

19. **for** u **in** U **do**

20.     **if** not is_face_test(u):

21.       U.remove(u)

22.     **endif**

23. **end** //for u

Figure 9. Preprocess before the tracklet-tracklet association

After each batch processing (the detection-tracking stage), we have a list of unknown-id tracklets that are needed to be assigned identifications. There are two types of tracklet after the batch processing: interrupted tracklets (ranging across batch intervals) and continuous tracklets (ranging within a batch interval). Among interrupted tracklets, we have head-interrupted and tail-interrupted ones. For example, assume our batch size is 64 and there is a tracklet starting from frame 55-th to frame 77-th. This batch approach will divide the original tracklet into two tracklets: the first part (tail-interrupted) starts from frame 55-th to 64-th and the second part (head-interrupted) starts from frame 65-th to 75-th. For each head-interrupted tracklet, if its respective tail-interrupted part (have the same track_id) has been assigned an identification, we merge it to the tail-interrupted part and treat the new tracklet as a known-id

18

437 tracklet, otherwise, we merge both into one tracklet and treat it as an unknown-id tracklet
438 waiting to be assigned in this stage. For each tail-interrupted tracklet, if it passes the
439 *is_face_test* (at least half of the number of boxes are from detection boxes) and the number of
440 detection boxes is smaller than 5, we delay the identification assignment for this tracklet and
441 wait for the information from next batch (wait for its head-interrupted part).

442 Before passing unknown-id tracklets to the assignment process, we filter out false positive
443 tracklets by the *is_face_test* function.

**Algorithm 2: Tracklet-tracklet association**

---

**Function Tracklet-tracklet association**

**Input**:

U: a list of unknown-id tracklets

A: a list of known-id tracklets

**Output**:

A: a list of known-id tracklets (updated)

1. **Begin**

2. *// compress tracklet for speed efficiency in extracting features*

3. **for** tracklet **in** U **do**

4.     ids = compress(tracklet)

5.     tracklet.features = deep_feature_extractor(tracklet, ids)

6. **end** //for tracklet

7.

8. *// Construct the graph*

9. graph = construct_graph(U, A)

10.

11. *// Leiden algorithm returns a list of groups*

12. groups = Leiden(graph)

13.

14. *// assign id for matched tracklet*

15. **for** group in groups **do**

16.     update_id_tracklet(group)

17. **end** //for group

---

**18. End** //Function Tracklet-tracklet association

Figure 10. Tracklet - Tracklet assignment

444

445 After preprocessing, we now have a list of unknown-id tracklets. We also have a list of
446 known-id tracklets. The two lists are input for the tracklet-tracklet association algorithm.

447     • First, we extract deep features for unknown-id tracklets as detailed above.

448     • We construct a graph with nodes being tracklets and edges being respective affinity
449       (computation process is detailed above)

450     • Apply Leiden algorithm on this graph, the output of this algorithm contains groups of
451       tracklets.

452     • If a group contains a known-id tracklet, we set that identification for others tracklet,
453       otherwise, we create a new identity for that group.

454 Following is some utility functions for calculating the weight for the graph.

**Utility procedures used in Algorithm 2**

1. lambda: weight between motion and visual

2.

3. **procedure** is_face_test(tracklet):

4. **Begin**

5.    nd = num_detections(tracklet)

6.    nt = num_tracks(tracklet)

7.    **if** nd > nt:

8.       **return** True

9.    **end if**

10.    **return** False

11. **End**

12.

