# Exploring Local Context for Multi-target Tracking in Wide Area Aerial Surveillance

Bor-Jeng Chen and Gérard Medioni

University of Southern California
{borjengc, medioni}@usc.edu

**Abstract.** Tracking many vehicles in wide coverage aerial imagery is crucial for understanding events in a large field of view. Most approaches aim to associate detections from frame differencing into tracks. However, slow or stopped vehicles result in long-term missing detections and further cause tracking discontinuities. Relying merely on appearance clue to recover missing detections is difficult as targets are extremely small and in grayscale. In this paper, we address the limitations of detection association methods by coupling it with a local context tracker (LCT), which does not rely on motion detections. On one hand, our LCT learns neighboring spatial relation and tracks each target in consecutive frames using graph optimization. It takes the advantage of context constraints to avoid drifting to nearby targets. We generate hypotheses from sparse and dense flow efficiently to keep solutions tractable. On the other hand, we use detection association strategy to extract short tracks in batch processing. We explicitly handle merged detections by generating additional hypotheses from them. Our evaluation on wide area aerial imagery sequences shows significant improvement over state-of-the-art methods.

**Keywords:** multi-target tracking, context tracker, wide area motion imagery

## 1 Introduction

Wide area motion imagery (WAMI) is acquired by high altitude unmanned aerial vehicles (UAV) and has made it possible to understand activities in a large area of interest. With current sensor and storage technologies, WAMI is typically captured in large format (tens $\sim$ hundreds of megapixels), low frame rate (1 $\sim$ 2 Hz), grayscale, and with ground sampling distance from 0.2 $\sim$ 0.5 meter/pixel. Because of these unique characteristics and its wide applications, it has gained attention in the computer vision field in recent years [1,2,3,4,5,6,7]. Tracking many vehicles in WAMI [6,7,8,9,10] is an essential component since it is needed for higher level activity analysis and scene understanding.

Detection association strategy has become standard for multi-target tracking with unknown number of targets [12,13,14]. By assigning detections in each frame into tracks, short-term missing detections are recovered by using motion interpolation. In the following, we call this type of tracker detection-based tracker
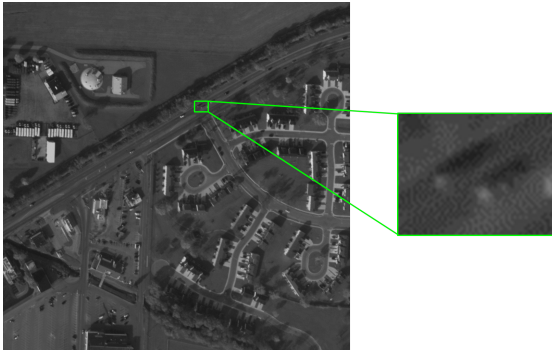
**Fig. 1.** A cropped view of WPFAB dataset [11]. The green rectangle shows details of two vehicles which are close to each other.

(DBT). Unique characteristics in WAMI bring additional challenges in DBT. First, low frame rate WAMI leads to a very large search space. Displacement of a target is up to 80 pixels in consecutive frames. Second, the size of a vehicle is extremely small (usually less than 20 pixels in length). Along with grayscale imagery, appearance information is less discriminating in WAMI than in many other scenarios. Fig. 1 shows a cropped view from a WAMI dataset. Notice that in this work, we assume that vehicles are the only moving targets, since other moving objects such as pedestrians are nearly invisible even for human beings.

Instead of using an appearance-based detector, current tracking approaches in WAMI [7,8,9,10] rely on motion detection. They first stabilize the imagery and then apply frame differencing methods. However, stopped or slow vehicles result in long-term missing detections, which can not be recovered using DBT only.

To alleviate the limitation of detectors, a new strategy that couples DBT with an appearance-based category free tracker (CFT) has been proposed in multi-target tracking [15,16]. CFT does not rely on a detector and it recovers missing detections at the ends of a track. However, appearance-based CFT is not robust in WAMI because targets are with weak appearance. Moreover, it is difficult to escape from local maximums when the frame rate is low. Prokaj and Medioni [6] proposed to run DBT and a regression based tracker in parallel. Though using a regressor is more efficient than using a classifier in low-frame-rate WAMI, the regression model is still not discriminating enough to avoid drifting.

Based on the above discussion, incorporating CFT with DBT in WAMI still remains very challenging. We conclude that there are three major issues that make it difficult in WAMI: 1) Discriminating ability of target appearance is very low because of the small target size and grayscale imagery. Illumination change can easily confuse a target appearance model with its local background or other neighboring targets. 2) Motion detections are imprecise especially at merged detections, which contain more than one target. Thus, it is difficult to have good
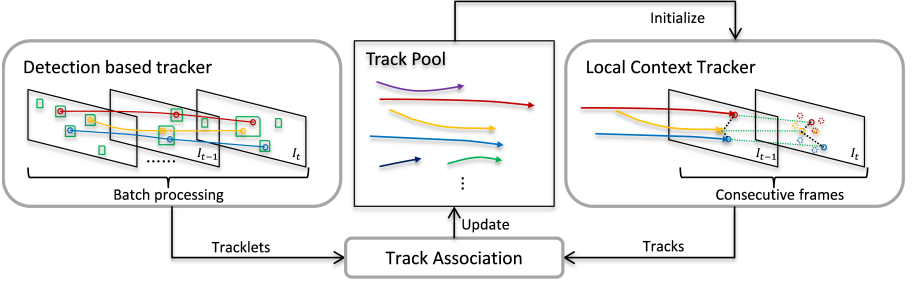
**Fig. 2.** Framework of our tracking system. Green rectangles represent motion detections. Solid circles show hypothesis centers that form tracks. Dashed circles represent hypothesis centers in LCT. Dashed black lines show context constraints in LCT.

initialization for CFT. 3) Displacement of a target between consecutive frames can be very large. A good strategy for hypothesis sampling in CFT is essential to keep the solution tractable.

