

A Rank Minimization Approach to Fast Dynamic Event Detection and Track Matching in Video Sequences

Tao Ding

Department of Electrical Engineering,
The Pennsylvania State University,
University Park, PA 16802.

Mario Sznaiier

Octavia Camps
Electrical and Comp. Engineering Department,
Northeastern University,
Boston, MA 02115.

Abstract—This paper addresses the problems of track stitching and dynamic event detection in a sequence of frames. The input data consists of tracks, possibly fragmented due to occlusion, belonging to multiple targets. The goals are to (i) establish track identity across occlusion, and (ii) detect points where the motion modalities of these targets undergo substantial changes. The main result of the paper is a simple, computationally inexpensive approach that allows achieving these goals in a unified way. Given a continuous track, the main idea is to detect changes in the dynamics by parsing it into segments according to the complexity of the model required to explain the observed data. Intuitively, changes in this complexity correspond to points where the dynamics change. In turn, the problem of estimating the complexity of the underlying model can be reduced to estimating the rank of a Hankel matrix constructed from the observed data, leading to a simple algorithm, computationally no more expensive than a sequence of SVDs. Proceeding along the same lines, fragmented tracks corresponding to multiple targets can be linked by searching for sets corresponding to minimal complexity joint models. As we show in the paper, this problem can be reduced to a semi-definite optimization and efficiently solved.

I. INTRODUCTION

Dynamic vision and imaging – the confluence of dynamics, computer vision, image processing and control – is uniquely positioned to enhance the quality of life for large segments of the general public. Aware sensors endowed with tracking and scene analysis capabilities can prevent crime, reduce time response to emergency scenes and allow elderly people to continue living independently. Moreover, the investment required to accomplish these goals is relatively modest, since a large number of imaging sensors are already deployed and networked. The challenge now is to develop a theoretical framework that allows for *robustly* processing this vast amount of information, within the constraints imposed by the need for real time operation in dynamic, partially stochastic scenarios. The objective of this paper is to illustrate the central role that control theory can play in achieving these goals. In particular, we concentrate in a problem that must be solved in all of the applications mentioned above: tracking objects in a sequence of frames, establishing, if needed, track identity across occlusion. Challenges in designing a robust tracking algorithm arise from several factors, e.g. changing appearances, changes in illumination, clutter and occlusion. During the past decade extensive research has

been carried out in this area, leading to several techniques that address these effects (see for instance [1], [7], [8], [17], [18], [9] and references therein). In particular, a class of dynamics based trackers has been developed that combine simple dynamic models of the target motion with optimal filtering –(unscented) Kalman, particle– [14], [10], [11] to track in the presence of occlusion. In order to further improve robustness, Camps *et. al.* [2] use interpolation theory to learn the dynamics of the target, thus removing a potential source of fragility arising from a mismatch between the assumed and actual dynamics. The implicit assumption in all these methods is that *the dynamics of the target do not change*, e.g. the underlying model is time invariant. In principle, as pointed out in [2], changes in this model could be detected by performing a model (in)validation step. The main advantage of this approach resides in its ability to unequivocally establish that a change in the dynamics has taken place. However, the entailed computational complexity is not small, specially in the case of long sequences. In addition, this approach cannot handle cases where the events occur while the target is occluded, which requires, as a pre-requisite, being able to match tracklets across the occlusion.

The problem of tracklet matching has been addressed in a number of recent papers, e.g. [12], [19], [16], [3]. However, while successful, these methods are fairly involved.

This paper shows that well known ideas from control theory can be exploited to obtain a computationally inexpensive approach that allows for both stitching tracklets across occlusion and detecting changes in the dynamics of the target. Given a continuous track, the main idea is to detect changes in the dynamics by parsing it into segments according to the complexity of the model required to explain the observed data. Intuitively, changes in this complexity correspond to points where the dynamics change. In turn, estimating the order of the underlying model reduces to estimating the rank of a Hankel matrix constructed from the observed data. Proceeding along the same lines, fragmented tracks corresponding to multiple targets can be linked by searching for sets corresponding to minimal complexity joint models. In this case, this is accomplished by estimating the missing points corresponding to the *lowest* order model that jointly explains a given set of candidate tracks. As shown in the paper, this leads to a rank minimization problem, which in turn can be relaxed to a semi-definite optimization.

This work was supported by NSF grants ECS-0221562, ITR-0312558 and ECS-050166, and AFOSR grant FA9550-05-1-0437.

