

GReTA – a novel Global and Recursive Tracking Algorithm in three dimensions

Alessandro Attanasi, Andrea Cavagna, Lorenzo Del Castello, Irene Giardina, Asja Jelić, Stefania Melillo, Leonardo Parisi, Fabio Pellacini, Edward Shen, Edmondo Silvestri, Massimiliano Viale

Abstract—Tracking multiple moving targets allows quantitative measure of the dynamic behavior in systems as diverse as animal groups in biology, turbulence in fluid dynamics and crowd and traffic control. In three dimensions, tracking several targets becomes increasingly hard since optical occlusions are very likely, i.e. two featureless targets frequently overlap for several frames. Occlusions are particularly frequent in biological groups such as bird flocks, fish schools, and insect swarms, a fact that has severely limited collective animal behavior field studies in the past. This paper presents a 3D tracking method that is robust in the case of severe occlusions. To ensure robustness, we adopt a global optimization approach that works on all objects and frames at once. To achieve practicality and scalability, we employ a divide and conquer formulation, thanks to which the computational complexity of the problem is reduced by orders of magnitude. We tested our algorithm with synthetic data, with experimental data of bird flocks and insect swarms and with public benchmark datasets, and show that our system yields high quality trajectories for hundreds of moving targets with severe overlap. The results obtained on very heterogeneous data show the potential applicability of our method to the most diverse experimental situations.

Index Terms—tracking, 3D, multi-object, multi-path, branching, global optimization, recursion, divide and conquer



1 INTRODUCTION

IN recent years there has been a growing interest in studying the motion of large groups of objects, both in two and in three dimensions: animals, humans, automotive vehicles, cells and microorganisms in field or laboratory experiments, as well as tracer particles in turbulent fluids flows [1], [2], [3], [4], [5], [6], [7], [8], [9], [10], [11], [12], [13], [14], [15], [16]. This kind of studies requires tracking, the automated process of following in space and time individual objects using visual information from single- or multi-camera video sequences. Sometimes experimental data contains information on the object’s features, so that for example color- or pattern-matching strategies can be exploited to simplify the problem. This, however, is not our case: we focus here on the three-dimensional tracking problem using stereometric information only. Examples of input and output data of our tracking algorithm are shown in Fig. 1.

There are two main reasons why tracking is hard. First, when the average inter-object distance in the images is small compared to the average displacement of the objects between two consecutive frames, ambiguities

arise when identifying individual objects in time. This is easily solved using cameras with a sufficiently high temporal resolution.

A second and far more serious difficulty arises when the average inter-objects distance in the images is small compared to their optical size, making optical occlusions highly likely. Each time an ambiguity due to an occlusion occurs, there is a high probability that the tracked trajectories of the objects involved are interrupted.

These interruptions are a minor problem when estimating velocity fields, but the situation becomes more problematic when we use the velocities to infer the inter-individual interactions within a group of animals. Several interruptions at any given time frame are equivalent to missing some of the individuals, which potentially biases the inferred interaction. The problem is even more serious when we measure observables that depend on the *entire* individual trajectories such as diffusion properties [17] or the kinematics of turning [18] in collective animal behavior. Even in turbulence studies, the lack of complete trajectories can introduce serious statistical biases on some physical observables [19].

In fact, interruptions are the best case scenario when we have many optical occlusions. The worst case is the introduction of non-existent trajectories that mix the identities of two different objects, especially problematic in physical and biological analysis.

