

Exploiting pedestrian interaction via global optimization and social behaviors

Laura Leal-Taixé, Gerard Pons-Moll, and Bodo Rosenhahn

Leibniz Universität Hannover, Appelstr. 9A, Hannover, Germany
{leal,pons,rosenhahn}@tnt.uni-hannover.de

Abstract. Multiple people tracking consists in detecting the subjects at each frame and matching these detections to obtain full trajectories. In semi-crowded environments, pedestrians often occlude each other, making tracking a challenging task. Tracking methods mostly work with the assumption that each pedestrian moves independently unaware of the objects or the other pedestrians around it. In the real world though, it is clear that when walking in a crowd, pedestrians try to avoid collisions, keep a close distance to a group of friends or avoid static obstacles in the scene.

In this paper, we present an approach which includes the interaction between pedestrians in two ways: first, including social and grouping behavior as a physical model within the tracking system, and second, using a global optimization scheme which takes into account all trajectories and all frames to solve the data association problem. Results are presented on three challenging publicly available datasets, showing our method outperforms state-of-the-art tracking systems. We also make a thorough analysis of the effect of the parameters of the proposed tracker as well as its robustness against noise, outliers and missing data.

1 Introduction

Multiple people tracking is a key problem for many computer vision tasks, such as surveillance, animation or activity recognition. In crowded environments occlusions and false detections are common, and although there have been substantial advances in the last years, tracking is still a challenging task. Tracking is often divided in two steps: detection, finding the objects of interest on every frame, and data association, matching the detections to form complete trajectories in time. Researchers have presented improvements on the object detector [1–3] as well as on the optimization techniques [4, 5] and even specific algorithms have been developed for tracking in crowded scenes [6, 7]. Though each object can be tracked separately, recent works have proven that tracking objects jointly and taking into consideration their interaction can give much better results in complex scenes. Current research is mainly focused on two aspects to exploit the interaction between pedestrians: the use of a global optimization strategy [8, 9] and a social motion model [10, 11]. The focus of this paper is to marry the concepts of global optimization and social and grouping behavior to obtain a robust tracker able to work in crowded scenarios. We extend the work presented in [12] to include more theoretical details, experimental results and details about the performance of the proposed method.

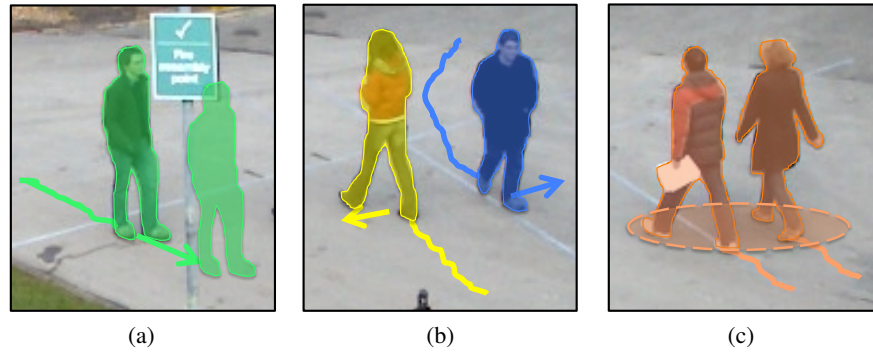


Fig. 1: Including social and grouping behavior to the network flow graph. (a) Constant velocity assumption. (b) Avoidance forces. (c) Group attraction forces.

1.1 Related work

The optimization strategy deals with the data association problem, which is usually solved on a frame-by-frame basis or one track at a time. Several methods can be used such as Markov Chain Monte Carlo (MCMC) [13], multi-level Hungarian [14], inference in Bayesian networks [15] or the Nash Equilibrium of game theory [16]. In [17] an efficient approximative Dynamic Programming (DP) scheme is presented, in which trajectories are estimated one after the other. This means that if a trajectory is formed using a certain detection, the other trajectories which are computed later will not be able to use that detection anymore. This obviously does not guarantee a global optimum for all trajectories. Recent works show that global optimization can be more reliable in crowded scenes as it solves the matching problem jointly for all tracks. The multiple object tracking problem is defined as a linear constrained optimization flow problem and Linear Programming (LP) is commonly used to find the global optimum. The idea was first used for people tracking in [18], although this method needs to know a priori the number of targets to track, which limits its application in real tracking situations. In [9], the scene is divided into identical cells, each represented by a node in the constructed graph. Using the information of the Probability Occupancy Map, the problem is formulated either as a max-flow and solved with Simplex, or as a min-cost and solved using k-shortest paths, which is a more efficient solution. Both methods show a far superior performance when compared to the same approach with DP [17]. The authors of [19] also define the problem as a maximum flow on an hexagonal grid, but instead of using matching individual detections, they make use of tracklets. This has the advantage that they can precompute the social forces for each of these tracklets, nonetheless, the fact that the tracklets are chosen locally, means the overall matching is not truly global, and if errors occur during the creation of the tracklets, these cannot be overcome by the global optimization. In [20], global and local methods are combined to match trajectories across cameras and across time. Finally, in [8] the tracking problem is formulated as a Maximum A-Posteriori (MAP) problem, which is mapped to a minimum-cost net-

work flow and then efficiently solved using LP. In this case, each node represents a detection, which means the graph is much smaller compared to [9, 19].

Most tracking systems work with the assumption that the motion model for each target is independent. This simplifying assumption is especially problematic in crowded scenes: imagine the chaos if every pedestrian followed his or her chosen path and completely ignored the other pedestrians in the scene. In order to avoid collisions and reach the chosen destination at the same time, a pedestrian follows a series of social rules or social forces. These have been defined in what is called the Social Force Model (SFM) [21], which has been used for abnormal crowd behavior detection [22], crowd simulation [23] and has only recently been applied to multiple people tracking: in [24], an energy minimization approach is used to predict the future position of each pedestrian considering all the terms of the social force model. In [10] and [25], the social forces are included in the motion model of the Kalman or Extended Kalman filter. In [26] a method is presented to detect small groups of people in a crowd, but it is only recently that grouping behavior has been included in a tracking framework [11, 27, 28]. In [28] groups are included in a graphical model which contains cycles and, therefore, Dual Decomposition [29] is needed to find the solution, which obviously is computationally much more expensive than using Linear Programming. Moreover, the results presented in [28] are only for short time windows. On the other hand, the formulations of [11, 27] are predictive by nature and therefore too local and unable to deal with trajectory changes (e.g. when people meet and stop to talk).

Social behavior models have only been introduced within a predictive framework, which are suboptimal due to the recursive nature of filtering. Therefore, in contrast to previous works, we propose to include social and grouping models into a global optimization framework which allows us to better estimate the true maximum a-posteriori probability of the trajectories.

1.2 Contributions

We present a novel approach for multiple people tracking which takes into account the interaction between pedestrians in two ways: first, using global optimization for data association and second, including social as well as grouping behavior. The key insight is that people plan their trajectories in advance in order to avoid collisions, therefore, a graph model which takes into account future and past frames is the perfect framework to include social and grouping behavior. We formulate multiple object tracking as a minimum-cost network flow problem, and present a new graph model which yields to better results than existing global optimization approaches. The social force model (SFM) and grouping behavior (GR) are included in an efficient way without altering the linearity of the problem. Results on several challenging public datasets show the improvement of the tracking results in crowded environments. Experiments with missing data, noise and outliers are also shown to test the robustness of the proposed approach. In this paper, we extend the work presented in [12] in three aspects : (i) more detailed theoretical explanations and background on Linear Programming for multiple object tracking; (ii) experimental results with different parameter values to see the effect of each of them on tracking results and (iii) detailed implementation details and computational aspects of the proposed method.

2 Multiple people tracking

Tracking is commonly divided in two steps: object detection and data association. First, the objects are detected in each frame of the sequence and second, the detections are matched to form complete trajectories. In this section we define the data association problem and describe how to convert it to a minimum-cost network flow problem, which can be efficiently solved using Linear Programming.

The idea is to build a graph in which the nodes represent the pedestrian detections. These nodes are fully connected to past and future observations by edges, which determine the relation between two observations with a cost. Thereby, the matching problem is equivalent to a minimum-cost network flow problem: finding the optimal set of trajectories is equivalent to sending flow through the graph so as to minimize the cost. This can be efficiently computed using the Simplex algorithm or k-shortest paths [30].