II. BACKGROUND RESULTS

The main idea underlying this paper is to model the evolution of target features as the output of (a possibly piecewise) linear time invariant model whose order must be estimated from the available experimental data. Specifically, following [2], we will collect the position of all relevant features of the target in a vector \mathbf{f} and assume that its evolution is governed by a model of the form:

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{e}_k \\ \mathbf{f}_k &= \mathbf{C}\mathbf{x}_k, \mathbf{y}_k = \mathbf{C}\mathbf{x}_k + \eta_k \end{aligned} \quad (1)$$

where $\mathbf{A} \in R^{n \times n}$, with $n \leq N_f$, the number of frames¹, and where $\mathbf{x}(\cdot)$, $\mathbf{e}(\cdot)$, $\mathbf{f}(\cdot)$ and $\mathbf{y}(\cdot)$ represent the states, an exogenous stochastic input with appropriate statistics, the actual value of the feature vector at time k , and its measurement, corrupted by additive noise η , respectively. Our goal is, given \mathbf{y} , to estimate the *minimum* n such that the model (1) holds. To this effect, we recall the following result [4]:

Theorem 1: Given an input/output sequence $\{\mathbf{e}_t, \mathbf{f}_t\}$ corresponding to the model (1), form the (Hankel) matrices:

$$\begin{aligned} \mathbf{H}_f(k, l) &\doteq \begin{bmatrix} \mathbf{f}_1 & \mathbf{f}_2 & \cdots & \mathbf{f}_l \\ \mathbf{f}_2 & \mathbf{f}_3 & \cdots & \mathbf{f}_{l+1} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{f}_k & \mathbf{f}_{k+1} & \cdots & \mathbf{f}_{k+l-1} \end{bmatrix} \\ \mathbf{H}_e(k, l) &\doteq \begin{bmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \cdots & \mathbf{e}_l \\ \mathbf{e}_2 & \mathbf{e}_3 & \cdots & \mathbf{e}_{l+1} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{e}_k & \mathbf{e}_{k+1} & \cdots & \mathbf{e}_{k+l-1} \end{bmatrix} \end{aligned} \quad (2)$$

$$\mathbf{H}_{f,e} \doteq [\mathbf{H}_f(k, l) \quad \mathbf{H}_e(k, l)]$$

where $l \geq k \gg n$. Then, if the input sequence $\{\mathbf{e}_t\}$ is sufficiently rich, i.e. $\text{rank}[\mathbf{H}_e(k, l)] \equiv k$, the order n of the model (1) satisfies:

$$\text{rank}[\mathbf{H}_{f,e}] = k + n \quad (3)$$

Remark 1: In the sequel, we will assume, by absorbing if necessary the dynamics of the stochastic input \mathbf{e} into the dynamics of the plant, that \mathbf{e} is an impulse, e.g. $\mathbf{e}_1 = 1$, $\mathbf{e}_i = 0$, $i > 1$. With this assumption (3) above reduces to

$$\text{rank}[\mathbf{H}_f(k, l)] = n \quad (4)$$

III. A RANK CRITERION FOR FAST EVENT DETECTION

Theorem 1 can be used to perform fast detection of changes in the motion modality of the target by simply searching for points where the rank of the Hankel matrix abruptly changes after having remained approximately constant. This corresponds to a formalization of the intuitive fact that trying to explain two different modalities (distinguished either by different dynamics or different inputs) using a

¹Note that this can be always assumed without loss of generality, since given N measurements of $\mathbf{f}(\cdot)$ and $\mathbf{e}(\cdot)$, there always exists a linear operator such that (1) is satisfied (Chapter 10 of [15])

single model will require considerable more complexity than that required to explain each modality alone. Note that the approach outlined above does not require explicitly finding the models (computationally expensive).

A potential difficulty here is that, rather than the actual feature positions \mathbf{f} , only the measurements $\mathbf{y} = \mathbf{f} + \eta$ corrupted by noise are available, and it is well known that rank computation is very sensitive to noise. To avoid this difficulty, begin by noting that the Hankel matrices corresponding to the actual and measured position are related by: $\mathbf{H}_y = \mathbf{H}_f + \mathbf{H}_\eta$, where \mathbf{H}_η denotes the Hankel matrix associated with the noise sequence $\eta(\cdot)$. Under ergodicity assumptions, $\mathbf{H}_\eta^T \mathbf{H}_\eta$ is an estimate of the covariance matrix of the noise [13]. Thus, noise measurements can be robustly handled by simply replacing $\text{rank}(\mathbf{H}_y)$ by $\text{NSV}_{\sigma_\eta}(\mathbf{H}_y)$, the number of singular values larger than σ_η , the standard deviation of the measurement noise. Dynamic events can then be robustly detected by searching for points where $\text{NSV}_{\sigma_\eta}(\mathbf{H}_y)$ changes. The efficiency of this approach is illustrated next with several examples.