1.1 Literature survey

Tracking algorithms differ according to how they exploit the information available in the images. In the last thirty years, several algorithms have been developed in the

- *The authors are with the Istituto Sistemi Complessi, Consiglio Nazionale delle Ricerche, UOS Sapienza, 00185 Rome, Italy. E-mail: see <http://www.cobbs.it>*
- *A. Attanasi, A. Cavagna, L. Del Castello, I. Giardina, A. Jelić, S. Melillo, and M. Viale are with the Dipartimento di Fisica, Università Sapienza, 00185 Rome, Italy.*
- *L. Parisi and F. Pellacini are with Dipartimento di Informatica, Università Sapienza, 00198 Rome, Italy.*
- *E. Shen is with Bublcam Technology Inc., Toronto, Canada.*
- *E. Silvestri is with the Dipartimento di Matematica e Fisica, Università Roma Tre, 00146 Rome, Italy.*

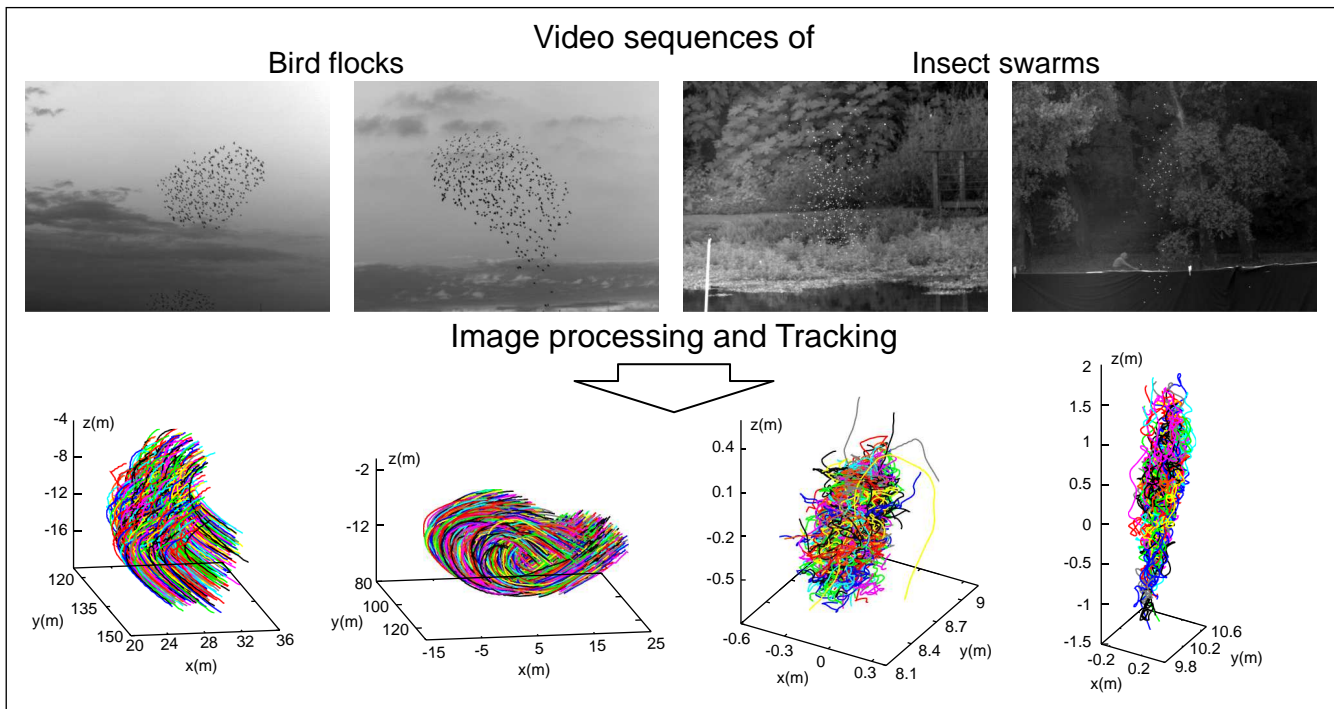


Fig. 1. Input and output data of our tracking algorithm. Top row: examples of original images extracted from the video sequences taken during four field experiments of flocking birds and swarming insects (from left to right, experimental events $E3$, $E6$, $E14$, and $E15$, respectively). See Section 5 and Table 2 therein). The original images are slightly cropped and enhanced for the sake of readability. Bottom row: the 3D reconstruction of the full trajectories for each experimental event.