2.1 Problem statement

Let $\mathcal{O} = \{\mathbf{o}_k^t\}$ be a set of object detections with $\mathbf{o}_k^t = (\mathbf{p}_k, t)$, where $\mathbf{p}_k = (x, y, z)$ is the 3D position and t is the time stamp. A trajectory is defined as a list of ordered object detections $T_k = \{\mathbf{o}_k^1, \mathbf{o}_k^2, \dots, \mathbf{o}_k^N\}$, and the goal of multiple object tracking is to find the set of trajectories $\mathcal{T}^* = \{T_k\}$ that best explains the detections. This is equivalent to maximizing the a-posteriori probability of \mathcal{T} given the set of detections \mathcal{O} . Assuming detections are conditionally independent, the objective function is expressed as:

$$\mathcal{T}^* = \underset{\mathcal{T}}{\operatorname{argmax}} P(\mathcal{T}|\mathcal{O}) = \underset{\mathcal{T}}{\operatorname{argmax}} \prod_k P(\mathbf{o}_k|\mathcal{T})P(\mathcal{T}) \quad (1)$$

$P(\mathbf{o}_k|\mathcal{T})$ is the likelihood of the detection. In order to reduce the space of \mathcal{T} , we make the assumption that the trajectories cannot overlap (i.e., a detection cannot belong to two trajectories), but unlike [8], we do not define the motion of each subject to be independent, therefore, we deal with a much larger search space. We extend this space by including the following dependencies for each trajectory T_k :

- Constant velocity assumption: the observation $\mathbf{o}_k^t \in T_k$ depends on past observations $[\mathbf{o}_k^{t-1}, \mathbf{o}_k^{t-2}]$
- Grouping behavior: If T_k belongs to a group, the set of members of the group $\mathcal{T}_{k,\text{GR}}$ has an influence on T_k
- Avoidance term: T_k is affected by the set of trajectories $\mathcal{T}_{k,\text{SFM}}$ which are close to T_k at some point in time and do not belong to the same group as T_k

The first and third dependencies are grouped into the SFM term. The sets $\mathcal{T}_{k,\text{SFM}}$ and $\mathcal{T}_{k,\text{GR}}$ are disjoint, i.e., a pedestrian can have an attractive effect or a repulsive effect on another pedestrian, but not both. Therefore, we can assume that these two terms are independent and decompose $P(\mathcal{T})$ as:

$$\begin{aligned} P(\mathcal{T}) &= \prod_{T_k \in \mathcal{T}} P(T_k \cap \mathcal{T}_{k,\text{SFM}} \cap \mathcal{T}_{k,\text{GR}}) \\ &= \prod_{T_k \in \mathcal{T}} P(\mathcal{T}_{k,\text{SFM}}|T_k)P(\mathcal{T}_{k,\text{GR}}|T_k)P(T_k) \end{aligned} \quad (2)$$

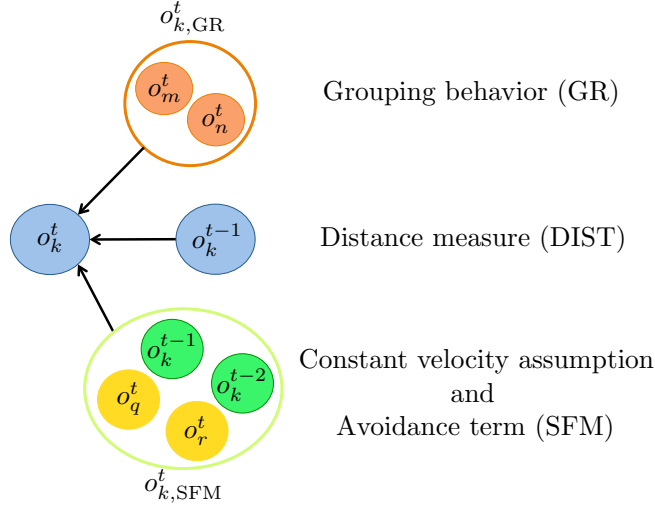


Fig. 2: Diagram of the dependencies for each observation \mathbf{o}_k^t .

where the trajectories are represented by a Markov chain:

$$\begin{aligned}
 P(\mathcal{T}) = \prod_{T_k \in \mathcal{T}} & P_{\text{in}}(\mathbf{o}_k^1) \dots P(\mathbf{o}_k^t | \mathbf{o}_k^{t-1}) \\
 & P_{k,\text{SFM}}(\mathbf{o}_k^t | \mathbf{o}_k^t, \mathbf{o}_k^{t-1}) P_{k,\text{GR}}(\mathbf{o}_k^t | \mathbf{o}_k^t, \mathbf{o}_k^{t-1}) \\
 & \dots P_{\text{out}}(\mathbf{o}_k^N)
 \end{aligned} \tag{3}$$

where $P_{\text{in}}(\mathbf{o}_k^t)$ is the probability that a trajectory is initiated with detection \mathbf{o}_k^t , $P_{\text{out}}(\mathbf{o}_k^t)$ the probability that the trajectory is terminated at \mathbf{o}_k^t and $P(\mathbf{o}_k^t | \mathbf{o}_k^{t-1})$ is the probability that \mathbf{o}_k^{t-1} is followed by \mathbf{o}_k^t in the trajectory. $P_{k,\text{SFM}}$ evaluates how well the social rules are kept if \mathbf{o}_k^t is matched to \mathbf{o}_k^{t-1} , and $P_{k,\text{GR}}$ describes how well the structure of the group is kept.

Let us assume that we are analyzing observation \mathbf{o}_k^t . In Figure 2 we summarize which observations influence the matching of \mathbf{o}_k^t . Typical approaches [8] only take into account distance (DIST) information, that is, the observation in the previous frame \mathbf{o}_k^{t-1} . We introduce the social dependencies (SFM) given by the constant velocity assumption (green nodes) and the avoidance term (yellow nodes). In this case, two observations, \mathbf{o}_q^t and \mathbf{o}_r^t that do not belong to the same group as \mathbf{o}_k^t , will be considered to create a repulsion effect on \mathbf{o}_k^t . On the other hand, the orange nodes which depict the grouping term (GR), are two other observations \mathbf{o}_m^t and \mathbf{o}_n^t which do belong to the same group as \mathbf{o}_k^t and therefore have an attraction effect on \mathbf{o}_k^t . Note that all these dependencies can only be modeled by high order terms, which means that either we use complex solvers [28] to find a solution in graphs with cycles, or we keep the linearity of the problem by using an iterative approach as we explain later on.

2.2 Tracking with Linear Programming

We linearize the objective function by defining a set of flow flags $f_{i,j} = \{0, 1\}$ which indicate if an edge (i, j) is in the path of a trajectory or not. In a minimum cost network flow problem, the objective is to find the values of the variables that minimize the total cost of the flows over the network. Defining the costs as negative log-likelihoods, and combining Equations (1), (2) and (3), the following objective function is obtained:

$$\begin{aligned} \mathcal{T}^* &= \underset{\mathcal{T}}{\operatorname{argmin}} \sum_{T_k \in \mathcal{T}} -\log P(T_k) - \log P(\mathcal{T}_{\text{SFM}}|T_k) \\ &\quad - \log P(\mathcal{T}_{\text{GR}}|T_k) + \sum_k -\log P(\mathbf{o}_k|\mathcal{T}) \\ &= \underset{\mathcal{T}}{\operatorname{argmin}} \sum_i C_{\text{in},i} f_{\text{in},i} + \sum_i C_{i,\text{out}} f_{i,\text{out}} \\ &\quad + \sum_{i,j} (C_{i,j} + C_{\text{SFM},i,j} + C_{\text{GR},i,j}) f_{i,j} + \sum_i C_i f_i \end{aligned} \quad (4)$$

subject to the following constraints:

- Edge capacities: we assume that each detection can only correspond to one trajectory, therefore, the edge capacities have an upper bound of $u_{ij} \leq 1$ and:

$$f_{\text{in},i} + f_i \leq 1 \quad f_{i,\text{out}} + f_i \leq 1 \quad (5)$$

- Flow conservation at the nodes:

$$f_{\text{in},i} + f_i = \sum_j f_{i,j} \quad \sum_j f_{j,i} = f_{i,\text{out}} + f_i \quad (6)$$

- Exclusion property:

$$f_{i,j} = \{0, 1\} \quad (7)$$

The condition in Eq. 7 requires us to solve an integer program, which is known to be NP-complete. Nonetheless, we can relax the condition to have the following linear equation:

$$0 \leq f_{i,j} \leq 1. \quad (8)$$

Now the problem is defined and can be solved as a linear program. If certain conditions are fulfilled, the solution \mathcal{T}^* will still be integer, and therefore will also be the optimal solution to the initial integer program. We discuss the integrality of the solution in more detail in Section 4.