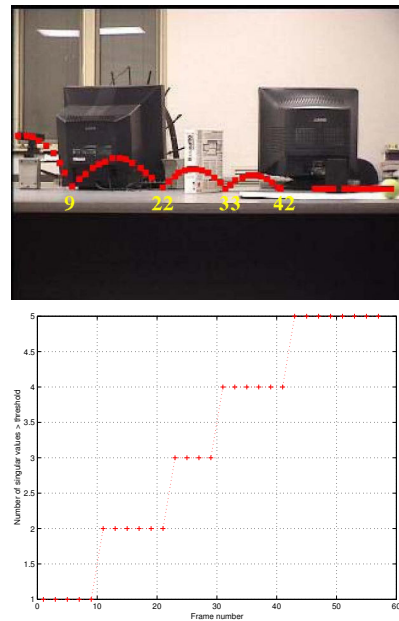


Fig. 1. Detecting events using the rank criterion. Top row: input sequence showing dynamic events at frames 10, 23, 34 and 43. Second row: corresponding NSV plot

Example 1: Consider the bouncing ball shown in Fig. 1. The dynamics change at frames 9, 22, 33 and 43, with the first 3 changes due to impact with the table and the last to the transition from bouncing to rolling motion. All these changes are clearly evidenced by the jumps in $\text{NSV}(\mathbf{H})$

Example 2: Change in human activity: This example consists of 78 frames of the sequence shown in Figure 2 of a moving person that abruptly switches from walking to jumping, starting at frame 51. As illustrated in the second row of the figure, this change is clearly shown in the plot of $\text{NSV}(\mathbf{H})$.

Example 3: Normal versus abnormal car slowdown. This example consists of two sequences showing a car undergoing

deceleration, as a result of a crash (Figure 3) (top row) and during normal braking (second row). The NSV plot corresponding to the crash exhibits a large jump starting around frame 45, indicating the occurrence of a dynamic event, while the NSV plot for normal deceleration, has a much smaller jump, around frame 50, as the car slows down to a stop.

IV. HANKEL MATRIX BASED TRACK MATCHING

In this section we turn our attention to the problem of establishing track identity across occlusion. As part of the process of solving this problem we develop an algorithm that allows for efficiently estimating the missing data that connects tracklets.

A. Track stitching: Estimating missing data.

Consider first the problem of estimating the missing data connecting two segments of the same track. Formally, this can be stated as:

Problem 1: Given two segments of a track, $\{\mathbf{y}_i^o\}$, $1 \leq i \leq r$ and $\{\mathbf{y}_j^o\}$, $s+1 \leq i \leq N_F$, with $r < s$ and $s-r \ll N_F$, estimate the missing values \mathbf{y}_k^* , $r+1 \leq k \leq s$ that are maximally consistent with the existing data, in the sense that the complete sequence is explained by the lowest possible order model.

From Theorem 1 it follows that the missing values \mathbf{y}^* can be optimally estimated by minimizing the rank of the corresponding Hankel matrix \mathbf{H}_y formed by combining \mathbf{y}^o and \mathbf{y}^* . A potential problem here is that rank minimization problems are known to be generically NP-hard. Thus, motivated by [5], [6], we will replace the rank minimization step by a convex relaxation that only entails solving a tractable, convex Linear Matrix Inequality (LMI) optimization, leading to the following Algorithm:

Algorithm 1: *HankelBasedTrackStitching*

Input: N observed values of \mathbf{y} ,

N_p , estimated number of missing points.

Output: Estimates \mathbf{y}^* of the missing data.