13. **procedure** compress(tracklet):

14. **Begin**

15.    light_features = tracklet.light_features

16.    clusters = Agglomerative_clustering(light_features, distance=L2)

17.    ids_choosen = []

18.    **for** cluster **in** clusters **do**

```
19.        ids_choosen.append(max_detection_score(cluster))
20.    end //for cluster
21.    return ids_choosen
22. End
23.
24. procedure motion_affinity(t1, t2):
25. Begin
26.    et1 = forward_extrapolate(t1)
27.    et2 = backward_extrapolate(t2)
28.    distances = []
29.    get_ratio_and_spatial_distance(et1, et2)
30.    return mean(distances)
31. End
32.
33. procedure visual_affinity(t1, t2):
34. Begin
35.    distances = []
36.    for f1 in t1.features do
37.      for f2 in t2.features do
38.        distances.append(Euclidean_distance(f1, f2))
39.      end //for f2
40.    end //for f1
41.    return mean(distances)
42. End
43. procedure affinity(t1, t2):
44. Begin
45.    motion = motion_affinity(t1, t2)
46.    visual = visual_affinity(t1, t2)
47.    return visual * lambda + (1 - lambda) * motion
48. End
```

455             Figure 11. Utility functions used in algorithm 2

456 **Contributions**

457 This proposed approach tackles challenges related to online approach above:

458       •   Instead of computing deep features for all faces of one tracklet as online approaches
459          do, we leverage light features (LBPHs) in the context of tracklet to efficiently
460          compute deep features (extracted by deep network) without compromising
461          representative power. In fact, the compressing method produces a more accurate
462          representation for a tracklet thanks to diversity and high detection quality (high-score
463          detected boxes).

464       •   Using this framework, we can tighten the constraints in the tracking-by-detection
465          stage so that the possibility of wrongly matching is low. Though having many
466          tracklets after the tracking-by-detection stage, these tracklets will be grouped in the
467          tracklet-tracklet association stage.

468       •   We do not have to assign identifications to new detections right away in the detection-
469          tracking stage but leave it to the tracklet-tracklet association stage. This way we can
470          filter out false positives efficiently in the pre-processing step.

471       •   The identification assignment step is tracklet-based; thus, we can take advantage of
472          temporal information of tracklets (co-extant tracklets belong to different individuals)

473       •   We also propose the trajectory difference metric to account for motion in tracklet-
474          tracklet comparison.

475 In application, we often have limited data so using a pre-trained model and finetuning on our
476 data is a reasonable choice. In this work, we show that simply adopting deep features
477 (extracted by Facenet) and employ Euclidean (or cosine) metric is not discriminative enough
478 in reference to real-life data. Therefore, we propose to apply Logistic discriminant metric
479 learning so that the new embedding space for real-life data is more discriminative.

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

483 # Results and Discussion

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

491 **The purpose of experiments on MobiFace and OTB datasets**

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

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

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

500 *Experiments on MobiFace dataset*

501 The MobiFace dataset

502 MobiFace dataset [44], as mentioned from the original paper, is the first dataset for single
503 face tracking in mobile situations. Due to the lack of engrossing face tracking datasets before
504 MobiFace, the performance of pioneer face trackers was reported on a few videos or on small
505 subsets of the OTB dataset, and the comparison between approaches was limited. The
506 introduced dataset provides a unified benchmark with different attributes for future
507 development in this field. Some samples of the dataset are illustrated in Figure 12.

508 The authors collected 80 unedited live-streaming mobile videos captured by 70 different
509 smartphone users in fully unconstrained environments and manually labeled over 95.000
510 bounding boxes on all frames. In order to cover typical usage of mobile device camera, the
511 authors fetched videos from YouTube mobile live-streaming channels. Most of the videos are
512 captured and uploaded under fully unconstrained environments without any extra video
513 editing or visual effects. 6021 videos were collected and discarded under strict criteria that
514 the target faces should appear at least in 10% of the video frames, and the target faces should
515 not always stay still to serve the purpose of visual tracking. Besides the common 8 attributes
516 in object tracking datasets, the authors proposed 6 additional attributes commonly seen in
517 mobile situations.