Our goal is to maximize the merit of both trackers and compensate their limitations in WAMI. Unlike [15,16,6] training an appearance model for CFT, we relax the dependency on appearance clue, which is not reliable in our scenario, by introducing *local context tracker* (LCT). LCT explores spatial relations for a target to avoid unreasonable model deformation in the next frame. We design two sampling strategies based on dense and sparse optical flow to overcome large search space in low frame rate aerial videos. In DBT, short tracks (tracklets) are produced by associating hypotheses from motion detections in a sliding temporal window. We explicitly handle merged detections by generating additional hypotheses from them. This step is important for combing DBT with LCT to ensure reasonable appearance and motion consistency. The track association step concludes results from both trackers and updates the "track pool", which stores all existing tracks and is used to initialize LCT in the next frame. Fig. 2 illustrates the framework of our system.

The contributions of this paper are:
1. We propose LCT that relaxes the dependency on frame differencing motion detection and appearance information.
2. We propose DBT that explicitly handles merged detections in detection association.
3. We propose a unified framework that couples LCT with DBT and takes advantages of both trackers.
4. Our performance shows significant improvement over state-of-the-art methods in two WAMI sequences.

The rest of the paper is organized as follows: We discuss related work in Section 2. We then illustrate LCT in Section 3. Our DBT is introduced in Section 4. The track association module is described in Section 5. In Section 6, we

show comparison results on two sequences from real WAMI datasets. Finally, we conclude this paper in Section 7.

## 2    Related Work

Multi-target tracking has been investigated in the computer vision society for many years. Joint probabilistic data association filter (JPDAF)[17] and multiple hypothesis tracking (MHT)[18] are two early successful approaches. However, the association step requires very high computational and memory cost in both methods; therefore, solutions are usually intractable in real-world applications. In practice, JPDAF has been incorporated with Kalman filter [19] and particle filter [20] to increase the efficiency. More recently, Rezatofighi [21] improves the efficiency of JPDAF by obtaining m-best solution using integer linear programming and shows state-of-the-art result. MHT usually introduces tree pruning strategies [22,23,24] to reduce the possible solution space. In recent years, network flow optimization has become popular in multi-target tracking approaches [25,26,27,28]. Despite that these methods have shown promising results in their scenario, they all use one-to-one matching assumption. Therefore, they are not suitable for WAMI where split-and-merge motion detections often occur.

Solving tracking problem with machine learning techniques has shown to be effective in boosting discriminative ability of appearance model for single target tracking [29,30,31] and multi-target tracking [32,15,33]. However, it is almost impossible to learn meaningful information in WAMI because the target size is extremely small and the imagery is typically in grayscale. Targets are often visually similar to each other and background patches.

Using context information is an appealing strategy for tracking against distracters and occlusion. This concept has been applied on single target tracking [34,35,36]. Recently, [37] incorporates spatial constraints with tracking-by-detection. Nevertheless, these trackers require accurate annotation for initialization. It is not trivial to directly apply these approaches in WAMI, where the number of targets varies with time and perfect initialization is not available.

Most WAMI tracking approaches focus on associating noisy motion detections into tracks. Motion detections are typically acquired by applying frame differencing methods to stabilized imagery. Perera *et al.* [8] propose to first generate short tracks using nearest-neighbor strategy and then handle split-then-merge situations in track linking. Reilly *et al.* [9] formulate the data association problem in Hungarian algorithm. Prokaj *et al.* [38] extract tracklet from detections by Bayesian network. Shi *et al.* [39] associate motion detections by rank-1 tensor optimization. Keck *et al.* [40] provide a real-time implementation for tracking based on multiple hypothesis tracking. The object-centric association method is proposed in [10] to relax the one-to-one matching assumption for motion detections. Additional context constraints are used to alleviate track intersection. Chen and Medioni [7] extract tracklets by finding the longest path through detection trees. These above trackers mainly rely on motion detections. Therefore, they cannot recover long-term missing detections from slow or stopped vehicles.
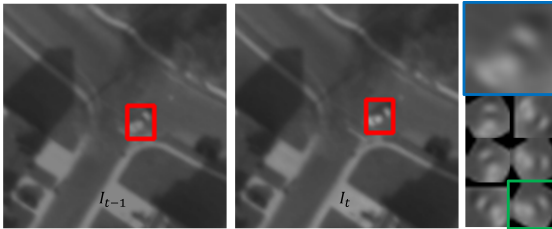
**Fig. 3.** A turning vehicle and rotation variants of $s_i$. Blue rectangle represents $s_i$ and the green rectangle represents the selected template with the highest NCC score.

Xiao *et al.* [41] propose to use appearance and shape templates to handle missing detections. To avoid drifting, they use road network information and consider pairwise spatial relation in optimization. However, road network information is not always available and considering spatial relations in Hungarian optimization is costly. Basharat *et al.* [42] apply an appearance-based tracker whenever detection association fails for a track or the motion of a target is slow. More recently, a hybrid approach that combines DBT with a regression-based tracker is proposed [6] to handle stop-then-go vehicles. Nevertheless, relying only on weak appearance information makes these trackers prone to drift and limits its ability in recovering missing detections.

## 3   Local Context Tracker

In this section, we present the details of our LCT. It explores context information to increase the robustness of tracker. Given a track pool $\mathcal{T}_e = \{T_1, T_2, ..., T_N\}$, which contains $N$ tracks from time 1 to $t-1$, our goal is to extend each track that ends at the previous frame $I_{t-1}$ to the current frame $I_t$. We introduce our hypothesis generation method using sparse and dense flow in Section 3.1. In Section 3.2, we formulate the tracking problem in graph optimization and find the optimal hypothesis tree considering appearance, motion, and context information.