To map this formulation into a cost-flow network, we define $G = (N, E)$ to be a directed network with a cost $C_{i,j}$ and a capacity u_{ij} associated with every edge $(i, j) \in E$. An example of such a network is shown in Figure 3; it contains two special nodes, the source s and the sink t ; all flow that goes through the graph starts at the s node and ends at the t node. Thereby, each flow represents a trajectory T_k and the path that each flow follows indicates which observations belong to each of the trajectories. Each

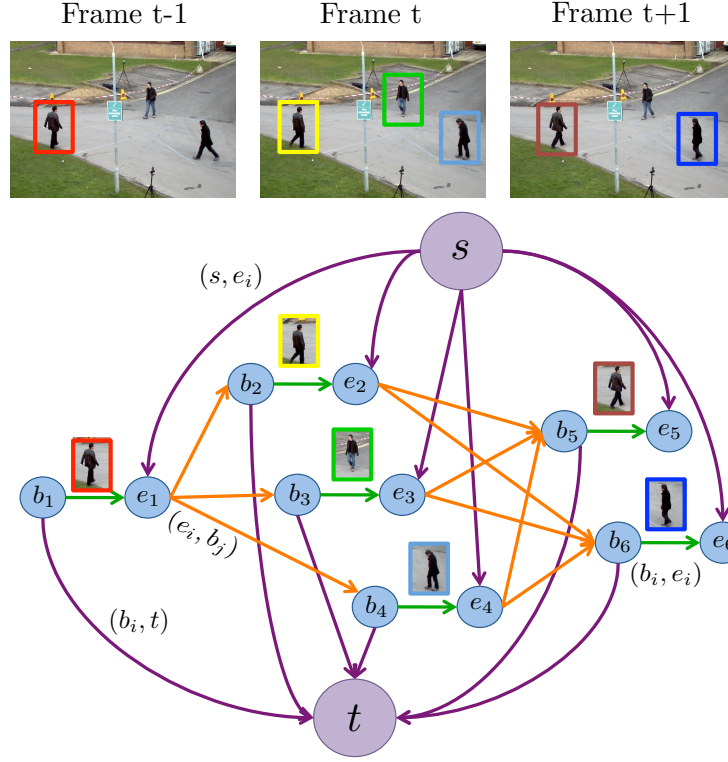


Fig. 3: Example of a graph with the special source s and sink t nodes, 6 detections which are represented by two nodes each: the beginning b_i and the end e_i .

observation \mathbf{o}_i is represented with two nodes, the beginning node $b_i \in N$ and the end node $e_i \in N$ (see Figure 3). A detection edge connects b_i and e_i .

Below we detail the three types of edges present in the graphical model and the cost for each type:

Link edges. The edges (e_i, b_j) connect the end nodes e_i with the beginning nodes b_j in following frames, with cost $C_{i,j}$ and flow $f_{i,j}$, defined as:

$$f_{i,j} = \begin{cases} 1, & \mathbf{o}_i \text{ and } \mathbf{o}_j \text{ belong to } T_k \text{ and } \Delta f \leq F_{\max} \\ 0, & \text{otherwise} \end{cases} \quad (9)$$

where Δf is the frame number difference between nodes j and i and F_{\max} is the maximum allowed frame gap.

The costs of the link edges represent the spatial relation between different subjects. Assuming that a subject cannot move a lot from one frame to the next, we define the costs to be a decreasing function of the distance between detections in successive

frames. The time gap between observations is also taken into account in order to be able to work at any frame rate, therefore velocity measures are used instead of distances. The velocities are mapped to probabilities with a Gauss error function as shown in Equation (10), assuming the pedestrians cannot exceed a maximum velocity V_{\max} . The effect of parameter V_{\max} is detailed in Section 5.1.

$$E(V_t, V_{\max}) = \frac{1}{2} + \frac{1}{2} \operatorname{erf} \left(\frac{-V_t + \frac{V_{\max}}{2}}{\frac{V_{\max}}{4}} \right) \quad (10)$$

As we can see in Figure 4, the advantage of using Equation (10) over a linear function is that the probability of lower velocities decreases more slowly, while the probability for higher velocities decreases more rapidly. This is consistent with the probability distribution of speed learned from training data.

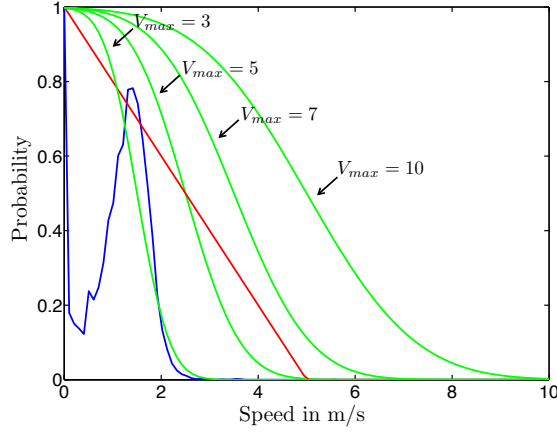


Fig. 4: *Blue* = normalized histogram of speeds learned from training data. *Red* = probability distribution if cost depends linearly on the velocity. *Green* = probability distribution if the relation of cost and velocities is expressed by Equation (10). An $V_{\max} = 7\text{m/s}$ is used in the experiments.

Therefore, the cost of a link edge is defined as:

$$\begin{aligned} C_{i,j} &= -\log(P(\mathbf{o}_j|\mathbf{o}_i)) + C(\Delta f) \\ &= -\log E \left(\frac{\|\mathbf{p}_j - \mathbf{p}_i\|}{\Delta t}, V_{\max} \right) + C(\Delta f) \end{aligned} \quad (11)$$

where $C(\Delta f) = -\log(B_j^{\Delta f-1})$ is the cost depending on the frame difference between detections.

Detection edges. The edges (b_i, e_i) connect the beginning node b_i and end node e_i , with cost C_i and flow f_i , defined as:

$$f_i = \begin{cases} 1, & \mathbf{o}_i \text{ belongs to } T_k \\ 0, & \text{otherwise} \end{cases} \quad (12)$$

If all the costs of the edges are positive, the solution to the minimum-cost problem is the trivial null flow. Consequently, we represent each observation with two nodes and a detection edge with negative cost:

$$C_i = \log(1 - P_{det}(\mathbf{o}_i)) + \log\left(\frac{BB_{\min}}{\|\mathbf{p}_{BB} - \mathbf{p}_i\|}\right). \quad (13)$$

The higher the likelihood of a detection $P_{det}(\mathbf{o}_i)$ the more negative the cost of the detection edge, hence, confident detections are likely to be in the path of the flow in order to minimize the total cost. If a map of the scene is available, we can also include this information in the detection cost. If a detection is far away from a possible entry/exit point, we add an extra negative cost to the detection edge, in order to favor that observation to be matched. The added cost depends on the distance to the closest entry/exit point \mathbf{p}_{BB} , and is only computed for distances higher than $BB_{\min} = 1.5m$. This is a probabilistic simple way of including other information present in the scene, such as obstacles or attraction points (shops, doors, etc).

Entrance and exit edges. The edges (s, e_i) connect the source s with all the end nodes e_i , with cost $C_{in,i}$ and flow $f_{in,i}$. Similarly, (b_i, t) connects the end node b_i with sink t , with cost $C_{i,out}$ and flow $f_{i,out}$. The flows are defined as:

$$f_{in,i} \text{ (or } f_{i,out}) = \begin{cases} 1, & T_k \text{ starts (or ends) at } \mathbf{o}_i \\ 0, & \text{otherwise} \end{cases} \quad (14)$$

In [8], the authors propose to create the opposite edges (s, b_i) and (e_i, t) , which means tracks entering and leaving the scene go through the detection node and therefore benefiting from its negative cost (see Figure 5(a)). If the costs C_{in} and C_{out} are then set to zero, a track will be started at each detection of each frame, because it will be cheaper to use the entrance and exit edges than the link edges. On the other hand, if C_{in} and C_{out} are very high, it will be hard for the graph to create any trajectory. Therefore, the choice of these two costs is extremely important. In [8], the costs are set according to the entrance and exit probabilities P_{in} and P_{out} , which are data dependent terms that need to be calculated during optimization.

In contrast, we propose to connect the s node with the end nodes and the t node to the begin nodes (as shown in Figure 5(b)). This way, we make sure that when a track starts (or ends) it does not benefit from the negative cost of the detection edge. Setting $C_{in} = C_{out} = 0$ and taking into account the flow constraints of Eqs. (5) and (6), we make sure the trajectories are only created with the information of the link edges.

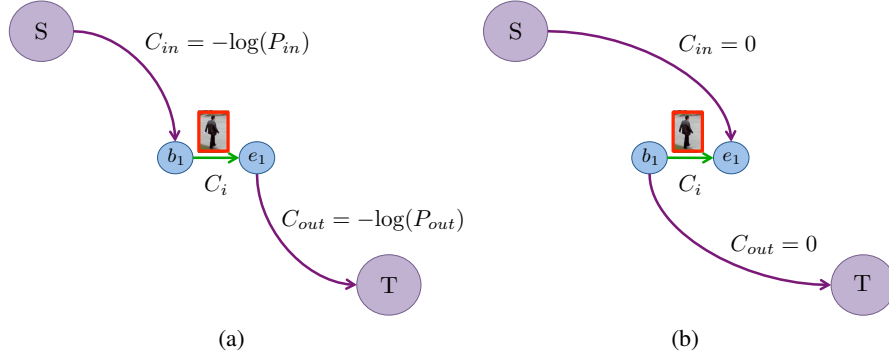


Fig. 5: (a) Graph structure as used in [8], which requires the computation of P_{in} and P_{out} in an Expectation-Maximization step during optimization. In contrast, the proposed graph structure in (b) allows us to get rid of these two extra parameters. The trajectories are found only with the information of the link and detection edges.

3 Modeling social behavior

If a pedestrian does not encounter any obstacles, the natural path to follow is a straight line. But what happens when the space gets more and more crowded and the pedestrian can no longer follow the straight path? Social interaction between pedestrians is especially important when the environment is crowded. In this section we consider how to include the social behavior [21], which we divide into the Social Force Model (SFM) and the Group behavior (GR), into our minimum-cost network flow problem.