1. Form a Hankel matrix $H_{\hat{y}}$, where \hat{y} is the sequence formed by combining y and y^* ,

$$\mathbf{H}_y \doteq \begin{bmatrix} \hat{y}_1 & \hat{y}_2 & \cdots & \hat{y}_{\frac{N_F}{2}} \\ \hat{y}_2 & \hat{y}_3 & \cdots & \hat{y}_{\frac{N_F}{2}+1} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{y}_{\frac{N_F}{2}} & \cdots & \cdots & \hat{y}_{N_F} \end{bmatrix},$$

where $N_F \doteq N + N_p$.

2. Obtain the best prediction \mathbf{y}^* by solving the LMI optimization problem as follows.

$$\text{minimize } Tr(\mathbf{Y}) + Tr(\mathbf{Z})$$

$$\text{subject to } \begin{bmatrix} \mathbf{Y} & \mathbf{H}_y \\ \mathbf{H}_y^T & \mathbf{Z} \end{bmatrix} \geq 0$$

$$\{\mathbf{y}^* \in \mathcal{R}^2\}$$

where $\mathbf{Y}^T = \mathbf{Y}$, $\mathbf{Z}^T = \mathbf{Z}$, and $\mathbf{H}_y \in \mathcal{R}^{2N_F \times \frac{N_F}{2}}$.

Example 4: Detecting dynamic events under occlusion. This example illustrates the ability of the proposed methods to detect event changes, even if these events occur while the target is occluded. In this example, a jumping ball exhibits different dynamics in frames 1–42 and 43–59. The available data consists of 4 tracklets, labeled \mathcal{W}_{1-4} in Fig. 4(a), with estimated noise level $\sigma_\eta = 7.75$. Applying *Algorithm 1* to stitch the track led to the connecting trajectories shown in red in the figure. Finally, the **NSV** plot shown in Figure 4(b), clearly indicates points at where dynamic events take place.

B. Multiple Track Matching and Stitching.

The ideas discussed above can also be used to match and stitch tracks across occlusion. The main idea is to group tracks according to the complexity of the *simplest* model that explain the joint data. Specifically, given two measurement matrices \mathcal{W}_i and \mathcal{W}_j corresponding to two tracklets, and where $N_{S_j} > N_{F_i}$ the starting and ending frame indexes in \mathcal{W}_j and \mathcal{W}_i , respectively, a similarity measure between tracks can be defined proceeding as follows:

- 1) Use Algorithm 1 to stitch the tracklets. Define $\mathcal{W}_{i,j} \doteq [\mathcal{W}_i \quad \mathcal{W}^* \quad \mathcal{W}_j]$ where \mathcal{W}^* denotes the estimates of the missing measurements.
- 2) The similarity measure $\Gamma_{i,j}$ between tracklets $\{i, j\}$ is defined as:

$$\Gamma_{i,j} \doteq \begin{cases} -\infty, & \text{if temporal conflict exists;} \\ \frac{\text{NSV}_{\sigma_\eta}(\mathbf{H}_{\mathcal{W}_i}) + \text{NSV}_{\sigma_\eta}(\mathbf{H}_{\mathcal{W}_j})}{\text{NSV}_{\sigma_\eta}(\mathbf{H}_{\mathcal{W}_{i,j}})} - 1 & \end{cases} \quad (5)$$

Intuitively, if \mathcal{W}_i and \mathcal{W}_j are samples of the same trajectory, then $\text{rank}(\mathbf{H}_{\mathcal{W}_i}) = \text{rank}(\mathbf{H}_{\mathcal{W}_j}) = \text{rank}(\mathbf{H}_{\mathcal{W}_{i,j}})$ and hence $\Gamma_{i,j} = 1$. On the other hand if \mathcal{W}_i and \mathcal{W}_j are uncorrelated, $\Gamma_{i,j} \approx 0$. The definition above formalizes this idea, using NSV_σ in lieu of $\text{rank}(\cdot)$ to improve robustness against measurement noise. Once $\Gamma_{i,j}$ is computed for all pairs that do not exhibit temporal conflicts (e.g. one track starting before the end of the second), tracks can be matched by simply looking for pair (i, j) corresponding to the largest values $\Gamma_{i,j}$. These ideas are summarized in *Algorithm 2*.

Algorithm 2: HANKEL MATRIX BASED TRACK MATCHING

Input: measurements matrices \mathcal{W}_i ;

n_T , total tracklet number;

noise standard deviation σ_η .

Output: Similarity matrix Γ

for all $i \neq j \in \{1, \dots, n_T\}$ **do**

if there exists temporal conflict

 Set $\Gamma_{i,j}$ as $-\infty$.

else

 Apply **Algorithm 1** to \mathcal{W}_i and \mathcal{W}_j to find $\mathcal{W}_{i,j}$.