### 3.1   Hypothesis Sampling

As mentioned in Section 1, the displacement of a target is in a very large range ($0 \sim 80$ pixels). In practice, it is not affordable to densely sample all possibilities. We propose a search mechanism which is based on motion information. We use the fact that visible targets are either fast enough to produce discontinuities in dense flow or slow enough to meet the small motion assumption of optical flow. For each track $T_i$, we construct a set of hypotheses $H_i$ for graph optimization in Section 3.2.
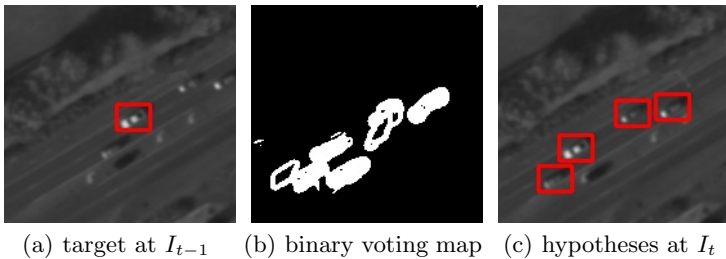
(a) target at $I_{t-1}$    (b) binary voting map    (c) hypotheses at $I_t$

**Fig. 4.** An example of generating hypotheses using a binary voting map.

**Sample from Dense Flow** In WAMI, most of the targets move fast and the small motion assumption fails in optical flow methods. Therefore, motion vectors from optical flow in a target region are usually noisy. Instead of using the motion vectors directly, we find discontinuities in dense flow for hypothesis sampling. We use $3 \times 3$ Sobel operators to get the gradient of flow and produce a binary voting map by thresholding the gradient magnitude. We search each target in a $160 \times 160$ window based on its location and size at previous frame. Here, we assume that the target size is consistent between consecutive frames because roads are typically flat and the UAV does not change its altitude drastically while flying. A valid hypothesis should have enough votes from the binary voting map. We set the threshold as one-fourth of the target size in all experiments.

In addition, we ensure that hypotheses match the target by using two kinds of templates including the target region at the previous frame $o_{t-1}^i$ and a "stable template" $s_i$ for the track $T_i$. We use normalized correlation coefficient (NCC) as template matching score. Non-maximum suppression is applied to reduce the number of candidates. The idea of stable template is to maintain a robust appearance model that only updates at high confidence to avoid gradually drifting [43]. In this work, we take the advantage of DBT for the confidence measure. We describe the update criterion of stable template in Section 5.

NCC template matching handles target motion in translation but not in rotation. Therefore, it fails when a target turns if we do not consider rotation of $s_i$, which may have different orientation compared with the target at the current frame. A collection of rotation variants of $s_i$ are used so that the tracker can deal with turning targets as shown in Fig. 3. We use 7 rotation variants ($-90$, $-60$, $-30$, 0, 30, 60, 90 degrees) in our implementation. The following rules are used to adopt hypotheses with strong appearance similarity:

$$NCC(h_t, o_{t-1}^i) > \phi, \tag{1}$$

$$MAX_j(NCC(h_t, s_i^j)) > \phi, \tag{2}$$

where $\phi$ is a constant which is set to 0.5, $h_t$ represents a template of a hypothesis candidate, and $s_i^j$ is the $j_{th}$ rotation variant of $s_i$. Fig. 4(c) shows an example of samples from the dense flow. Note that generating hypotheses directly from

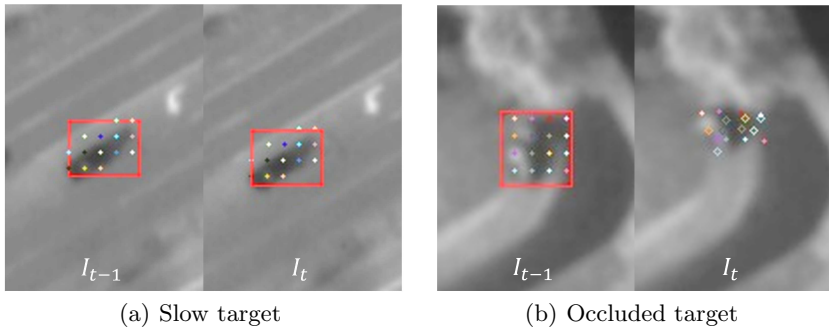(a) Slow target                    (b) Occluded target

**Fig. 5.** Different situations of sampling from the sparse flow. Points with the same color represent correspondences from the sparse optical flow. Hollow rhombuses show failure points in forward-backward check. (a) Optical flow tracks most points successfully when the target is slow and visible. A valid hypothesis in $I_t$ is shown in the red rectangle. (b) Most points fail in forward-backward check when the target is occluded. No hypothesis is generated in this case.

connected components of the voting map is not suitable since merged blobs appear frequently in practice. Fig. 4(b) shows an example of this situation.

**Sample from Sparse Flow** When a target moves slowly or stops, there is no discontinuity in dense flow between the target region and its local background region. Fortunately, small motion assumption of optical flow is valid in this case. We attempt to generate a hypothesis from densely sampled sparse flow to handle the situation that a target moves relatively slow or stops. For each pixel in the target bounding box at the previous frame, we use Lucas-Kanade optical flow [44] to track them at the current frame. Note that target hypotheses with fast motion are already covered by using dense flow in the previous section. Thus, we only use a relatively small window for sparse flow. In all our experiments, this window is set to $15 \times 15$ in pixels.

The main reason of using sparse flow is that we can apply a forward-backward check [31,33] efficiently to remove inconsistent correspondences. The median flow vector of remaining points is used to predict the hypothesis region. We avoid introducing false alarms by rejecting hypotheses without enough number of valid correspondences as shown in Fig. 5 (b). The threshold is set to $1/8$ of the target size at $t - 1$. Additionally, two NCC template matching policies, which are illustrated in equation 1 and 2, are applied to select valid hypotheses.