3.1 Social Force Model

The social force model states that the motion of a pedestrian can be described as if they were subject to "social forces". There are three main terms that need to be considered: the desire of a pedestrian to maintain a certain speed, the desire to keep a comfortable distance from other pedestrians and the desire to reach a destination. Since we cannot know a priori the destination of the pedestrian in a real tracking system, we focus on the first two terms.

Constant velocity assumption. The pedestrian tries to keep a certain speed and direction, therefore we assume that in $t + \Delta t$ we have the same speed as in t and predict the pedestrian's position in $t + \Delta t$ accordingly.

$$\tilde{\mathbf{p}}_i^{t+\Delta t} = \mathbf{p}_i^t + \mathbf{v}_i^t \Delta t$$

Avoidance term. The pedestrian also tries to avoid collisions and keep a comfortable distance from other pedestrians. We model this term as a repulsion field with an exponential distance-decay function with value α learned from training data.

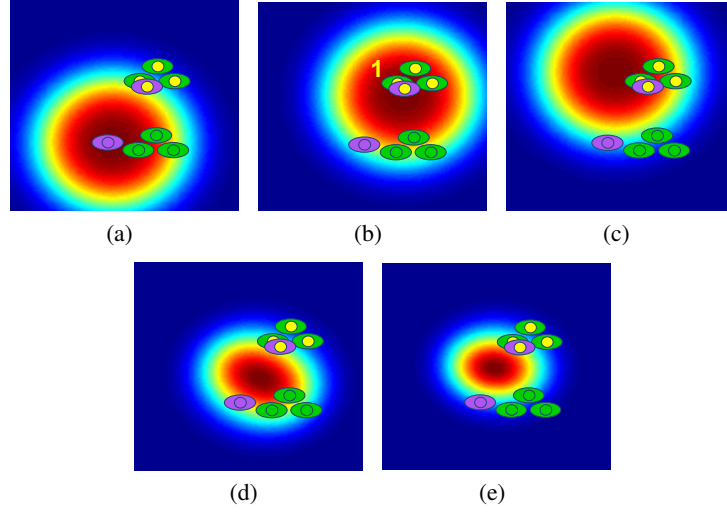


Fig. 6: Three green pedestrians walk in a group, the predicted positions in the next frame are marked by yellow heads. The purple pedestrian's linearly predicted position (yellow head) clearly interferes with the trajectory of the group. Representation of the probability (blue is 0 red is 1) distribution for the purple's next position using: 6(a) only distances, 6(b) only SFM (constant velocity assumption and avoidance term), 6(c) only GR (considering the purple pedestrian belongs to the group), 6(d) distances+SFM and 6(e) distances+SFM+GR.

$$\mathbf{a}_i^{t+\Delta t} = \sum_{g_m \neq g_i} \exp \left(-\frac{\|\tilde{\mathbf{p}}_i^{t+\Delta t} - \tilde{\mathbf{p}}_m^{t+\Delta t}\|}{\alpha \Delta t} \right) \quad (15)$$

If we are computing the cost of edge (i, j) , we use the constant velocity assumption to predict the position of \mathbf{o}_i and \mathbf{o}_j as well as the rest of pedestrians $\tilde{\mathbf{p}}_m^{t+\Delta t}$, and compute the repulsion acceleration each pedestrian has on i . The only pedestrians that have this repulsion effect on subject i are the ones which do not belong to the same group as i and $\|\tilde{\mathbf{p}}_i^{t+\Delta t} - \tilde{\mathbf{p}}_m^{t+\Delta t}\| \leq 1m$. The different avoidance terms are combined linearly.

Now the prediction of the pedestrian's next position is also influenced by the avoidance term (acceleration) from all pedestrians:

$$\tilde{\mathbf{p}}_i^{t+\Delta t} = \mathbf{p}_i^t + (\mathbf{v}_i^t + \mathbf{a}_i^{t+\Delta t} \Delta t) \Delta t \quad (16)$$

The distance between prediction and real measurements is used to compute the cost:

$$C_{\text{SFM},i,j} = -\log E \left(\frac{\|\tilde{\mathbf{p}}_i^{t+\Delta t} - \mathbf{p}_j^{t+\Delta t}\|}{\Delta t}, V_{\max} \right) \quad (17)$$

where the function E is detailed in Eq. (10).

In Figure 6 we plot the probability distribution computed using different terms. Note, this is just for visualization purposes, since we do not compute the probability for each point on the scene, but only for the positions where the detector has fired. There are 4 pedestrians in the scene, the purple one and 3 green ones walking in a group. As shown in 6(b), if we only use the predicted positions (yellow heads) given the previous speeds, there is a collision between the purple pedestrian and the green marked with a 1 collide. The avoidance term shifts the probability mode to a more plausible position.

3.2 Group Model

The social behavior [21] also includes an attraction force which occurs when a pedestrian is attracted to a friend, shop, etc. We model the attraction between members of a group. Before modeling group behavior we determine which tracks form each group and at which frame the group begins and ends (to deal with splitting and formation of groups). The idea is that if two pedestrians are close to each other over a reasonable period of time, they are likely to belong to the same group. From the training sequence in [10], we learn the distance and speed probability distributions of the members of a group P_g vs. individual pedestrians P_i . If m and n are two trajectories which appear on the scene at $t = [0, N]$, we compute the flag $G_{m,n}$ that indicates if m and n belong to the same group.

$$G_{m,n} = \begin{cases} 1, & \sum_{t=0}^N P_g(m, n) > \sum_{t=0}^N P_i(m, n) \\ 0, & \text{otherwise} \end{cases} \quad (18)$$

For every observation \mathbf{o}_i , we define a group label g_i which indicates to which group the observation belongs to, if any. If several pedestrians form a group, they tend to keep a similar speed, therefore, if i belongs to a group, we can use the mean speed of all the other members of the group to predict the next position for i :

$$\tilde{\mathbf{p}}_i^{t+\Delta t} = \mathbf{p}_i^t + \sum_{g_m=g_i} \mathbf{v}_m^t \Delta t \quad (19)$$

The distance between this predicted position and the real measurements is used in (10) to obtain the cost for the grouping term.

An example is shown in Figure 6(c), where we can see that the maximum probability provided by the group term keeps the group configuration. In Figure 6(d) we show the combined probability of the distance and SFM information, which narrows the space of probable positions. Finally, Figure 6(e) represents the combined probability of DIST, SFM and GR. As we can see, the space of possible locations for the purple pedestrian is considerably reduced as we add the social and grouping behaviors, which means we have less ambiguities for data association. This is specially useful to decrease identity switches as we present in Section 5.

4 Implementation details

To compute the SFM and grouping costs, we need to have information about the velocities of the pedestrians, which can only be obtained if we already have the trajectories. We solve this chicken-and-egg problem iteratively as shown in Algorithm 1; on the first iteration, the trajectories are estimated only with the information defined in Section 2.2, for the rest of iterations, the SFM and GR is also used. The algorithm stops when the trajectories do not change or when a maximum number of iterations M_i is reached.

Algorithm 1 Iterative optimization

```

while  $\mathcal{T}_i \neq \mathcal{T}_{i-1}$  and  $i \leq M_i$  do
  if  $i == 1$  then
    1.1. Create the graph using only DIST information
  else
    1.2. Create the graph using DIST, SFM and GR information
  end if
  2. Solve the graph to find  $\mathcal{T}_i$ 
  3. Compute velocities and groups given  $\mathcal{T}_i$ 
end while

```

Linear Programming solvers The minimum cost solution is found using the Simplex algorithm [30], with the implementation given in [31]. Though Simplex has an exponential worst-case complexity, we are able to track most sequences in just a few seconds; this is because each node represents one detection, and therefore the dimension of the graph is quite small. For larger graphs [9] or more crowded environments, we can use the k-shortest paths solver [9, 32] which has a worst case complexity of $O(k(m + n \cdot \log(n)))$. For more details on network flows and Simplex we refer the reader to [33], and to [34] for more information on the k-shortest path algorithm.

Integrality of the solution When defining the program to be solved, we saw that Eq. (7) defined an integer program, which is known to be NP-complete. We relaxed the condition into Eq. (8) in order to use efficient Linear Programming solvers to find the optimum solution to our problem. If the solution to the relaxed version of the program is integer, then we know it is an optimal solution of the original problem [33]. The question is, can we guarantee that the solution will be always integer?

Let us assume the conditions of the Linear Program are expressed as: $Ax = b$. If all entries of A and b are integer, as it is our case, we can determine that $Ax = b$ has an integer solution by Cramer's rule:

$$Ax = b \Leftrightarrow x = A^{-1}b \Leftrightarrow \forall i : x_i = \frac{\det(A^i)}{\det(A)} \quad (20)$$

where A^i is equal to A except on the i -th column where it is equal to b . From here, we can determine that x will be integer when $\det(A)$ is equal to $+1$ or -1 . A matrix $A \in \mathbb{Z}^{m \times n}$ is *totally unimodular* if the determinant of all the sub-square matrices of A is either 0, $+1$ or -1 .

Theorem 1: If A is totally unimodular, every vertex solution of $Ax \leq b$ is integer.

A well-known case of totally unimodular matrices are the node arc incidence matrices N of a directed network. Therefore, our defined constraint matrix is totally unimodular, and the solutions we will obtain will always be integer.