field of fluid dynamics [1], [2], [3]. In particular, the algorithm by Ouellette *et al.* [3], in case of ambiguities, optimizes the solution for all tracked particles locally in time. In the field of collective animal behavior, different algorithms have been developed to reconstruct the 3D-positions of individual animals in groups [4], [5], [6]. The technical shortcomings and limited results of these initial studies have been the catalyst for the successive empirical investigations on collective animal phenomena, and lead to the development of tracking algorithms tested on various animals as fruit flies [11], [15], [20], mosquitoes [21], [22], bees [23], bats [10], [14], [24], and fish [25].

A recent significant breakthrough in the field is represented by the Multiple Hypothesis Tracking (MHT) approach [26], which finds objects correspondences across multiple views through an NP-hard multidimensional assignment problem. MHT methods based on global optimization over the space of the tracked objects and/or time have been implemented with greedy approaches. Betke *et al.* [10] developed an algorithm based on a multi-dimensional assignment problem solved with a greedy approximation. Zou *et al.* [11] implemented a tracking algorithm which uses a global correspondence selection scheme, and applies Gibbs sampling locally in time to reduce the complexity of the algorithm. H.S. Wu *et al.* [12], [13] implemented a different algorithm based on three linear assignment problems, making use of ghost objects to partially solve the problem of short-

term optical occlusions. Liu *et al.* [15] proposed a very efficient algorithm in the framework of particle filter able to deploy weak visual information to distinguish the identities of the tracked objects. More interesting is the approach proposed by Z. Wu *et al.* [14], who recognized the importance of a global optimization over the full temporal sequence and over all the tracked objects, posing the problem in the form of a weighted set-cover.

Our efforts focus on 3D tracking of large groups of featureless objects for long temporal sequences, and the long-term optical occlusions typical of our experimental data need to be addressed in a different way. Therefore we aim at the globally optimal solution of the problem, and not at efficient and fast ways to approximate it with greedy approaches. The bottleneck of this strategy is the computational complexity which grows exponentially with the duration of the acquisition.

1.2 Our tracking approach

We propose here a novel Global and Recursive Tracking Algorithm (GR_eTA), an approach which dramatically reduces the computational complexity of the global optimization problem thanks to a recursive divide and conquer strategy. Within this new framework, we can optimize the solution globally over longer temporal sequences. In order to preserve the global scope, we introduce a way to extend the temporal horizon over which ambiguous choices are made within the divide and conquer scheme.

Thanks to this method, we are able to resolve optical occlusions lasting up to dozens of consecutive frames, and therefore to distinguish the identities of the tracked objects without creating interruptions, even when the optical density in the images is very large. The reconstructed trajectories have negligible fragmentation even in the presence of large optical density and frequent occlusions.

We validate our approach using synthetic data as ground-truth, and we test its potential by applying it to original experimental field data of flocking birds and swarming insects.

The rest of the paper is divided into 6 sections. Section 2 explains the algorithm. In Section 3 we analyze the problem complexity. Sections 4 and 5 report the validation of our algorithm with synthetic data and experimental field data. In Section 6 we present a comparison with prior works, and the conclusions in Section 7.

2 METHODS

The main idea of our method is that when an occlusion occurs, we assign multiple temporal links and we use these links to create all possible paths running through the occlusion. Many of these paths will be non-physical, but certainly the paths corresponding to the real objects will also be there. Then, the information from all cameras is assembled and the selection of the physically meaningful paths, namely those that optimize multi-camera coherence, is performed globally in space (over the tracked objects) and in time (over the temporal sequence), making use of a recursive divide and conquer scheme.

2.1 The basic steps of a tracking system

The goal is to track individual objects in time while reconstructing their positions in 3D space. We use stereoscopic video-sequences of the target objects acquired via a synchronized and calibrated three-camera system. The data gathering procedures we used in our experiments are described in Appendix A.