### 3.2   Optimization in Hypothesis Graph

Assuming independence between all tracks causes ambiguities between neighboring targets with similar appearance in WAMI. While considering every linking
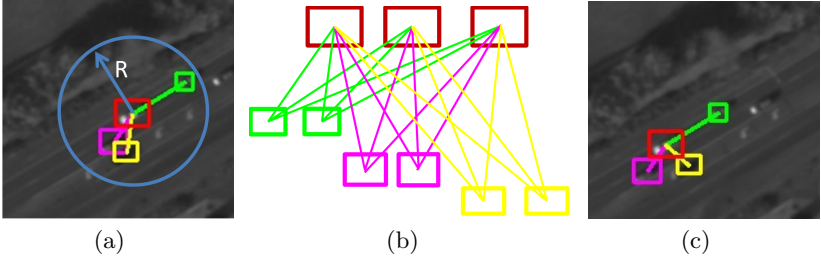
**Fig. 6.** (a) Context relation at $I_{t-1}$ (b) Hypothesis graph at time $t$ (c) The tracking result of LCT at $I_t$

possibilities between all hypotheses and all tracks is computationally expensive, we argue that each target trajectory is mainly affected by other targets with a similar motion direction in a local neighborhood. The movement of each vehicle is constrained to avoid collisions with its neighbors. For each target track $T_i \in \mathcal{T}_e$, we find its neighboring tracks with similar moving direction at $t-1$ as its "motion neighbors". We resolve the tracking problem for each $T_i$ and its motion neighbors at the same time. Fig. 6 (a) shows the neighborhood of the target (red). In our implementation, we use a search radius $R = 50$ pixels.

Given a target track $T_i$, its motion neighbors, and their hypothesis sets, we construct a hypothesis graph $G = (V, E)$ as shown in Fig. 6 (b). Each sample in a hypothesis set is represented by a node, while edges are constructed between every node from target track and every node from each of its motion neighbor. Let $k = \{k_t^i|_{i=1}^{|V|}, k_t^i \in V\}$ be a set of nodes that forms a tree, we formulate the tracking problem as maximizing the following objective function:

$$O(k) = \lambda \sum_{k_t^i \in V} U(k_t^i) + \sum_{(k_t^i, k_t^j) \in E} B(k_t^i, k_t^j), \qquad (3)$$

where $U(k_t^i)$ is the unary score at node $k_t^i$, and $B(k_t^i, k_t^j)$ is the binary score between nodes that forms an edge in $E$. We use $\lambda$, which is set to 3 in our experiments to leverage weighting between these two scores. Using the trajectory of $T_i$, we define $U(k_t^i)$ as:

$$U(k_t^i) = C_{app}(k_t^i) \cdot C_{mot}(k_t^i), \qquad (4)$$

where $C_{app}(k_t^i)$ is the appearance measurement based on NCC calculated between the template of hypothesis and the template of its corresponding target at the previous frame. $C_{mot}(k_t^i)$ is computed by multiplying the velocity similarity with the acceleration similarity. Here, we adopt the same velocity and acceleration measure as in [7], which uses different Gaussian kernels for magnitude and orientation components.

We calculate a binary score $B(k_t^i, k_t^j)$ according to the local context. This provides discriminating power against id-switches and false alarms. Given a target
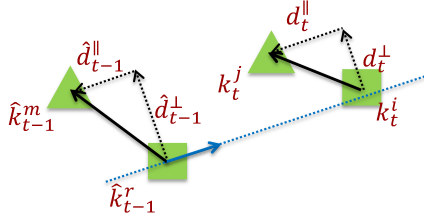
**Fig. 7.** Context constraint model at time $t-1$ and one of the hypothesis pairs at time $t$. Triangles represent the motion neighbor and its hypothesis while square represents the target and its hypothesis. The blue arrow shows the latest velocity vector of the target.

track and one of its motion neighbors, we observe that the relative displacement between them is more flexible along the moving direction of the target than along its normal direction. Therefore, we decompose the displacement according to the target velocity vector and use different Gaussian kernels for structure consistency in different components. Fig. 7 shows our local context modeling. $\hat{d}_{t-1}^{\parallel}$ and $\hat{d}_{t-1}^{\perp}$ are spatial relations between the target and the neighbor while $d_t^{\parallel}$ and $d_t^{\perp}$ represent one of the hypothesis pairs between them. The binary score is defined as:

$$B(k_t^i, k_t^j) = C_{\parallel}(k_t^i, k_t^j) \cdot C_{\perp}(k_t^i, k_t^j), \tag{5}$$

$$C_{\parallel}(k_t^i, k_t^j) = exp(-\alpha \cdot Abs(\hat{d}_{t-1}^{\parallel} - d_t^{\parallel})), \tag{6}$$

$$C_{\perp}(k_t^i, k_t^j) = exp(-\beta \cdot Abs(\hat{d}_{t-1}^{\perp} - d_t^{\perp})), \tag{7}$$

where $Abs(\cdot)$ returns the absolute value. We set $\alpha = 0.01$ and $\beta = 0.05$ as constants in our experiment to penalize deformation in different components.

The optimization can be solved by using dynamic programming efficiently in polynomial time [45]. Fig. 6 (c) shows an example LCT tracking result. By considering the local context as well as the trajectory history, LCT does not drift to other targets or false alarms.

## 4   Detection Based Tracker

One of the difficulties in coupling DBT and LCT is: LCT assumes that the observation at the end of a track contains only one target. This assumption fails when a merged detection occurs. The appearance of these detections accounts for multiple targets. The center of them may be far from any target. This makes both appearance and motion information unreliable for LCT initialization. Unlike most DBT methods that do not handle this situation in the detection association level, we address this problem by generating additional hypotheses rather than merely adopt motion detections directly.
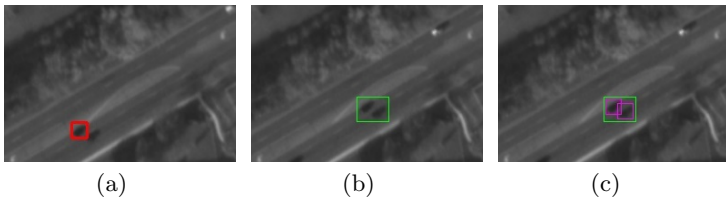
**Fig. 8.** Generate additional nodes from merged detection. (a) parent node (red) (b) child node (a merged detection) (green) (c) child node with additional nodes (purple).