Computationally reduction To reduce the computational cost, we prune the graph using the physical constraints represented by the edge costs. If any of the costs C_{ij} , $C_{\text{SFM},i,j}$ or $C_{\text{GR},i,j}$ is infinite, the two detections i and j are either too far away to belong to the same trajectory or they do not match according to social and grouping rules, therefore the edge (i, j) is erased from the graphical model. For long sequences, we divide the video into several batches and optimize for each batch. For temporal consistency, the batches have an overlap of $F_{\text{max}} = 10$ frames. With our non-optimized code, the runtime for a sequence of 800 frames (114 seconds), 4837 detections, batches of 100 frames and 6 iterations is 30 seconds on a 3GHz machine.

5 Experimental results

In this section we show the tracking results of our method on three publicly available datasets and compare with existing state-of-the-art tracking approaches using the CLEAR metrics [35], which split the measuring scores into *accuracy* and *precision*:

- **Detection Accuracy (DA):** measures how many detections were correctly found and therefore is based on the count of missed detections m_t and false alarms f_t for each frame t .

$$DA = 1 - \frac{\sum_{t=1}^{N_f} m_t + f_t}{\sum_{t=1}^{N_f} N_G^t}$$

where N_f is the number of frames of the sequence and N_G^t is the number of ground truth detections in frame t . A detection is considered to be correct when it is found within 50 pixels from the ground truth and the bounding boxes of both ground truth and detection have some overlap.

- **Tracking Accuracy (TA):** similar to DA but also including the identity switches i_t . In this case, the measure does not penalize identity switches as much as a missing detection or a false alarm as we use a \log_{10} weight.

$$DA = 1 - \frac{\sum_{t=1}^{N_f} m_t + f_t + \log_{10}(1 + i_t)}{\sum_{t=1}^{N_f} N_G^t}$$

- **Detection Precision (DP):** precision measurements represent how well the bounding box detections match the ground truth. For this, an overlap measure between bounding boxes is used:

$$Ov^t = \sum_{i=1}^{N_{\text{mapped}}^t} \frac{|G_i^t \cap D_i^t|}{|G_i^t \cup D_i^t|}$$

where N_{mapped}^t is the number of mapped objects in frame t , i.e., the number of detections that are matched to some ground truth object. G_i^t is the i th ground truth object of frame t and D_i^t the detected object matched to G_i^t . The DP measure is then expressed as:

$$DP = \frac{\sum_{t=1}^{N_f} \frac{Ov^t}{N_{\text{mapped}}^t}}{N_f}$$

- **Tracking Precision (TP):** measures the spatiotemporal overlap between ground truth trajectories and detected ones, taking into account also split and merged trajectories.

$$TP = \frac{\sum_{i=1}^{N_{\text{mapped}}^t} \sum_{t=1}^{N_f} \frac{|G_i^t \cap D_i^t|}{|G_i^t \cup D_i^t|}}{\sum_{t=1}^{N_f} N_{\text{mapped}}^t}$$

All experiments except the ones in Section 5.1 are performed with 6 iterations, a batch of 100 frames, $V_{\text{max}} = 7m/s$, $F_{\text{max}} = 10$, $\alpha = 0.5$ and $B_j = 0.3$.

5.1 Analysis of the effect of the parameters

All parameters defined in previous sections are learned from training data; in our case we use one sequence of the publicly available dataset [10]. In this section we study the effect of the few parameters needed in our implementation, and show the proposed graph works well for a wide range of these parameters and therefore no parameter tuning is needed to obtain a good performance. The analysis is done on two publicly available datasets: a crowded town center [36] and the well-known PETS2009 dataset [37], to see the different effects of each parameters on each dataset.

Number of iterations. The first parameter we analyze is the number of iterations M_i that we allow. This determines how many times the loop between computing social forces and computing trajectories is performed as explained in Algorithm 1. Looking

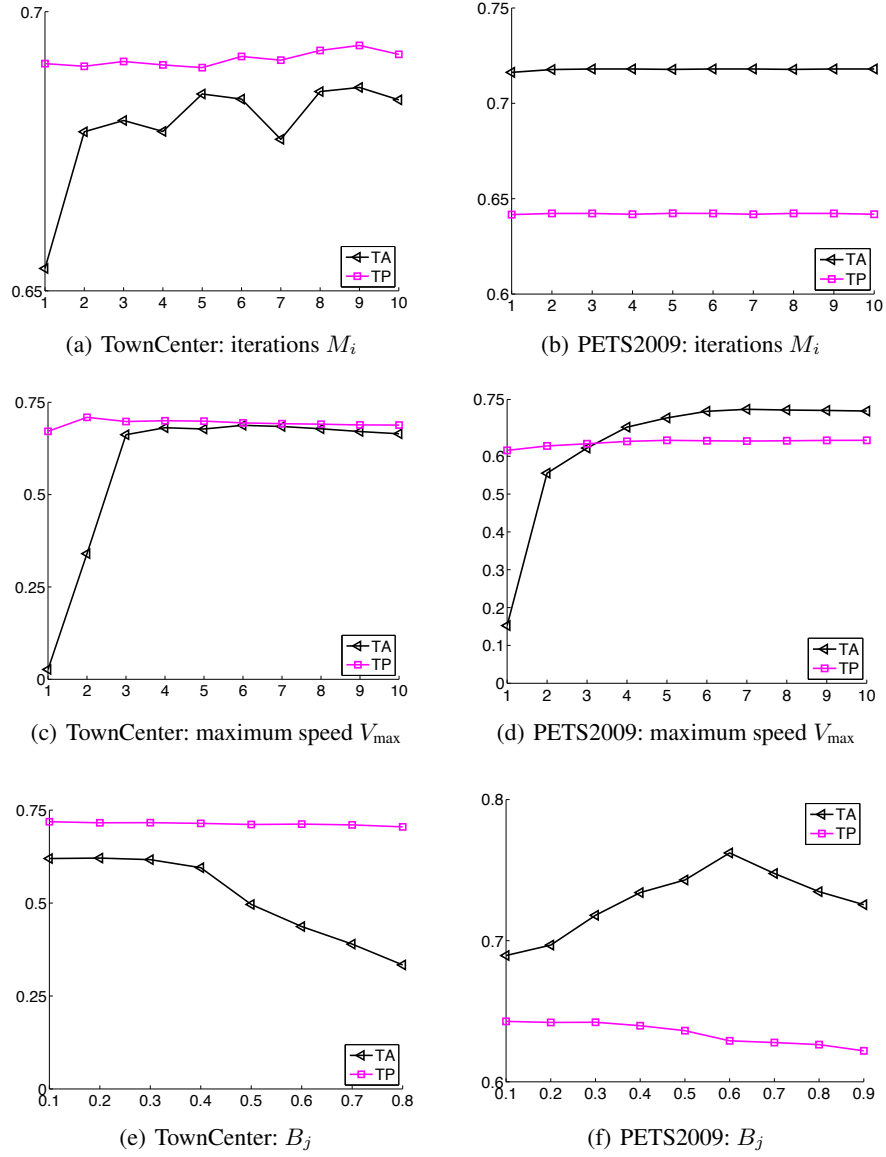


Fig. 7: Tracking accuracy (black) and precision (magenta) obtained for the Town Center dataset (left column) and the PETS 2009 dataset (right column) given varying parameter values.

at the results on the PETS 2009 dataset in Figure 7(b), we can see that after just 2 iterations the results remain very stable. Actually, the algorithm reports no changes in the trajectories after 3 iterations, and therefore stops even though the maximum number of iterations allowed is higher. The result with 1 and 2 iterations is also not very different, which means the social and grouping behavior do not significantly improve the results for this particular dataset. This is due to the fact that this dataset is very challenging from a social behavior point of view, with subjects often changing direction and groups forming and splitting frequently. More details and comments on these results can be found in Section 5.3. On the other hand, we observe a different effect on the TownCenter dataset, shown in Figure 7(a). In this case, there is a clear improvement when using social and grouping behavior (i.e. the result improves when we use more than one iteration). We also observe a pattern on how the Tracking Accuracy of the dataset evolves: there is a cycle of 3 iterations for which the accuracy increases and decreases in a similar pattern. This means that the algorithm is jumping between two solutions and will not converge to neither one of them. This happens when pedestrians are close together for a long period of time but are not forming a group, which means that even with social forces, it is hard to say which paths they will follow.

Maximum speed. This is the parameter that determines the maximum speed of the pedestrians that we are observing. In this case, we can see in Figures 7(c) and 7(d) a clear trend in which the results are very bad when we force the pedestrians to walk more slowly than they actually do, since we are artificially splitting trajectories. The results converge when the maximum speed allowed is around 3m/s - 5m/s, which is the reported mean speed of pedestrians in a normal situation. More interestingly, we observe that the results are kept constant when using higher maximum speed values. This is a positive effect of the global optimization framework, since we can use a much higher speed limit and this will still give us good results and will allow us to track a person running through the scene, a case of panic when people start running, etc.