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

521 *Setup the experiments*

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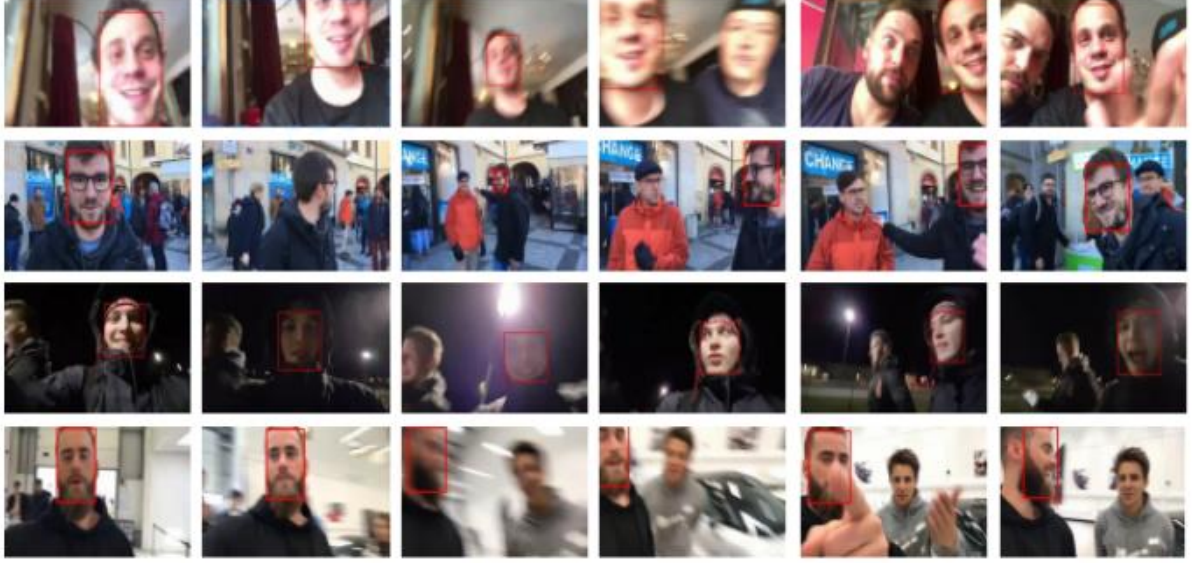528 results of that target only.

529

Figure 12. Some example frame from the MobiFace dataset [44]. Red ground truth bounding boxes are annotated by the authors.

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

534 We consider the three metrics proposed in the dataset. As most of the metrics are in plot
535 form, we will explain the way to extract an important metric from the plot, the area under the
536 curve (AUC). With N is the number of thresholds used to draw the plot, we have $n =$
537 $1, 2, 3, \ldots, N$. The curve was drawn from points with coordinate $(t_n, f_n)$, $t_n$ is the threshold
538 value at that point and $f_n$ is the evaluated value of our algorithm at that threshold, i.e. location
539 error of precision plot, overlap score of success plot. The AUC is then calculated by

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

(6)

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

547 *Success plot*: Overlap score is also another commonly used metric in the tracking field. Given
548 a ground truth bounding box $r_{gt}$ of the target, the predicted bounding box of our algorithm is
549 $r_p$. Then we can compute the overlap score by the intersection over union (IoU) of those two
550 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
551 rectangles, respectively. The success plot reflects the percentage of frames in which the
552 intersection over union (IoU) of the predicted and ground truth bounding box is greater than a
553 given threshold. Usually, the average success rate at 0.5 threshold is enough for evaluation.
554 In addition, the area under the curve (AUC), which is the accumulated success rate can also

555 be used for measurement. We can use those metrics interchangeably to summarize the
556 performance.