Image segmentation. The first step of a tracking algorithm is the detection of the objects in the images, done by image segmentation, see Fig. 2, first and second rows. Several approaches may be used to perform the segmentation, and the choice strongly depends on the type of objects. Our approach to image segmentation is not an essential part of the tracking system we propose, so we leave its description to Appendix B.

Stereoscopic linking. The second step is to compute the stereoscopic linking of the detected objects, which consists of matching the individual objects across the images acquired by different cameras at the same time, see Fig. 2, third row. We assign multiple stereoscopic links between the object images as seen by three cameras using standard trifocal geometry [27]. The details of

the linking method do not matter, and we describe the exact procedure in Appendix C.

Temporal linking. The third step of our algorithm is to assign multiple temporal links for each object, which consists of matching individual objects from one frame to the next one, as shown on the same third row of Fig. 2. We use different prediction strategies according to the specific data we process. The precise details of the linking methods are not essential, and the exact procedures are described in Appendix C.

Tracking. Recent global optimization approaches, as the one we propose here, rely on the assumption that objects may be linked to several other objects (multi-linking instead of one-to-one linking), and the global optimization is performed over the space of these links to select the matches corresponding to real 3D trajectories. In the following sections, we explain the method and we present the formalisms of our tracking approach.

2.2 Multi-path branching

Let us consider the example shown in Fig. 3, which illustrates a partial temporal sequence of two objects A and B as seen by two cameras. The two objects overlap in the image of the left camera for three frames. Most prior work would assign only one temporal link to each detected object, therefore the points of occlusion belong to only one uninterrupted trajectory, the second recovered trajectory being broken. This results in a fragmented trajectory. Furthermore, assigning temporal links using information that is purely local has the drawback that the identities of occluding objects might be lost.

To tackle this concerns, we use a path branching approach with global optimization. For the example in Fig. 3, we assign multiple temporal links and create all possible paths running through the occlusion in the left camera view. In this case, there are four paths, of which two are real (AA and BB) and two have hybrid object identities (AB and BA). In order to build the set of all possible paths through the segmented objects, the temporal links are propagated for each camera view to build the temporal graph of each camera. We then have to solve the problem of how to select the correct paths in the 2D graph of each camera and match them across cameras. The advantage of our approach is that at an early stage each object can have more than one path, which is what is needed to handle occlusions.

2.3 Global optimization

The selection of the correct matches between the 2D paths across cameras is the core of the tracking problem. We create all the possible 2D paths propagating the temporal links in the image space of each camera, while