 Use (5) to calculate $\Gamma_{i,j}$.

end if

end for

Find the best correspondence guided by Γ .

Example 5: Multi-target track matching under occlusion.

TABLE I
SIMILARITY MATRIX FOR THE TWO BALLS EXAMPLE.

i	$\Gamma_{i,1}$	$\Gamma_{i,2}$	$\Gamma_{i,3}$	$\Gamma_{i,4}$
1	NA			
2	$-\infty$	NA		
3	1^\dagger	-0.17	NA	
4	-0.29	0.14^\dagger	$-\infty$	NA

TABLE II
SIMILARITY MATRIX FOR THE CAR AND BALL EXAMPLE.

i	$\Gamma_{i,1}$	$\Gamma_{i,2}$	$\Gamma_{i,3}$	$\Gamma_{i,4}$
1	NA			
2	$-\infty$	NA		
3	1^\dagger	-0.38	NA	
4	0	0.33^\dagger	$-\infty$	NA

Figure 5 shows 34 frames of a partially occluded sequence of two balls with different dynamic behavior. Applying *Algorithm 2* to the 4 tracklets (using $\sigma_\eta = 3.5$) yields $\mathbf{NSV}(\mathcal{W}_1) = \mathbf{NSV}(\mathcal{W}_3) = 1$, $\mathbf{NSV}(\mathcal{W}_2) = \mathbf{NSV}(\mathcal{W}_4) = 4$, $\mathbf{NSV}(\mathcal{W}_{1,3}) = 1$, $\mathbf{NSV}(\mathcal{W}_{2,4}) = 7$, $\mathbf{NSV}(\mathcal{W}_{1,4}) = 7$, and $\mathbf{NSV}(\mathcal{W}_{2,3}) = 6$. The resulting similarity matrix Γ is given in Table I. As shown in Fig. 5 grouping tracks according to this matrix indeed leads to the correct assignments.

The next example illustrates the ability of the method to exploit dynamical information to match partially overlapping tracks. It consists of 49 frames of the sequence shown in Figure 6, containing two moving objects: a ball and a car, the latter appearing only after frame 16. The similarity matrix Γ shown in Table II shows that $\mathcal{W}_{1,3}$ and $\mathcal{W}_{2,4}$ are the most consistent tracklets.

V. CONCLUSIONS AND FURTHER RESEARCH

In this paper we addressed the problem of multi-target dynamical event detection using fragmented tracks. In order to solve this problem, we introduced algorithms for (i) establishing track identity across occlusion, (ii) estimating missing data and (iii) analyzing the (reconstructed) track to establish points where dynamical events took place. The underlying idea in all cases is that tracks corresponding to a single target who is not undergoing dynamic events can be explained by a model whose complexity is lower than that required to jointly explain different dynamics. The latter situation can be due for instance to having different targets or a single target performing different activities. In turn, by exploiting well known results from Systems Theory, the problem of estimating the order of the model, can be reduced to computing the rank of a Hankel matrix constructed from the experimental data. This observation leads to fast, computationally simple algorithms that do not require finding explicit models. The effectiveness of this technique was illustrated using several examples.

These results point out to the central role that control theory can play in developing a comprehensive framework leading to robust dynamic vision systems. In turn, dynamic

vision can provide a rich environment to both draw inspiration from and test new developments in systems theory. For instance, the applications addressed in this paper point out, among others, to the need for further research into low complexity nonlinear identification methods and the development of worst-case identification methods for switched systems that are not necessarily ℓ^2 stable (to allow for parsing video sequences into different activities).