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

564 *Experiment results*

565 Evaluation metrics of our method and state-of-the-art methods are illustrated in Figure 13,
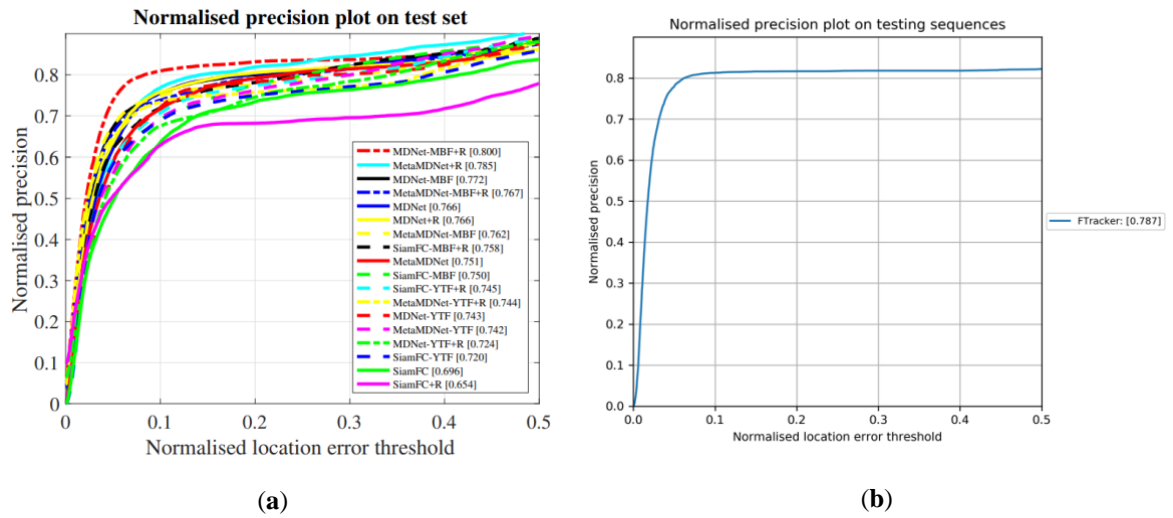566 Figure 14 and a detailed comparison is showed in Table 1:



(**a**)   (**b**)

567 Figure 13. Evaluation results of trackers on MobiFace test set: (**a**) results from MobiFace paper [44], (**b**) results
568 on our method
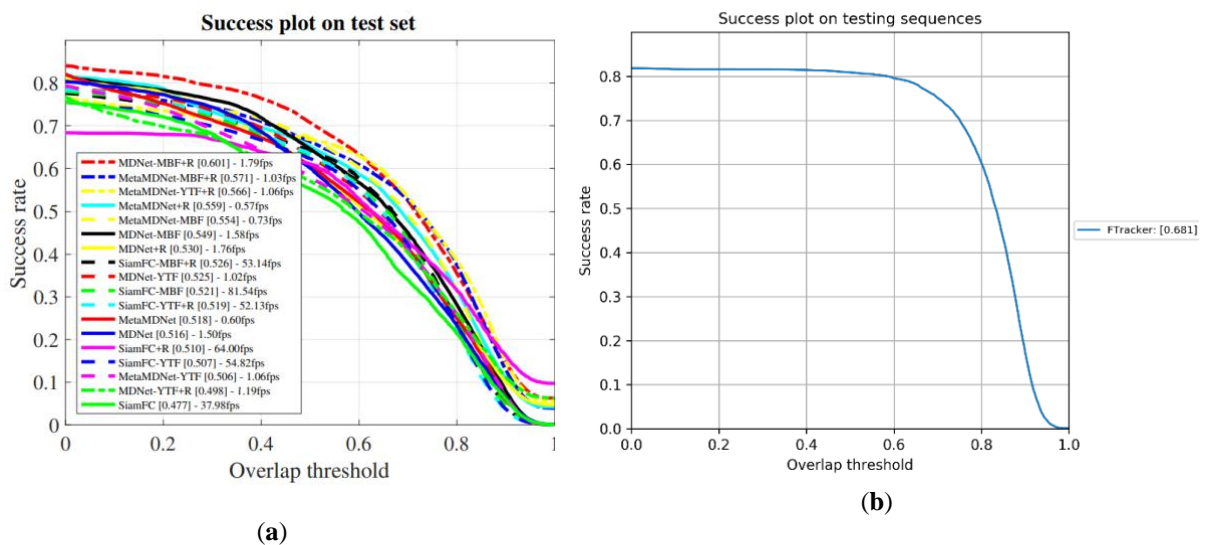