Cost for the frame difference. The last parameter, B_j , appears in Eq. (12) and represents the penalization term that we apply when the frame difference between two detections that we want to match is larger than 1. This term is used in order to give preference to matches that are close in time. Here we can again see different effects on the two datasets. In Figure 7(e), we see that the results are stable until a value of 0.4. The lower the value, the higher is the penalization cost for the frame difference, which means it is more difficult to match those detections which are more than 1 frame apart. When the value of B_j is higher than 0.4, there are more ambiguities in the data association process because it is easier to match detections which are many frames apart. In the TownCenter dataset, there is no occluding object in the scene, which means missing detections are sporadic within a given trajectory. In this scenario, a lower value for B_j is better, since small gaps can be filled and there are less ambiguities. Nonetheless, we see different results in the PETS 2009 dataset in Figure 7(f), since here there is a clear occluding object in the middle of the scene (see Figure 8) which occludes the pedestrians for longer periods of time. In this case, a higher value of B_j allows to overcome these large gaps of missing data, and that is why the best value for this dataset is around 0.6.

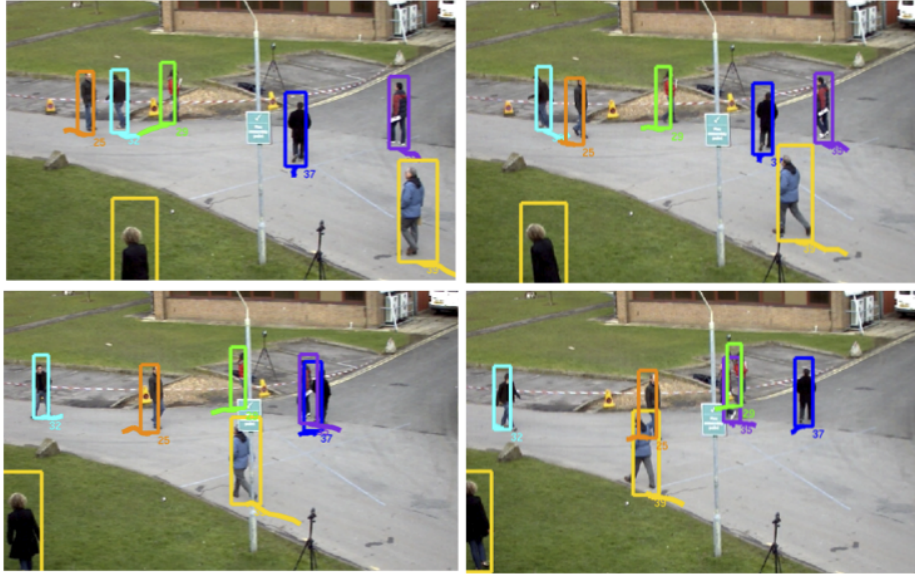


Fig. 8: Four frames of the PETS2009 sequence (separation of 9 frames), showing several occlusions, both created by the obstacle on the scene and between pedestrians. All the occlusions can be recovered with the proposed method.

5.2 Evaluation with missing data, noise and outliers

We evaluate the impact of every component of the proposed approach with one of the sequences of the dataset [10], which contains images from a crowded public place, with several groups as well as walking and standing pedestrians. The sequence is 11601 frames long and contains more than 300 trajectories. First of all, we evaluate our group detection method on the whole sequence with ground truth detections: 61% are correctly detected, 26% are only partially detected, 13% are not found and an extra 7% groups are detected wrongly.

Using the ground truth (GT) pedestrian positions as the baseline for our experiments, we perform three types of tests, missing data, outliers and noise, and compare the results obtained with:

- DIST: proposed network model with distances
- SFM: adding the Social Force Model (Section 3.1)
- SFM+GR: adding SFM and grouping behavior (Section 3.2)

Missing data. This experiment shows the robustness of our approach given missed detections. This is evaluated by randomly erasing a certain percentage of detections from the GT set. The percentages evaluated are $[0, 4, 8, 12, 16, 20]$ from the total number of detections over the whole sequence. As we can see in Figure 10, both SFM and SFM+GR increase the tracking accuracy when compared to DIST.

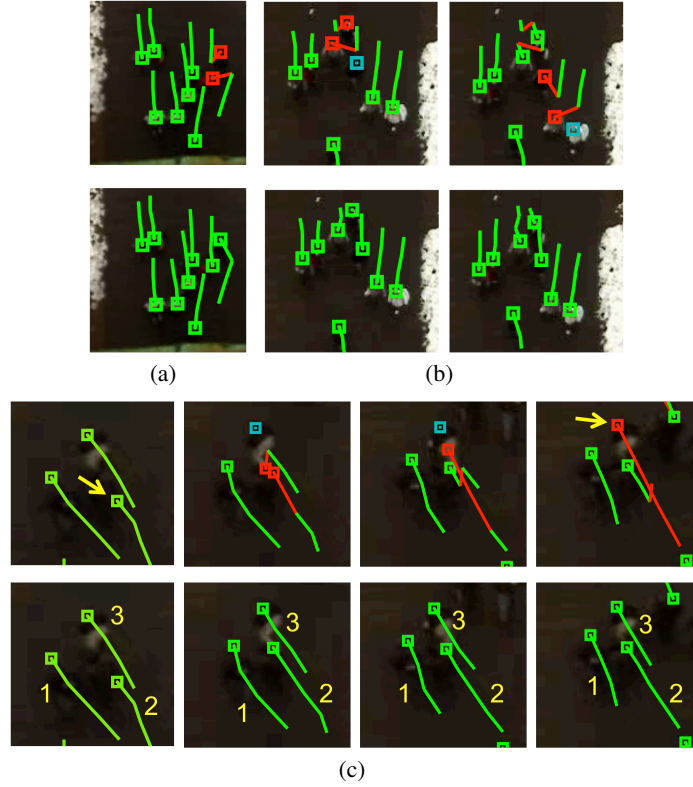


Fig. 9: *Top row*: Tracking results with only DIST. *Bottom row*: Tracking results with SFM+GR. *Green* = correct trajectories, *Blue* = observation missing from the set, *Red* = wrong match. 9(a) Wrong match with DIST, corrected with SFM. 9(b) Missing detections cause the matches to shift due the global optimization; correct result with SFM. 9(c) Missed detection for subject 3 on two consecutive frames. With SFM, subject 2 in the first frame (yellow arrow) is matched to subject 3 in the last frame (yellow arrow), creating an identity switch; correct result with grouping information.

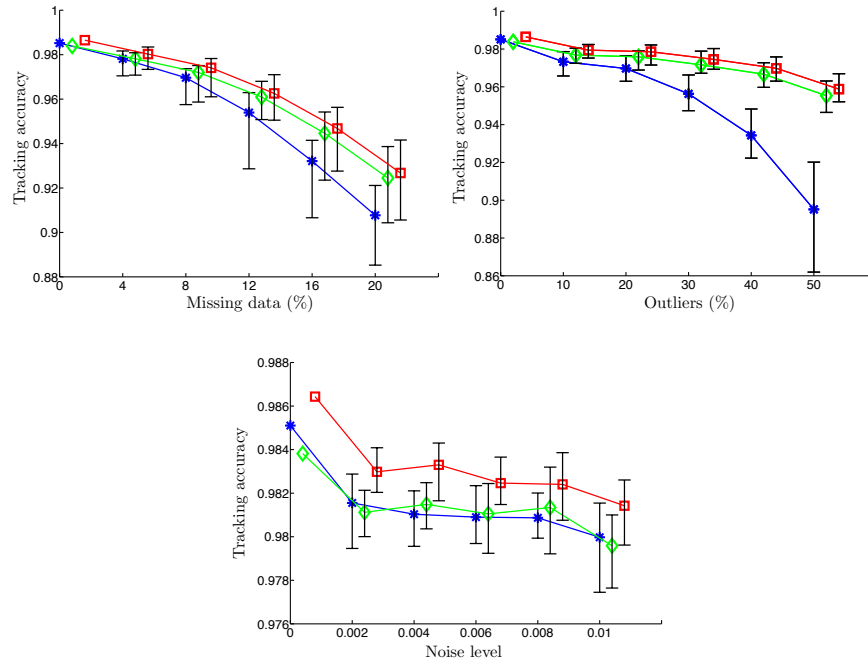


Fig. 10: Experiments are repeated 50 times and average result, maximum and minimum are plotted. *Blue star* = results with DIST, *Green diamond* = results with SFM, *Red square* = results with SFM+GR. *From left to right:* Experiment with simulated missing data, with outliers, and with random noise.

Outliers. With an initial set of detections of GT with 2% missing data, tests are performed with $[0, 10, 20, 30, 40, 50]$ percentage of outliers added in random positions over the ground plane. In Figure 10, the results show that the SFM is especially important when the tracker is dealing with outliers. With 50% of outliers, the identity switches with SFM+GR are reduced 70% w.r.t the DIST results.

Noise. This test is used to determine the performance of our approach given noisy detections, which are very common mainly due to small errors in the 2D-3D mapping. From the GT set with 2% missing data, random noise is added to every detection. The variances of the noise tested are $[0, 0.002, 0.004, 0.006, 0.008, 0.01]$ of the size of the scene observed. As expected, group information is the most robust to noise; if the position of pedestrian A is not correctly estimated, other pedestrians in the group will contribute to the estimation of the true trajectory of A.