Our DBT is based on the motion proprogation approach proposed in [7]. It builds detection trees layer-by-layer from motion detections at each frame in a sliding temporal window, which is set to 8 frames in our implementation. Each tree node represents a detection. Edges are constructed iteratively based on the best path from the root to each node at the previous frame. Each detection tree produces at most one tracklet by finding an optimal path from root to leaf based on appearance and motion consistency.

Using this framework, we identify abnormal changes of detection size for each edge in a detection tree. Then, we estimate additional hypotheses from these cases to improve tracking accuracy. Let $d_{t-1}$ be a parent detection of a child node, $d_t$. If $1.5 \cdot Size(d_{t-1}) < Size(d_t)$, we call $d_t$ a "potential merged detection". $Size(d)$ returns the size of a detection $d$. Fig. 8 (a) shows an example of $d_{t-1}$ and (b) shows an example of $d_t$.

Appearance consistency is used to generate hypotheses from a potential merged detection. We scan the region with the template of the parent node and calculate NCC score. Local maximums with scores larger than $\phi$ produce new tree nodes. We insert the new nodes into the detection tree and add a link between each new node and $d_{t-1}$ to update the detection tree. In Fig. 8 (c), purple rectangles show two additional nodes from a potential merged detection in the green rectangle. Their parent node is shown in Fig. 8 (a) in the red rectangle. By inserting these hypotheses, the optimal solution avoids choosing merged detections. This is because the additional hypothesis at the correct location leads to higher appearance and motion consistency than a merged detection, which is suboptimal.

## 5   Track Association

With new tracklets from DBT and tracking results from LCT, we describe how to generate the final result in this section. This result is further used to update the track pool. We find associations between new tracklets $\mathcal{T}_n$ from DBT and existing tracks $\mathcal{T}_e$ in the track pool in two steps. Let $t_{st}^n, t_{ed}^n, t_{st}^e, t_{ed}^e$ be the first frame index and the end frame index of a new tracklet $T_n$ and an existing track $T_e$. We define the association score $A(T_e, T_n) = S_p(T_e, T_n) \times S_v(T_e, T_n)$,

where $S_p(T_e, T_n)$ and $S_v(T_e, T_n)$ represents the position and velocity similarity. Since most successfully associated pairs are with overlapping time, given $T_n$, the first step is to find the best match among existing tracks with $t_{st}^e \geq t_{ed}^n$. In this case, we calculate $S_p(T_e, T_n)$ as:

$$S_p(T_e, T_n) = Match(T_e, T_n)/(t_{st}^e - t_{ed}^n + 1), \tag{8}$$

where $Match(T_e, T_n)$ returns the number of matched observation pairs in the overlapping period. Here, we define an observation as a bounding box region in a track or a tracklet. Given a time index, if the center of an observation from a tracklet lies in the observation from an existing track or vice versa, we treat them as matched observation. $S_v(T_e, T_n)$ is based on the Euclidean distance between the velocities from $T_e$ and $T_n$ at time $t_{st}^n$ with a Gaussian kernel. The constant threshold $\zeta = 0.6$ is used to adopt valid association.

If the successful association cannot be found in overlapping existing tracks, we interpolate the first observation of $T_n$ to time $t_{ed}^e$ of each non-overlapping track using a linear motion model and compute $S_p(T_e, T_n) = exp(-0.01 * Dist(T_e, T_n))$, where $Dist(T_e, T_n)$ is the Euclidean distance between the interpolated observation and the last observation of $T_e$. $S_v(T_e, T_n)$ is calculated using the Euclidean distance between velocities at $t_{st}^n$ in $T_n$ and at $t_{ed}^e$ in $T_e$. Again, $\zeta$ is used to accept successful association pairs. Unassociated tracklets initialize new tracks in the track pool.

Given a track in the track pool, if a tracklet is associated with it and LCT also tracks the target successfully, we append the observation at time $t$ with larger NCC score compared with the observation at time $t - 1$ in the tracklet. If only one of the trackers produce a valid result, we will extend the track using this result.

Updating template robustly is difficult in tracking problems. A template has to adapt to changes of target appearance. At the same time, the template should not drift to background gradually because of the update. Fortunately, DBT tracklet provides strong evidence in target existence without using template information. Thus, we update the stable template of a track whenever a valid tracklet is associated with it. The latest observation of the track is used to update the rotation variants of the stable template.

## 6   Experiments

### 6.1   Setup

We compare our methods with state-of-the-art trackers [9,38,6,21,7] on two WAMI sequences. In order to get a fair comparison, we use the same motion detection result as input for all trackers. The detection approach is based on background subtraction with 3-D stabilization [46], which reduces most false alarms from parallax effect.

We obtain executables of [9] and [6] from the authors of [6]. The programs of [38,7] are provided by their authors. We do not change parameters for the

| Method | Recall | Precision | FP/F | FP/GT | MODA | S/T | B/T | MOTA |
|---|---|---|---|---|---|---|---|---|
| Reilly *et al.* [9] | 0.573 | 0.94 | 0.887 | 0.037 | 0.536 | 0.851 | 1.293 | 0.522 |
| Prokaj *et al.* [38] | 0.504 | 0.985 | 0.18 | 0.007 | 0.497 | 0.249 | 1.515 | 0.493 |
| Prokaj *et al.* [6] | 0.539 | 0.96 | 0.548 | 0.023 | 0.516 | 0.237 | 1.022 | 0.512 |
| Rezatofighi *et al.* [21] | 0.44 | 0.74 | 3.746 | 0.155 | 0.285 | 0.529 | 2.855 | 0.276 |
| Chen *et al.* [7] | 0.55 | 0.987 | 0.171 | 0.007 | 0.543 | 0.2 | 0.5 | 0.54 |
| Ours (DBT only) | 0.537 | 0.985 | 0.195 | 0.008 | 0.529 | 0.061 | 0.573 | 0.528 |
| Ours (DBT+LCT) | **0.606** | **0.99** | **0.145** | **0.006** | **0.6** | **0.015** | **0.317** | **0.599** |

**Table 1.** Comparison of tracking results of the WPFAB 2009 sequence. Best results in each indexes are shown in bold text.

above methods except for the metadata information, which includes frame rate and the ground sampling distance of imagery. Evaluation of [21] is based on the MATLAB code provided by its authors. We tune parameters based on their helpful advice in order to apply it on low-frame-rate WAMI.