(**b**)

(**a**)

569 Figure 14. Evaluation results of trackers on MobiFace test set: (**a**) results from MobiFace paper [44], (**b**) results
570 on our method

571 Table 1. A detailed comparison between our method and MobiFace evaluated results

| Tracker | Normalised Precision plot (AUC) | Success plot (AUC) | FPS |
|---|---|---|---|
| MDNet-MBF+R | **0.800** | **0.601** | 1.79 |
| MetaMDNet-MBF+R | 0.767 | **0.571** | 1.03 |
| MetaMDNet-YTF+R | 0.744 | 0.566 | 1.06 |
| MDNet-MBF | **0.772** | 0.549 | 1.58 |
| SiamFC-MBF+R | 0.758 | 0.526 | **53.14** |
| SiamFC-MBF | 0.750 | 0.521 | **81.54** |
| Our framework | **0.787** | **0.681** | 44.38[1] |

572 *Discussion*

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

579 As shown in the above plot, our method has an advantage in the success plot, but not the
580 precision plot. The precision plot affected by the Euclidean distance of centers of ground
581 truth bounding boxes and our tracking boxes. When the predicted box is drifted from the
582 face, we terminate the tracklet instantly; therefore, with high normalised error, our tracker
583 performs the same as with low normalised error while other trackers yield noticeably
584 different results with different normalised errors.

585 The success plot may have more practical usages because the IoU decide how we can make
586 use of the information we extracted. The success plots of trackers evaluated in MobiFace
587 dataset are started from very high, but the slope is very steep. Starting from above 0.8 success
588 rate for threshold 0, to threshold 0.5, they drop to below 0.7 success rate. The steep slope
589 indicates predicted boxes of those trackers are not always aligned with ground truth boxes.
590 Our starting point is somewhere below 0.8 success rate but maintains the success rate over
591 the overlap threshold change. At threshold 0.5, our approach still has a high success rate,
592 above 0.7, indicating our boxes is closely aligned with ground truth boxes. At 0.5 threshold,
593 the predicted boxes cover most of the track target and can be well used in application.

---

[1] We 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).

594 Besides, as the main target of ours is for practical usages, a good success plot and success
595 rate at 0.5 threshold - while keeping the speed - are acceptable.

596 *Experiments on OTB (Object Tracking Benchmark) Dataset*

597 About the dataset

598 OTB Dataset [45] is one of the most famous datasets specifically used for benchmarking the
599 object trackers since its appearance. The authors worked to collect and annotate most of the
600 common tracking sequences from different datasets. They also classified those sequences into
601 multiple categories by challenges as in Table 2 and selected 50 difficult and representative
602 ones in the TB-50 dataset for an in-depth analysis. The full dataset contains more sequences
603 of human (36 body and 26 face/head videos) than other categories because human target
604 objects have the most practical usages, some samples of the dataset is illustrated in Figure 15.

605

606 Table 2. Annotated Sequence Attributes with the Threshold Values in the Performance Evaluation from OTB
607 Dataset [45]