REFERENCES

- [1] M. J. Black and A. D. Jepson. Eigentracking: Robust matching and tracking of articulated objects using a view-based representation. *International Journal of Computer Vision*, 26(1):63–84, 1998.
- [2] O. Camps, H. Li, M. C. Mazzaro, and M. Sznaiier. A caratheodory-fejer approach to robust multiframe tracking”, computer vision. In *Proceedings of ICCV 2003*, volume 2, pages 1048–1055. IEEE, 2003.
- [3] M. T. Chan, A. Hoogs, R. Bhotika, A. Perera, J. Schmiederer, and G. Doretto. Joint recognition of complex events and track matching. In *Proceedings of CVPR 2006*, volume 2, pages 1615–1622. IEEE, 2006.
- [4] B. de Moor, M. Moonen, L. Vandenbergh, and J. Vandewalle. A geometrical approach for the identification of state space models with singular value decomposition. In *Proceedings of ICASSP–88*, pages 2244–2247, 1998.
- [5] M. Fazel, H. Hindi, and S. Boyd. Rank minimization and applications in system theory. In *Proceedings of American Control Conf. 2004*, volume 4, pages 3273–3278. AACC, 2004.
- [6] M. Fazel, H. Hindi, and S. P. Boyd. Log-det heuristic for matrix rank minimization with applications to hankel and euclidean distance matrices. In *Proceedings of American Control Conf. 2003*, volume 3, pages 2156–2162. AACC, 2003.
- [7] W. E. L. Grimson, C. Stauffer, R. Romano, and L. Lee. Using adaptive tracking to classify and monitor activities in a site. In *Proceedings of CVPR 1998*, volume 1, pages 22–29. IEEE, 1998.
- [8] G. Hager and P. Belhumeur. Efficient region tracking with parametric models of geometry and illumination. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 20(10):1025–1039, 1997.
- [9] M. Irani and P. Anandan. Unified approach to moving object detection in 2d and 3d scenes. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 20(6):577–589, 1998.
- [10] M. Isard and A. Blake. Condensation - conditional density propagation for visual tracking. *International Journal of Computer Vision*, 29(1):5–28, 1998.
- [11] S. Julier, J. Uhlmann, and H. F. Durrant-Whyte. A new approach for filtering nonlinear systems. In *Proceedings of ACC 1995*, volume 1, pages 1628–1632. IEEE, 1995.
- [12] R. Kaucic, A. G. A. Perera, G. Brooksby, J. Kaufhold, and A. Hoogs. A unified framework for tracking through occlusions and across sensor gaps. In *Proceedings of CVPR 2005*, volume 1, pages 990–997. IEEE, 2005.
- [13] R. Lublinerman, M. Sznaiier, and O. Camps. Dynamics based robust motion segmentation. In *Proceedings of CVPR 2006*, volume 1, pages 17–22. IEEE, 2006.
- [14] B. North, A. Blake, M. Isard, and J. Rittscher. Learning and classification of complex dynamics. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 22(9):1016–1034, 2000.
- [15] R. Sanchez Pena and M. Sznaiier. *Robust Systems Theory and Applications*. Wiley & Sons, Inc., 1998.
- [16] M. Piccardi and E. D. Cheng. Multi-frame moving object track matching based on an incremental major color spectrum histogram matching algorithm. In *Proceedings of CVPR 2005*, volume 3, pages 19–24. IEEE, 2005.
- [17] J. Shi and C. Tomasi. Good features to track. In *Proceedings of CVPR 1994*, volume 1, pages 593–600. IEEE, 1994.
- [18] C. R. Wen, A. Azarbayejani, T. Darrell, and A. P. Pentland. Pfunder: real-time tracking of the human body. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 19(7):780–785, 1997.
- [19] Y. Zhou and H. Tao. A background layer model for object tracking through occlusion. In *Proceedings of ICCV 2003*, volume 1, pages 1079–1085. IEEE, 2003.