These results corroborate that having good behavioral models becomes more important as the observations deteriorate. In Figure 9 we plot the tracking results of a sequence with 12% simulated missing data. Only using distance information can see identity switches as shown in Figure 9(a). In Figure 9(b) we can see how missing data

affects the matching results. The matches are shifted, this chain reaction is due to the global optimization. In both cases, the use of SFM allows the tracker to interpolate the necessary detections and find the correct trajectories. Finally, in Figure 9(c) we plot the wrong result which occurs because track 3 has two consecutive missing detections. Even with SFM, track 2 is switched for 3, since the switch does not create extreme changes in velocity. In this case, the grouping information is key to obtaining good tracking results. More results are shown in Figure 13, first row.

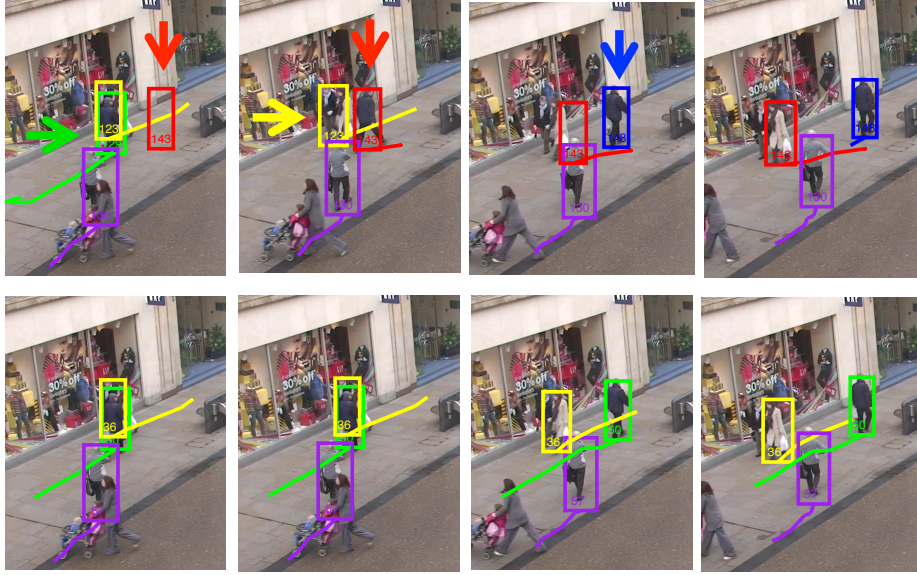


Fig. 11: Predictive approaches [10, 11] (first row) vs. Proposed method (second row)

5.3 Tracking results

We evaluate the proposed algorithm on two publicly available datasets: a crowded town center [36] and the well-known PETS2009 dataset [37]. We compare results with:

- [36]: using the results provided by the authors for full pedestrian detections. The HOG detections are also given by the authors and used as input for all experiments.
- [8]: globally optimum tracking based on network flow linear programming, for which we use our own implementation.
- [10]: tracker based on Kalman Filter which includes social behavior, using the code provided by the authors.
- [11]: tracker based on Kalman Filter which includes social and grouping behavior, using our own implementation.

For a fair comparison, we do not use appearance information for any method. The methods [10, 11, 36] are online, while [8] processes the video in batches.

Town Center dataset We perform tracking experiments on a video of a crowded town center [36]. To show the importance of social behavior and the robustness of our algorithm at low frame rates, we track at 2.5fps (taking one every tenth frame). We show detection accuracy (DA), tracking accuracy (TA), detection precision (DP) and tracking precision (TP) measures as well as the number of identity switches (IDsw).

	DA	TA	DP	TP	IDsw
HOG Detections	63.1	—	71.9	—	—
Benfold et al. [36]	64.9	64.8	80.5	80.4	259
Zhang et al. [8]	66.1	65.7	71.5	71.5	114
Pellegrini et al. [10]	64.1	63.4	70.8	70.7	183
Yamaguchi et al. [11]	64.0	63.3	71.1	70.9	196
Proposed	67.6	67.3	71.6	71.5	86

Table 1: Town Center sequence.

Note, the precision reported in [36] is about 9% higher than the input detections precision; this is because the authors use the motion estimation obtained with a KLT feature tracker to improve the exact position of the detections, while we use the raw detections. Still, our algorithm reports 64% less ID switches. As shown in Table 1, our algorithm outperforms [10], which includes social behavior, and [11], which includes also grouping information, by almost 4% in accuracy and with 50% less ID switches. In Figure 11 we can see an example where [10, 11] fail. The errors are created in the greedy phase of predictive approaches, where people fight for detections. The red false detection in the first frame takes the detection in the second frame that should belong to the green trajectory (which ends in the first frame). In the third frame, the red trajectory overtakes the yellow trajectory and a new blue trajectory starts where the green should have been. None of the resulting trajectories violate the SFM and GR conditions. On the other hand, our global optimization framework takes full advantage of the SFM and GR information and correctly recovers all the trajectories. More results of the proposed algorithm can be seen in Figure 13, last row.

Results on the PETS2009 dataset In addition, we perform monocular tracking on the PETS2009 sequence L1, Views 1,5,6,7,8 and obtain the detections using the Mixture of Gaussians (MOG) background subtraction method. We compare the results with the previously described methods plus the monocular result of View 1 presented in [9], where the detections are obtained using the Probabilistic Occupancy Map (POM) and the tracking is done using k-shortest paths.

The first observation that we make is that the linear programming methods (LP and Proposed) clearly outperform predictive approaches in accuracy. This is because this dataset is very challenging from a social behavior point of view, because the subjects often change direction and groups form and split frequently. Since our approach is based on a probabilistic framework, it is better suited for unexpected behavior changes (like

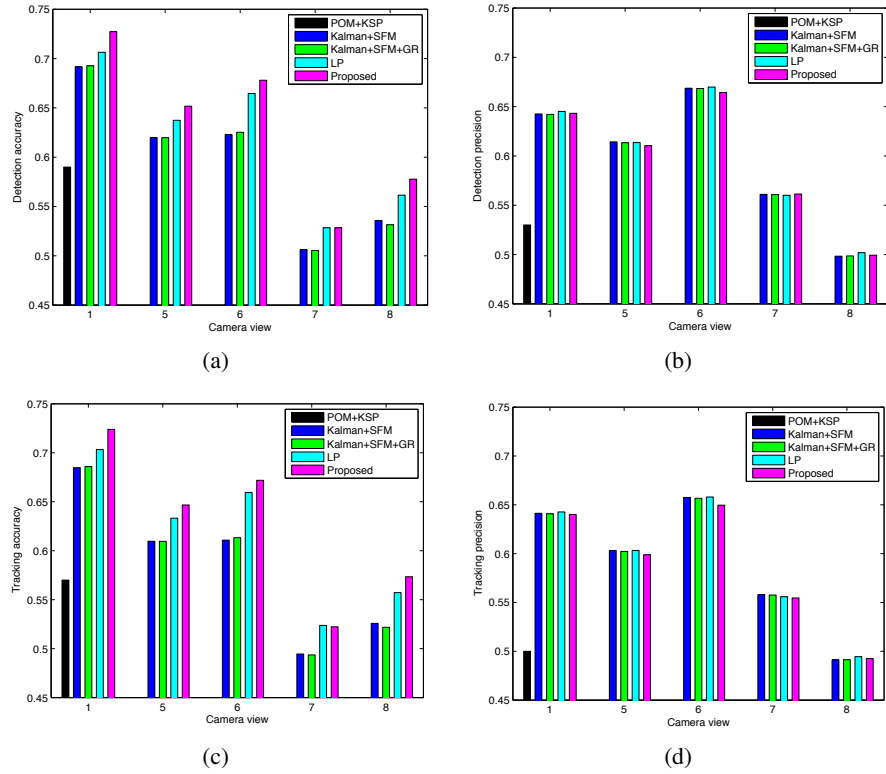


Fig. 12: Results of the proposed method on the PETS2009 dataset views 1,5,6,7,8. (a) Detection accuracy, DA. (b) Detection precision, DP. (c) Tracking accuracy, TA. (d) Tracking precision.

destination changes), where other predictive approaches fail [10, 11]. We can also see that the Proposed method has a higher accuracy in most views that the LP method, which does not take into account social and grouping behavior. The grouping term is specially useful to avoid identity switches between member of a group (see an example in Figure 13, third row, the cyan and green pedestrian who walk together). Precision is similar for all methods since the same detections have been used for all the experiments and we do not apply smoothing or correction of the bounding boxes. In general, views 7 and 8 are hard for tracking, due to 2D-3D calibration errors and a low field of view which means it is impossible to keep the identities and many small separate trajectories are created.