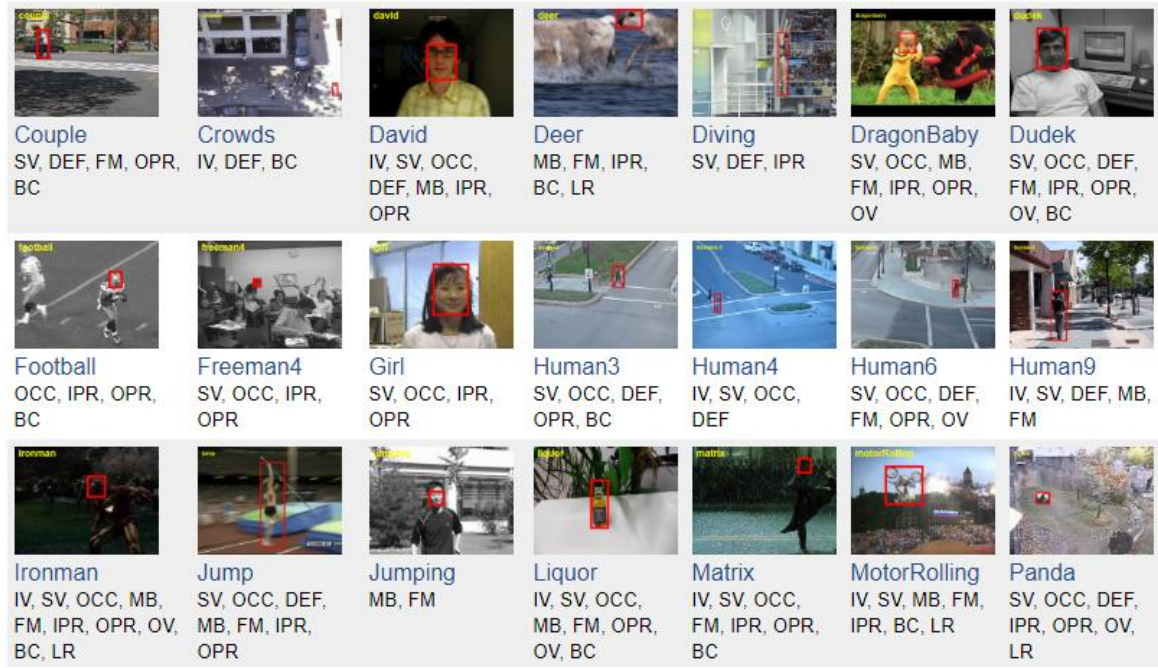| Attribute | Description |
| --- | --- |
| **IV** | Illumination Variation - The illumination in the target region is significantly changed |
| **SV** | Scale Variation - The ratio of the bounding boxes of the first frame and the current frame is out of range. $\left[\frac{1}{t_s}, t_s\right], t_s > 1 (t_s = 2)$ |
| **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 $t_r$ ($t_r = 400$) |

608

609 Figure 15. Some example sequences from the OTB Dataset [45]

610 Before the introduction of MobiFace dataset, face tracking methods can only be evaluated on
611 small self-collected datasets or a subset of OTB dataset. The whole dataset is designed for the
612 object tracking algorithms, but we selectively pick out the sequences with faces to conduct
613 experiments and compare with those methods mentioned before. The chosen face subset is
614 described in Table 3, the top 10 sequences are referred to as the difficult set and top 15 is the
615 normal set [46]:

616 Table 3. Chosen sequences and their attributes

| # | 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 |

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

620 *Set up the experiments*

621 Because the authors of MobiFace dataset inherit a lot of legacy from OTB dataset, in general,
622 the setup stage and evaluation stage for OTB Dataset are the same as the MobiFace dataset.