### 6.2   Evaluation Metrics

Our quantitative evaluation is based on commonly used metrics which are also adopted in [6,7], including *recall*: number of true positive detection/number of ground truth detection; *precision*: number of true positive detection/number of detection in tracks; *false positive per frame (FP/F)*; *false positive per ground truth detection(FP/GT)*; *multiple object detection accuracy (MODA)*; *number of swaps (id-switches) per track (S/T)*; *number of breaks per track (B/T)*; *multiple object tracking accuracy (MOTA)*. The definition of *MODA* and *MOTA* can be found in [47].

### 6.3   Results on WPAFB 2009 Sequence

The sequence is selected and preprocessed by the authors of [6] from a public WAMI dataset [11]. This dataset is recorded around 1 Hz and it provides ground truth labels for vehicles. The sequence covers a 429 m × 429 m suburb area in OH, USA with 1408 pixels × 1408 pixels. The ground sampling distance is around 0.3 meters per pixel. There are 1025 frames and 410 tracks. Several stop-then-go situation happens in the scene. Many merged detections appear when vehicles are close to each other.

Fig. 9 shows the qualitative comparison between our methods with DBT only and with DBT + LCT. The target slows down then stops to make a left turn. Motion detection fails from frame 11 to frame 21. Without LCT, the track breaks into two tracks as shown in the first row of Fig. 9, the yellow arrow indicates the missing detection at frame 17. On the contrary, DBT + LCT successfully recovers missing detections and continue the track.

Table 1 shows the quantitative results. By inserting additional hypotheses from potential merged blobs in the optimization, our DBT reduces the number
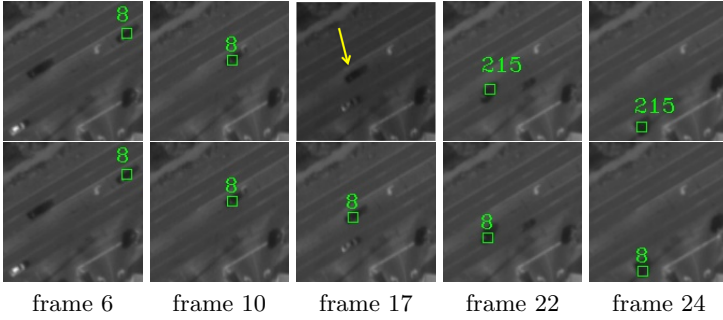
frame 6          frame 10          frame 17          frame 22          frame 24

**Fig. 9.** Snapshots of results from our two methods. First row shows the result using our DBT only. Second row shows the result by coupling DBT + LCT. Tracking results of other targets are not shown for clarity.

| Method | Recall | Precision | FP/F | FP/GT | MODA | S/T | B/T | MOTA |
|---|---|---|---|---|---|---|---|---|
| Reilly *et al.* [9] | 0.562 | 0.841 | 2.344 | 0.106 | 0.455 | 0.545 | 1.591 | 0.444 |
| Prokaj *et al.* [38] | 0.519 | 0.953 | 0.563 | 0.025 | 0.493 | 0.205 | 1 | 0.489 |
| Prokaj *et al.* [6] | 0.569 | 0.905 | 1.323 | 0.06 | 0.509 | 0.773 | 1.455 | 0.493 |
| Rezatofighi *et al.* [21] | 0.554 | 0.842 | 2.292 | 0.104 | 0.45 | 0.227 | 1.205 | 0.445 |
| Chen *et al.* [7] | 0.5 | 0.973 | 0.302 | 0.014 | 0.486 | 0.023 | 0.636 | 0.486 |
| Ours (DBT only) | 0.497 | **0.999** | **0.01** | **0.001** | 0.497 | **0** | 0.568 | 0.497 |
| Ours (DBT + LCT) | **0.761** | **0.999** | **0.01** | **0.001** | **0.761** | **0** | **0.159** | **0.761** |

**Table 2.** Comparison of tracking results of the Rochester sequence. Best results in each indexes are shown in bold text.

of *S/T* significantly. Combining LCT with DBT, we recover many missing detections between DBT. This increases *Recall* and reduces *B/T*. Note that false alarms and id-switches further decrease because LCT reduces the gap between DBT tracks and partially avoid errors from linear interpolation. Our method outperforms state-of-the-art methods in all indexes.

### 6.4   Results on Rochester Sequence

We select a sequence with 650 pixels × 650 pixels (250 m × 250 m) from Rochester dataset, which is captured from the city of Rochester, NY, USA. This imagery is recorded at 2 Hz and the ground sampling distance is around 0.38 meters per pixel. We manually label ground truth for each target from the frame it starts to move to the frame right before it leaves the scene. The sequence contains 96 frames and 44 tracks. Since it is an urban-view dataset, many vehicles stop at intersections for a long period.

The quantitive results are shown in Table 2. We do not produce any swaps in this sequence in both of our methods. Using LCT further improves more than

|           | Reilly *et al.* [9] | Prokaj *et al.* [38] | Prokaj *et al.* [6] | Chen *et al.* [7] | Ours (DBT) | Ours (DBT+LCT) |
|-----------|---------------------|----------------------|---------------------|-------------------|------------|----------------|
| WPFAB     | 0.121               | 0.094                | 9.413               | 0.21              | 0.124      | 1.009          |
| Rochester | 0.093               | 0.026                | 15.166              | 0.064             | 0.018      | 0.361          |

**Table 3.** Comparison of computation time in seconds per frame in the WPFAB 2009 sequence and the Rochester sequence.

19% in *recall* from move-then-stop situations. Methods that rely on only motion detections [9,38,21,7] fail in these cases.