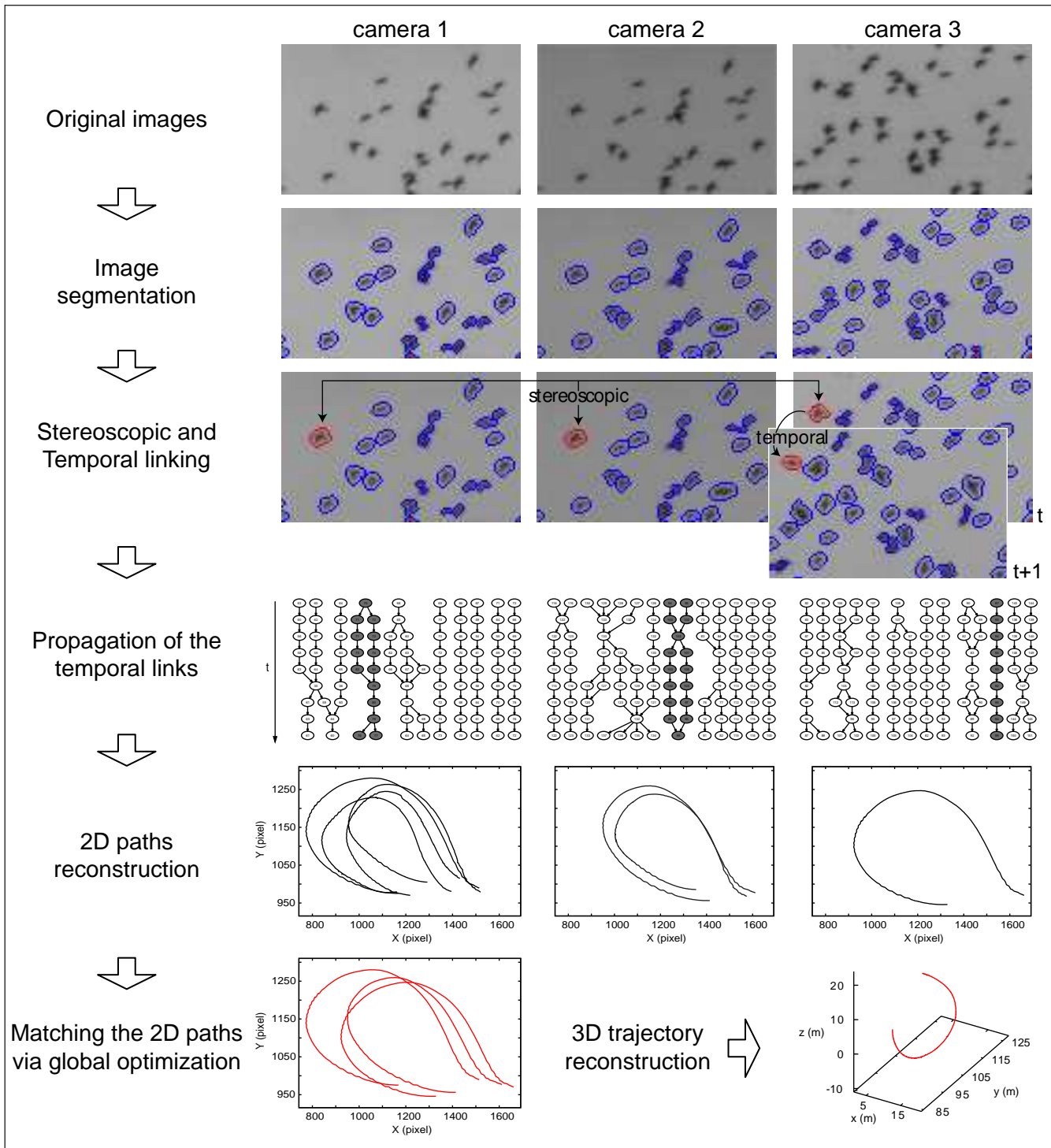


Fig. 2. Scheme illustrating the main steps of the full tracking system (image processing and tracking algorithms). A small crop of the original images extracted from the video sequences of three synchronized cameras are shown on the first row. The segmented images are shown on the second row using blue color for the object borders, and red for their centers of mass. On the third row, we show one example of a (trifocal) stereoscopic link connecting the views of the same object in the three cameras, and one example of temporal link connecting the same object between subsequent frames in each camera sequence. In the fourth row, we show a crop of the temporal graphs for each camera view, which represents a useful visualization of the full set of temporal links assigned for each camera. The figures show only a small crop representing a few objects for 9 time frames, and we indicated with grey color a cluster of paths stereoscopically linked across the three views. The fifth row illustrates the 2D paths reconstructed in the image space of each camera, obtained by simple propagation of the temporal links. The algorithm outputs all possible 2D paths at this step. Global optimization is used to match the correct 2D paths between the camera views, as shown on the sixth and last row, from which we finally obtain the 3D the trajectory using standard stereoscopic geometry.