623 *Experimental results*

624 Evaluation metrics of our method and state-of-the-art methods are illustrated in Figure 16,
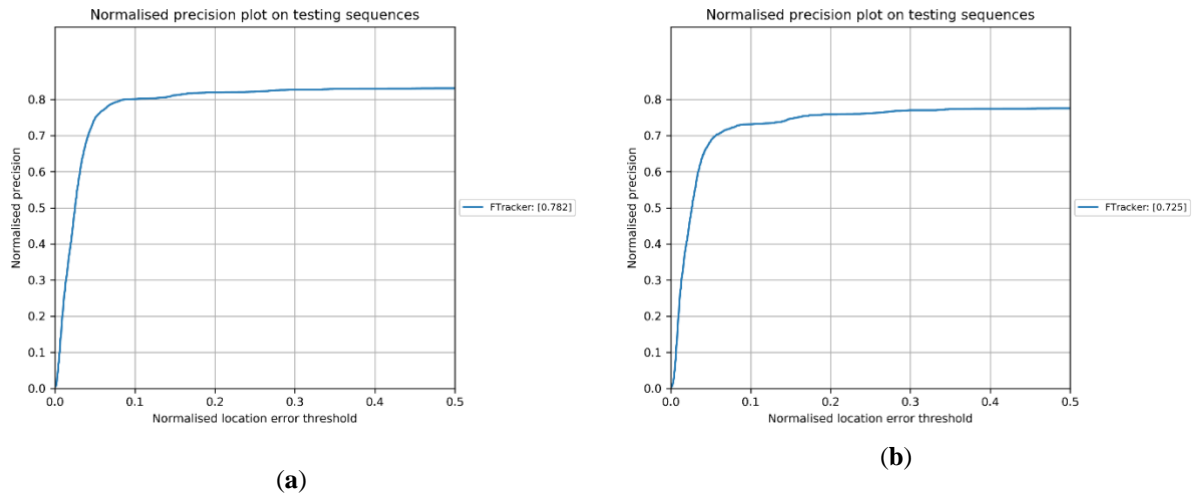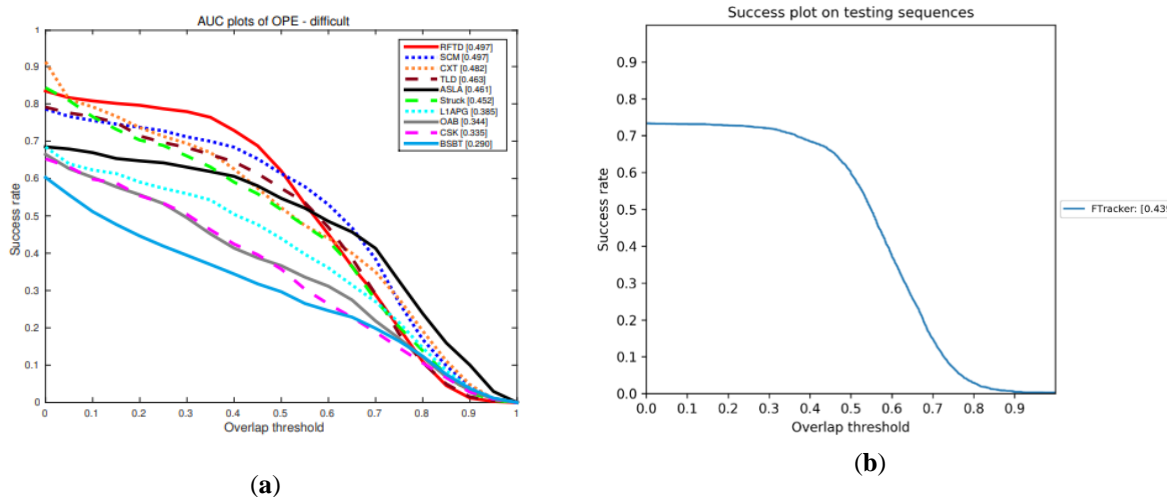625 Figure 17, and a detailed comparison is showed in Table 4 and Table 5.

626



Figure 16. Our normalised precision plot on OTB Dataset face subsets (a) normal set(b) difficult set

Figure 17. Success plots of trackers on OTB Dataset face subset (difficult set): **(a)** results from RFTD paper[46] **(b)** results on our method

Table 4. 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 | **71.3** |
| Struck | **55.9** | 67.6 |
| SCM | **58.3** | **72.6** |
| ASLA | 53.8 | 62.9 |
| CSK | 48.0 | 56.8 |
| L1APG | 50.7 | 59.7 |
| OAB | 42.6 | 48.9 |
| TLD | 51.8 | 67.3 |
| CXT | **57.3** | 65.7 |
| BSBT | 40.6 | 47.0 |
| **Our framework** | 51.9 | **68.3** |

Table 5. Top tracker comparison on OTB Dataset face subsets (difficult set). Evaluated results are from RFTD paper [46].