This sequence has higher ground sampling distance and lower contrast than the WPFAB 2009 sequence. These factors make the appearance of a target less discriminating from the background. Therefore, although [6] performs the second best among all methods in *recall*, it introduces many false alarms because the regression based tracker may drift to visually similar background patches. On the contrary, LCT recovers missing detections by exploring stronger evidence based on context information. We maintain high precision compared with other methods. Again, DBT + LCT is clearly the leader in all indexes.

### 6.5   Computation Time

Our approach is implemented using C++ on a desktop with 3.6GHz CPU, 16GB memory and a NVIDIA GeForce GTX 580 GPU. We use GPU only for dense flow computation using FlowLib [48]. Table 3 shows the computational cost of methods with C++ implementation. Our DBT achieves similar computation time compared with other detection association methods [9,38,7].

The major limitation of our DBT + LCT is the higher computational cost compared with other DBTs [9,38,7]. In the WPFAB 2009 sequence, dense flow calculation takes 0.482 seconds which accounts for nearly half of computation time in our DBT + LCT. However, we are still more efficient than the state-of-the-art hybrid approach [6]. In the Rochester sequence, since the image size is smaller, the flow computation time reduced to 0.013 sec/frame, and our DBT + LCT takes 0.361 sec/frame. Notice that the computation time of [6] increases because appearance models have to update whenever regression trackers lose track. Failure of the regression tracker happens more frequently in the Rochester sequence where target appearance is less discriminating.

## 7   Conclusions

Existing multi-target tracking approaches in WAMI have limitations in recovering long-term missing detections from slow or stopped targets. We propose a unified approach which couples LCT with DBT. Instead of merely relying on an appearance model, LCT explores context information between a target and its

motion neighbors so that it is robust against neighboring distracters and background clutter. Hypotheses are generated by dense and sparse flow efficiently. We reduce id-switches by handling merged detections in DBT with additional hypotheses. Our DBT+LCT significantly improves tracking results compared with state-of-the-art methods on two WAMI sequences.

Our future work is to reduce the computation time by using parallel programming. Note that both DBT and LCT can be processed in parallel for each target; therefore, we expect obvious reduction in computation time. Furthermore, we want to infer the association between tracks under long-term occlusion, which often happens in urban view scenarios.

# References

1. Reilly, V., Solmaz, B., Shah, M.: Shadow casting out of plane (scoop) candidates for human and vehicle detection in aerial imagery. Int. J. Comput. Vision **101**(2) (January 2013) 350–366
2. Chen, B.J., Medioni, G.: 3-d mediated detection and tracking in wide area aerial surveillance. In: Applications of Computer Vision (WACV), 2015 IEEE Winter Conference on. (Jan 2015) 396–403
3. Liao, H.H., Lin, Y., Medioni, G.: Aerial 3d reconstruction with line-constrained dynamic programming. In: Computer Vision (ICCV), 2011 IEEE International Conference on. (Nov 2011) 1855–1862
4. Kang, Z., Medioni, G.: 3d urban reconstruction from wide area aerial surveillance video. In: Applications and Computer Vision Workshops (WACVW), 2015 IEEE Winter. (Jan 2015) 28–35
5. Shu, T., Xie, D., Rothrock, B., Todorovic, S., Zhu, S.C.: Joint inference of groups, events and human roles in aerial videos. In: CVPR. (2015)
6. Prokaj, J., Medioni, G.: Persistent tracking for wide area aerial surveillance. In: Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on. (June 2014) 1186–1193
7. Chen, B.J., Medioni, G.: Motion propagation detection association for multi-target tracking in wide area aerial surveillance. In: Advanced Video and Signal Based Surveillance (AVSS), 2015 12th IEEE International Conference on. (Aug 2015) 1–6
8. Perera, A., Srinivas, C., Hoogs, A., Brooksby, G., Hu, W.: Multi-object tracking through simultaneous long occlusions and split-merge conditions. In: Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on. Volume 1. (June 2006) 666–673
9. Reilly, V., Idrees, H., Shah, M.: Detection and tracking of large number of targets in wide area surveillance. In: Proceedings of the 11th European Conference on Computer Vision Conference on Computer Vision: Part III. ECCV'10, Berlin, Heidelberg, Springer-Verlag (2010) 186–199
10. Saleemi, I., Shah, M.: Multiframe many—many point correspondence for vehicle tracking in high density wide area aerial videos. Int. J. Comput. Vision **104**(2) (September 2013) 198–219
11. AFRL: WPAFB 2009 dataset. `https://www.sdms.afrl.af.mil/index.php?collection=wpafb2009`
12. Andriluka, M., Roth, S., Schiele, B.: People-tracking-by-detection and people-detection-by-tracking. In: Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on. (June 2008) 1–8
13. Breitenstein, M., Reichlin, F., Leibe, B., Koller-Meier, E., Van Gool, L.: Robust tracking-by-detection using a detector confidence particle filter. In: Computer Vision, 2009 IEEE 12th International Conference on. (Sept 2009) 1515–1522
14. Collins, R.: Multitarget data association with higher-order motion models. In: Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on. (June 2012) 1744–1751
15. Yang, B., Nevatia, R.: Online learned discriminative part-based appearance models for multi-human tracking. In: Proceedings of the 12th European Conference on Computer Vision - Volume Part I. ECCV'12, Berlin, Heidelberg, Springer-Verlag (2012) 484–498

16. Wang, W., Nevatia, R., Yang, B.: Beyond pedestrians: A hybrid approach of tracking multiple articulating humans. In: Applications of Computer Vision (WACV), 2015 IEEE Winter Conference on. (Jan 2015) 132–139

17. Fortmann, T.E., Bar-Shalom, Y., Scheffe, M.: Sonar tracking of multiple targets using joint probabilistic data association. Oceanic Engineering, IEEE Journal of **8**(3) (Jul 1983) 173–184

18. Reid, D.: An algorithm for tracking multiple targets. Automatic Control, IEEE Transactions on **24**(6) (Dec 1979) 843–854

19. Kang, J., Cohen, I., Medioni, G.: Soccer player tracking across uncalibrated camera streams. In: In: Joint IEEE International Workshop on Visual Surveillance and Performance Evaluation of Tracking and Surveillance (VS-PETS) In Conjunction with ICCV. (2003) 172–179