6 Conclusions

In this paper, we argued for integrating pedestrian behavioral models in a linear programming framework. Our algorithm finds the MAP estimate of the trajectories total

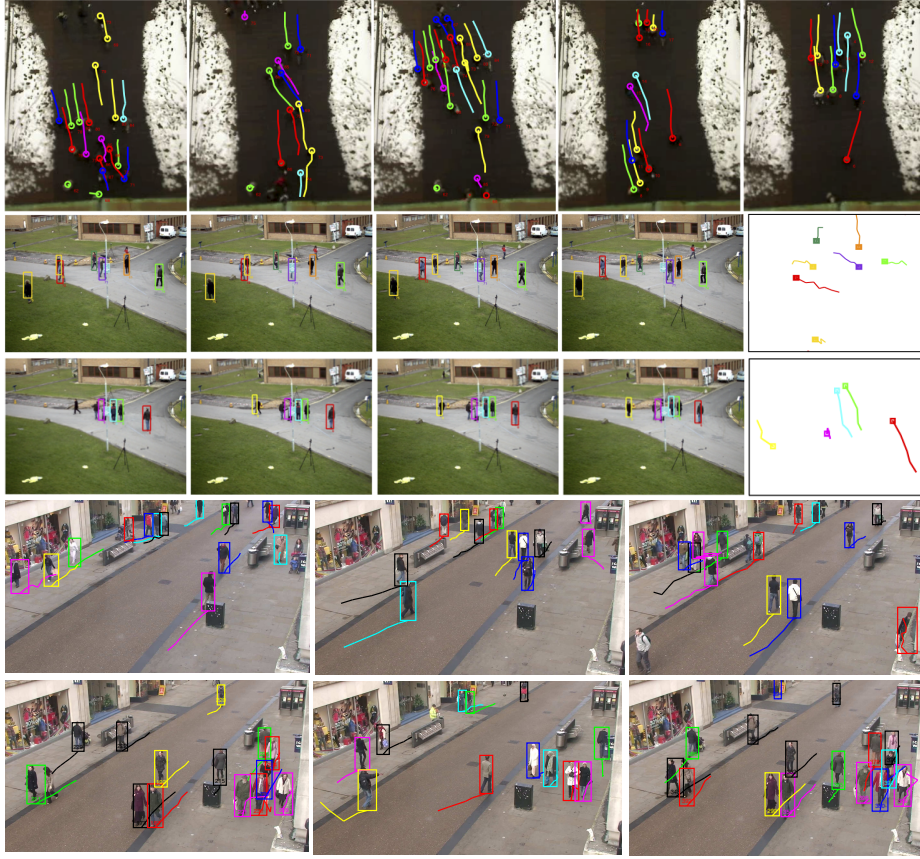


Fig. 13: *First row*: Results on the BIWI dataset (Section 5.2). The scene is heavily crowded, social and grouping behavior are key to obtaining good tracking results. *Second and third rows*: Results on the PETS2009 dataset (Section 5.3). *Last two rows*: Results on the Town Center dataset (Section 5.3).

posterior including social and grouping models using a minimum-cost network flow with an improved novel graph structure that outperforms existing approaches. People interaction is persistent rather than transient, hence the proposed probabilistic formulation fully exploits the power of behavioral models as opposed to standard predictive and recursive approaches such as Kalman filtering. Experiments on three public datasets reveal the importance of using social interaction models for tracking in difficult conditions such as in crowded scenes with the presence of missed detections, false alarms and noise. We present an extensive analysis of the effect of the parameters to show the robustness of our method. Results show that our approach is superior to state-of-the-art multiple people trackers. As future work, we plan on working on the optimization itself in order to find an efficient optimization method that keeps the linearity of the problem and at the same time does not require to iterate between computing the social forces and

computing the data association. On the other hand, we also plan to extend our approach to even more crowded scenarios where individuals cannot be detected and therefore features might be used as in [38]. This will be a first step to bridge macroscopic and microscopic approaches for crowd analysis.

Acknowledgements. This work was partially funded by the German Research Foundation, DFG projects RO 2497/7-1 and RO 2524/2-1.

References

1. Gall, J., Yao, A., Razavi, N., van Gool, L., Lempitsky, V.: Hough forests for object detection, tracking, and action recognition. *TPAMI* (2011)
2. Breitenstein, M., Reichlin, F., Leibe, B., Koller-Meier, E., van Gool, L.: Robust tracking-by-detection using a detector confidence particle filter. *ICCV* (2009)
3. Wu, B., Nevatia, R.: Detection and tracking of multiple, partially occluded humans by bayesian combination of edgelet part detectors. *IJCV* **75**(2) (2007)
4. Leibe, B., Schindler, K., Cornelis, N., van Gool, L.: Coupled detection and tracking from static cameras and moving vehicles. *TPAMI* **30**(10) (2008)
5. Kaucic, R., Perera, A., Brooksby, G., Kaufhold, J., Hoogs, A.: A unified framework for tracking through occlusions and across sensor gaps. *CVPR* (2005)
6. Ali, S., Shah, M.: Floor fields for tracking in high density crowded scenes. *ECCV* (2008)
7. Rodriguez, M., Sivic, J., Laptev, I., Audibert, J.: Data-driven crowd analysis in videos. *ICCV* (2011)
8. Zhang, L., Li, Y., Nevatia, R.: Global data association for multi-object tracking using network flows. *CVPR* (2008)
9. Berclaz, J., Fleuret, F., Türetken, E., Fua, P.: Multiple object tracking using k-shortest paths optimization. *TPAMI* (2011)
10. Pellegrini, S., Ess, A., Schindler, K., van Gool, L.: You'll never walk alone: modeling social behavior for multi-target tracking. *ICCV* (2009)
11. Yamaguchi, K., Berg, A., Ortiz, L., Berg, T.: Who are you with and where are you going? *CVPR* (2011)
12. Leal-Taixé, L., Pons-Moll, G., Rosenhahn, B.: Everybody needs somebody: Modeling social and grouping behavior on a linear programming multiple people tracker. *ICCV Workshops. 1st Workshop on Modeling, Simulation and Visual Analysis of Large Crowds* (2011)
13. Khan, Z., Balch, T., Dellaert, F.: Mcmc-based particle filtering for tracking a variable number of interacting targets. *TPAMI* (2005)
14. Leal-Taixé, L., Heydt, M., Rosenhahn, A., Rosenhahn, B.: Automatic tracking of swimming microorganisms in 4d digital in-line holography data. *IEEE Workshop on Motion and Video Computing (WMVC)* (2009)
15. Nillius, P., Sullivan, J., Carlsson, S.: Multi-target tracking - linking identities using bayesian network inference. *CVPR* (2006)
16. Yang, M., Yu, T., Wu, Y.: Game-theoretic multiple target tracking. *ICCV* (2007)
17. Berclaz, J., Fleuret, F., Fua, P.: Robust people tracking with global trajectory optimization. *CVPR* (2006)
18. Jiang, H., Fels, S., Little, J.: A linear programming approach for multiple object tracking. *CVPR* (2007)
19. Andriyenko, A., Schindler, K.: Globally optimal multi-target tracking on an hexagonal lattice. *ECCV* (2010)

20. Wu, Z., Kunz, T., Betke, M.: Efficient track linking methods for track graphs using network-flow and set-cover techniques. *CVPR* (2011)
21. Helbing, D., Molnár, P.: Social force model for pedestrian dynamics. *Physical Review E* **51** (1995) 4282
22. Mehran, R., Oyama, A., Shah, M.: Abnormal crowd behavior detection using social force model. *CVPR* (2009)
23. Pelechano, N., Allbeck, J., Badler, N.: Controlling individual agents in high-density crowd simulation. *Eurographics/ACM SIGGRAPH Symposium on Computer Animation* (2007)
24. Scovanner, P., Tappen, M.: Learning pedestrian dynamics from the real world. *ICCV* (2009)
25. Luber, M., Stork, J., Tipaldi, G., Arras, K.: People tracking with human motion predictions from social forces. *ICRA* (2010)
26. GE, W., Collins, R., Ruback, B.: Automatically detecting the small group structure of a crowd. *WACV* (2009)
27. Choi, W., Savarese, S.: Multiple target tracking in world coordinate with single, minimally calibrated camera. *ECCV* (2010)
28. Pellegrini, S., Ess, A., van Gool, L.: Improving data association by joint modeling of pedestrian trajectories and groupings. *ECCV* (2010)
29. Bertsekas, D.: *Nonlinear programming*. Athena Scientific (1999)
30. Dantzig, G.: *Linear programming and extensions*. Princeton University Press, Princeton, NJ (1963)
31. Makhorin, A.: Gnu linear programming kit (glpk). <http://www.gnu.org/software/glpk/> (2010)
32. Pirsivash, H., Ramanan, D., Fowlkes, C.: Globally-optimal greedy algorithms for tracking a variable number of objects. *CVPR* (2011)
33. Ahuja, R., Magnanti, T., Orlin, J.: *Network flows: Theory, algorithms and applications*. Prentice Hall (1993)
34. Suurballe, J.: Disjoint paths in a network. *Networks* **4** (1974) 125–145
35. Kasturi, R., Goldgof, D., Soundararajan, P., Manohar, V., Garofolo, J., Boonstra, M., Kozhova, V., Zhang, J.: Framework for performance evaluation for face, text and vehicle detection and tracking in video: data, metrics, and protocol. *TPAMI* **31**(2) (2009)
36. Benfold, B., Reid, I.: Stable multi-target tracking in real-time surveillance video. *CVPR* (2011)
37. Ferryman, J.: *Pets 2009 dataset: Performance and evaluation of tracking and surveillance*. (2009)
38. Brostow, G., Cipolla, R.: Unsupervised detection of independent motion in crowds. *CVPR* (2006)