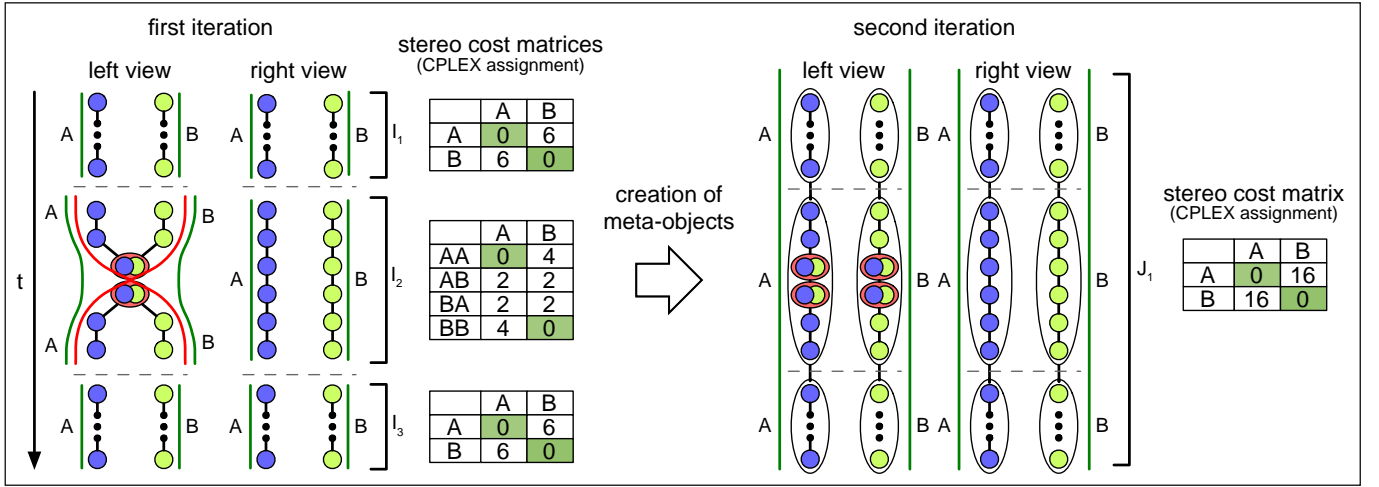


Fig. 3. Scheme illustrating a real case of confined data. Two objects A and B occlude each other for two frames in the left view, while they always appear as separate objects in the right view. During the first iteration of the recursive divide and conquer approach, the event is divided into three intervals, I_1 , I_2 , and I_3 . In the first and in the third intervals, there are no tracking ambiguities. The propagation of the temporal links results only in the two correct 2D paths, A and B , in each camera. We define the cost of a pair of 2D paths as the sum of the costs of the links between them, i.e. the number of missing stereoscopic links. The global optimization selects the correct matches, (A, A) and (B, B) . In the second interval, there are still not ambiguities in the right view: the only 2D paths, A and B , are correct. Instead, in the left view, we propagate the temporal links and we create 4 2D paths; the two correct ones (green lines) AA and BB , and the two wrong ones (red lines) AB and BA . Two of the possible set-cover solutions Γ are: the correct one $G \equiv \{(AA, A), (BB, B)\}$ with a cost equal to 0, and the wrong one $\{(AB, A), (BA, B)\}$ with a cost equal to 4. Here, the global optimization is essential to select the correct solution. All the matched 2D paths are then analyzed at the second iteration as meta-objects. At this iteration, there are no occlusions. Propagating the temporal links, we create two 2D paths in each camera view, and the tracking problem is correctly solved for the entire duration of the event.

we choose how to match them using stereoscopic links. The assumption here is that the 2D paths representing the same 3D object are strongly linked stereoscopically. On the contrary, 2D paths corresponding to different 3D objects are loosely linked stereoscopically. We define a measure of the stereoscopic quality of each match, i.e. a cost function, and we use a global optimization approach to retrieve the set of the correct matches. This is an NP-hard multidimensional assignment problem, and we solve it using Integer Linear Programming (ILP) techniques in order to find the globally optimal solution [28].