20. Schulz, D., Burgard, W., Fox, D., Cremers, A.: Tracking multiple moving targets with a mobile robot using particle filters and statistical data association. In: Robotics and Automation, 2001. Proceedings 2001 ICRA. IEEE International Conference on. Volume 2. (2001) 1665–1670 vol.2

21. Rezatofighi, S.H., Milan, A., Zhang, Z., Shi, Q., Dick, A., Reid, I.: Joint probabilistic data association revisited. In: ICCV. (2015)

22. Cox, I.J., Hingorani, S.L.: An efficient implementation and evaluation of reid's multiple hypothesis tracking algorithm for visual tracking (1994)

23. Chenouard, N., Bloch, I., Olivo-Marin, J.: Multiple hypothesis tracking for cluttered biological image sequences. Pattern Analysis and Machine Intelligence, IEEE Transactions on **35**(11) (Nov 2013) 2736–3750

24. Kim, C., Li, F., Ciptadi, A., Rehg, J.M.: Multiple hypothesis tracking revisited. In: Computer Vision (ICCV), IEEE International Conference on, IEEE (December 2015)

25. Berclaz, J., Fleuret, F., Turetken, E., Fua, P.: Multiple object tracking using k-shortest paths optimization. Pattern Analysis and Machine Intelligence, IEEE Transactions on **33**(9) (Sept 2011) 1806–1819

26. Zhang, L., Li, Y., Nevatia, R.: Global data association for multi-object tracking using network flows. In: Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on. (June 2008) 1–8

27. Butt, A., Collins, R.: Multi-target tracking by lagrangian relaxation to min-cost network flow. In: Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on. (June 2013) 1846–1853

28. Pirsiavash, H., Ramanan, D., Fowlkes, C.C.: Globally-optimal greedy algorithms for tracking a variable number of objects. In: Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on. (June 2011) 1201–1208

29. Ross, D.A., Lim, J., Lin, R.S., Yang, M.H.: Incremental learning for robust visual tracking. Int. J. Comput. Vision **77**(1-3) (May 2008) 125–141

30. Babenko, B., Yang, M.H., Belongie, S.: Visual tracking with online multiple instance learning. IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI) (August 2011)

31. Kalal, Z., Mikolajczyk, K., Matas, J.: Tracking-learning-detection. Pattern Analysis and Machine Intelligence, IEEE Transactions on **34**(7) (July 2012) 1409–1422

32. Kuo, C.H., Huang, C., Nevatia, R.: Multi-target tracking by on-line learned discriminative appearance models. In: Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on. (June 2010) 685–692

33. Xiang, Y., Alahi, A., Savarese, S.: Learning to track: Online multi-object tracking by decision making. In: International Conference on Computer Vision. (2015)

34. Dinh, T.B., Vo, N., Medioni, G.: Context tracker: Exploring supporters and distracters in unconstrained environments. In: Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on. (June 2011) 1177–1184
35. Grabner, H., Matas, J., Van Gool, L., Cattin, P.: Tracking the invisible: Learning where the object might be. In: Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on. (June 2010) 1285–1292
36. Yang, M., Wu, Y., Hua, G.: Context-aware visual tracking. Pattern Analysis and Machine Intelligence, IEEE Transactions on **31**(7) (July 2009) 1195–1209
37. Zhang, L., van der Maaten, L.: Preserving structure in model-free tracking. Pattern Analysis and Machine Intelligence, IEEE Transactions on **36**(4) (April 2014) 756–769
38. Prokaj, J., Duchaineau, M., Medioni, G.: Inferring tracklets for multi-object tracking. In: Computer Vision and Pattern Recognition Workshops (CVPRW), 2011 IEEE Computer Society Conference on. (June 2011) 37–44
39. Shi, X., Ling, H., Xing, J., Hu, W.: Multi-target tracking by rank-1 tensor approximation. In: Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on. (June 2013) 2387–2394
40. Keck, M., Galup, L., Stauffer, C.: Real-time tracking of low-resolution vehicles for wide-area persistent surveillance. In: Applications of Computer Vision (WACV), 2013 IEEE Workshop on. (Jan 2013) 441–448
41. Xiao, J., Cheng, H., Sawhney, H., Han, F.: Vehicle detection and tracking in wide field-of-view aerial video. In: Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on. (June 2010) 679–684
42. Basharat, A., Turek, M., Xu, Y., Atkins, C., Stoup, D., Fieldhouse, K., Tunison, P., Hoogs, A.: Real-time multi-target tracking at 210 megapixels/second in wide area motion imagery. In: Applications of Computer Vision (WACV), 2014 IEEE Winter Conference on. (March 2014) 839–846
43. Matthews, I., Ishikawa, T., Baker, S.: The template update problem. Pattern Analysis and Machine Intelligence, IEEE Transactions on **26**(6) (June 2004) 810–815
44. yves Bouguet, J.: Pyramidal implementation of the lucas kanade feature tracker. Intel Corporation, Microprocessor Research Labs (2000)
45. Felzenszwalb, P., Zabih, R.: Dynamic programming and graph algorithms in computer vision. Pattern Analysis and Machine Intelligence, IEEE Transactions on **33**(4) (April 2011) 721–740
46. Chen, B.J., Medioni, G.: Persistent 3d stabilization for aerial imagery. In: Applications of Computer Vision (WACV), 2015 IEEE Winter Conference on. (Mar 2016)
47. Kasturi, R., Goldgof, D., Soundararajan, P., Manohar, V., Garofolo, J., Bowers, R., Boonstra, M., Korzhova, V., Zhang, J.: Framework for performance evaluation of face, text, and vehicle detection and tracking in video: Data, metrics, and protocol. IEEE Transactions on Pattern Analysis and Machine Intelligence **31**(2) (Feb 2009) 319–336
48. Zach, C., Pock, T., Bischof, H.: A duality based approach for realtime TV-$L^1$ optical flow. In: Proceedings of the 29th DAGM conference on Pattern recognition. (2007)