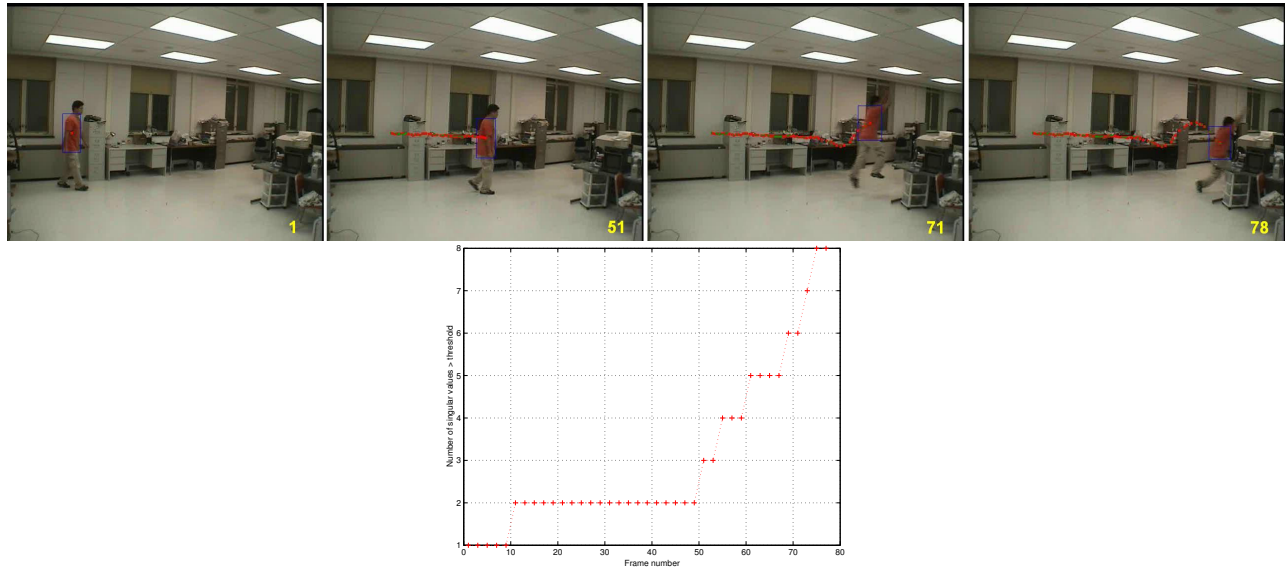


Fig. 2. Event detection. Top: transition from walking to jumping. Bottom: Corresponding NSV plot.