Definition of trajectory. Consider a system of three cameras, and denote a trajectory γ as a triplet of matched 2D paths, $\gamma = (\gamma_1, \gamma_2, \gamma_3)$. Each 2D path γ_i represents a temporal sequence of 2D objects detected in the images of the i -th camera, and connected by temporal links. Moreover the triplet $(\gamma_1, \gamma_2, \gamma_3)$ is stereoscopically linked for at least one frame. Let Γ be the set of all the possible trajectories. The goal is to find the correct subset of trajectories $G \subseteq \Gamma$, see Fig. 3.

Cost function. We evaluate the quality of each subset $\hat{\Gamma} \subseteq \Gamma$ by defining a cost function $C(\hat{\Gamma})$. Let $C(\hat{\Gamma}) = \sum_{\gamma \in \hat{\Gamma}} c(\gamma)$, where $c(\gamma)$ is a cost associated to the trajectory γ and based on the stereoscopic coherence (the higher the quality, the lower the cost), see Fig. 3.

Let us formally define the cost function. Considering a three-camera system, for each $\gamma \in \Gamma$ and at each instant

of time, the cost function $c(\gamma(t)) = c(\gamma_1(t), \gamma_2(t), \gamma_3(t))$ is defined as the trifocal distance [27] in the case of matched triplets, and as the epipolar distance [27] in the case of pairs (corresponding to miss-detection in one camera). Whenever the cost exceeds a threshold value c_{max} or in the case of absence of a stereoscopic link, we set $c = c_{max}$. The cost of a trajectory γ is then defined as the temporal average:

$$c(\gamma) = \frac{\sum_{t \in T_\gamma} c(\gamma(t))}{|T_\gamma|}, \quad (1)$$

where T_γ is the set of time frames, t , for which $c(\gamma(t))$ is defined.

2.3.1 Formalization of the tracking problem

Let us distinguish between two different types of input data.

Confined data: the objects are in the common field-of-view of the camera system at least for a short temporal sequence. Each segmented object belongs to at least one trajectory, and the solution of the tracking problem is a cover for the set of all the objects.

Non-confined data: one or more objects never appear in the common field-of-view of the camera system, but they are seen by one camera only. Therefore they are far from the objects of interest in three dimensions, and they should not be matched as they do not belong to any trajectory. A

typical example is represented by pollen particles passing in front of one camera only, and appearing as large blurred objects. The problem becomes more complex, and the covering condition needs to be relaxed to exclude these objects.

Confined data, joint weighted set-cover. When applied to confined data, the global optimization approach is equivalent to a joint weighted set-cover (as in [14]). The tracking problem can be formulated as:

$$c(\Gamma_{opt}) = \min_{\{x\}} \sum_{\gamma \in \Gamma} c(\gamma)x_{\gamma}, \quad (2)$$

with the constraint:

$$\forall p, \quad \sum_{\gamma \in \Gamma_p} x_{\gamma} \geq 1, \quad (3)$$

where x_{γ} is a boolean variable associated to γ , p is a 2D object in the image space of a camera, and Γ_p is the set of all trajectories passing by p . The retrieved set $\Gamma_{opt} \equiv \{\gamma \mid x_{\gamma} = 1\}$ covers with the best weight the full set of segmented objects.

It can be proven that, under suitable conditions, the global optimization approach finds the correct solution. Indeed, in the case of confined data, when all the correct temporal and stereoscopic links are known and when some particular ambiguities are forbidden (for a formal definition, see Appendix D), the correct solution of the tracking problem is the only set-cover minimizing the cost defined by Eq. 2 with the constraint in Eq. 3. We refer the reader to Theorem 1 in Appendix D for the exact list of hypotheses holding this statement, together with its proof.