| Face Tracker | Success Plot AUC | Success plot Threshold (0.5) |
|---|---|---|
| RFTD | **49.7** | **62.0** |
| Struck | 45.2 | 51.7 |
| SCM | **49.7** | **61.3** |
| ASLA | 46.1 | 54.7 |
| CSK | 33.5 | 52.2 |
| L1APG | 38.5 | 43.9 |
| OAB | 34.4 | 36.6 |
| TLD | 46.3 | 57.4 |
| CXT | **48.2** | 52.2 |
| BSBT | 29.0 | 29.7 |
| **Our framework** | 43.9 | **59.7** |

634 *Discussion*

635 The precision plots in Figure 16 are good. The overall results are quite good, and the slope is
636 shallow as predicted after witnessing above experiments. However, we have no data from
637 other works to have an in-depth comparison.

638 As first sight from the metric Table 4 and Table 5, our framework has average AUC while
639 the slope of our framework is also shallow as predicted. The main reason here is because
640 when the predicted box is drifted from the face, we terminate the tracklet instantly; therefore,
641 with high normalised error, our tracker performs the same as with low normalised error while
642 other trackers yield noticeably different results with different normalised errors. The initial
643 modest success rate leads to a modest average value. The success rate at threshold 0.5 is still
644 good, ranking third in that section in both subsets.

## Conclusions

646 In this work, we proposed a method for face tracking problem in semi-online manner - the
647 online process with some minor delay. The comparing experiments are conducted on two
648 datasets: MobiFace dataset and OTB dataset with many state-of-the-arts works in the field.
649 The results show that our method can produce robust accuracy while keeping a good speed.
650 With that, the effectiveness of adding the tracklet-tracklet association stage after detection
651 stage in semi-online manner is proven. The manipulation of appearance affinity and motion
652 affinity have brought us the accuracy of the framework, while the workload division and
653 information sharing of the two main stages make our process lighter and achieve better speed.
654 With the improvements, all the disadvantages pointed out in section 1 are solved.

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

661 Future works following this work can dig in many ways for the better. First, try other
662 combinations of related techniques (detector, tracker, feature extractor) to achieve better
663 results. Second, exploit the concept of semi-online manner (use some delay for better results)
664 is also a good point. Third, the work only targets the continuous video stream, but in other
665 fields, many other attributes can be used to improve the results. Finally, a public multi-face
666 dataset will be a major contribution to this field.

## Data Availability

668 The OTB and MobiFace dataset supporting this study are from previously reported studies
669 and datasets, which have been cited. The processed data used to support the findings of this
670 study are available from the corresponding author upon request.

## Conflicts of Interest

The author(s) declare(s) that there is no conflict of interest regarding the publication of this paper.

## Funding Statement

## Acknowledgments

## Supplementary Materials

## 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. Rehg, "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. Rehg, "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 IEEECVF Conf. Comput. Vis. Pattern Recognit. Workshop CVPRW*, pp. 1519–151909, 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.

33

[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 Erdös-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.

777  [43]    Jianbo Shi and Tomasi, "Good features to track," in *1994 Proceedings of IEEE Conference on*
778  *Computer Vision and Pattern Recognition*, 1994, pp. 593–600.

779  [44]    Y. Lin, S. Cheng, J. Shen, and M. Pantic, "MobiFace: A Novel Dataset for Mobile Face Tracking in the
780  Wild," *ArXiv E-Prints*, p. arXiv:1805.09749, May 2018.

781  [45]    Y. Wu, J. Lim, and M.-H. Yang, "Object Tracking Benchmark," *IEEE Trans. Pattern Anal. Mach.*
782  *Intell.*, vol. 37, pp. 1–1, 2015.

783  [46]    F. Comaschi, S. Stuijk, T. Basten, and H. Corporaal, "Robust online face tracking-by-detection," in
784  *2016 IEEE International Conference on Multimedia and Expo (ICME)*, 2016, pp. 1–6.

785