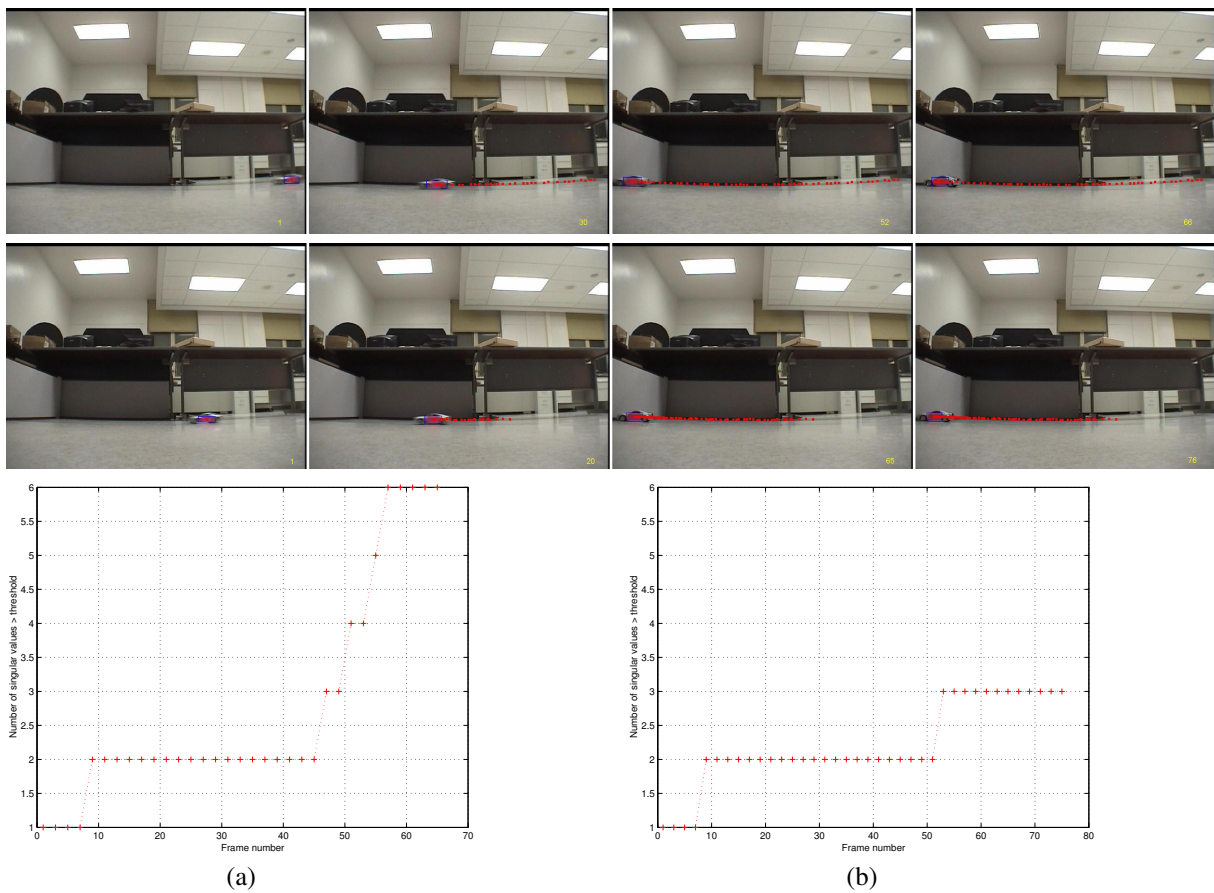


Fig. 3. Event detection. Top row: car crashing. Middle row: normal deceleration. Bottom row. (a) NSV plot for the crash case. (b) NSV plot for normal deceleration.

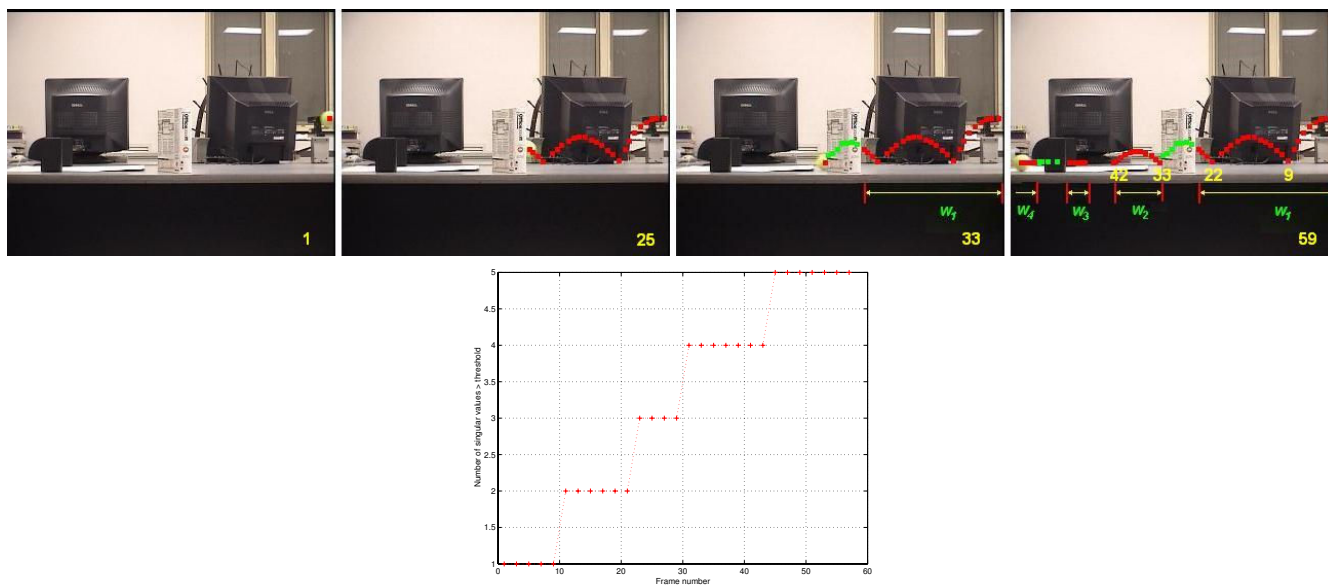


Fig. 4. Event detection: (top) Detecting occluded events. (bottom) NSV plot showing events.

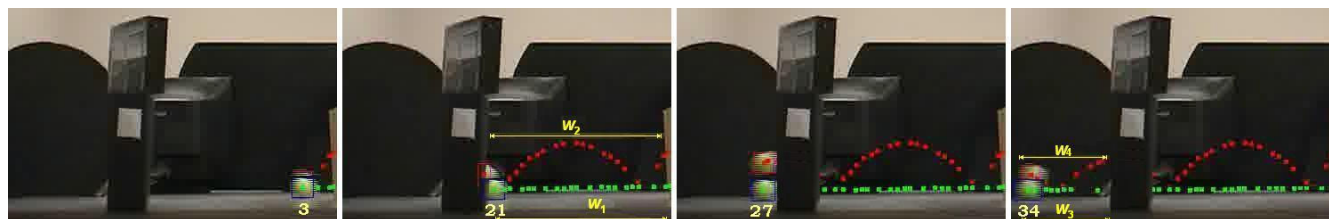


Fig. 5. Multi-target track matching.



Fig. 6. Track matching with partially overlapping tracks.