Non-confined data, relaxed joint weighted set-cover. In the case of non confined data, not all the segmented objects can be tracked. We need to discard those objects not appearing in all the three cameras, therefore lacking stereoscopic correspondance, because they do not belong to the group of interest. To this aim, we need to relax the covering constraint in Eq. 3. Let us introduce for each detected object p a new boolean variable y_p . The relaxed joint weighted set-cover problem is then formalized as:

$$\min_{\{x,y\}} \left[\sum_{\gamma \in \Gamma} c(\gamma)x_{\gamma} + \frac{\lambda}{T} \sum_p (1 - y_p) \right], \quad (4)$$

where T is the event duration, and with the constraint:

$$\forall p, \quad \sum_{\gamma \in \Gamma_p} x_{\gamma} \geq y_p. \quad (5)$$

The contribution of a discarded object, for which $y_p = 0$, to the global cost of the solution is λ/T . Assigning to λ a value lower than the highest cost assigned to the stereoscopic links (the threshold value c_{max}), we manage to exclude from the solution the objects detected only in

one camera view. We experimentally choose $\lambda = 0.9c_{max}$. See the example sketched in Fig. 4.

In our implementation, the optimization problem is solved using linear programming, for which we use the library in [29].

2.4 Recursive divide and conquer

The computational complexity of global optimization problem strongly limits the size of the datasets that can be processed. In order to reduce the complexity, the full temporal sequence can be divided into shorter intervals over which smaller optimization problems can be solved. A well-known method used to join the subtrajectories constructed within limited time windows is the sliding window approach [24]. This approach matches the subtrajectories of the first interval with the ones of the second one, then the ones of the second with the ones of the third interval, repeating the procedure until the full trajectories are recovered. Such approach is very efficient and extremely powerful when applied to sparse data or when the tracked objects can be identified using features like shape, pattern, or color. Its weakness resides in the optimization which is not performed globally in time. For this reason, the identities of the objects can easily be lost whenever treating dense data of featureless objects.

The GReTA algorithm we propose here is based on a recursive divide and conquer strategy. We divide the acquisition into temporal intervals with length $\tau_1 < T$. The optimization described in the previous section is

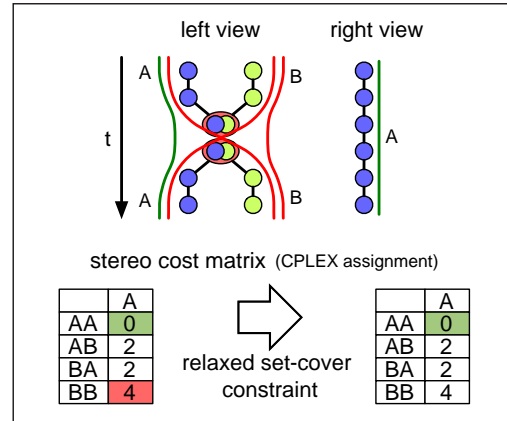


Fig. 4. Scheme illustrating a real case of non-confined data, i.e. data corrupted by the images of objects in one camera view which do not belong to the group of interest and do not appear in the common field-of-view of all cameras (e.g., insects or birds flying in front of one camera only, or a pollen particle passing in front of a camera lens). A tracked object A and a pollen particle B occlude each other in the left view, while only A is visible in the right view. There are only two possible set-cover solutions, both characterized by the same cost equal to 4: $\{(AA, A), (BB, A)\}$ and $\{(AB, A), (BA, A)\}$. Both solutions would produce wrong trajectories, and the algorithm would fail to find the correct solution. Relaxing the set-cover constraint, the correct solution is found as the trajectory (AA, A) , while the objects belonging only to the 2D path BB in the left view are discarded.

performed in each time interval. Each path – a temporal sequence of linked 2D objects – belonging to a selected trajectory becomes then a *meta-object*, which can be linked stereoscopically and in time to other meta-objects (paths). The procedure is then iterated. A new time interval τ_2 is selected, where now τ_2 enumerates the number of intervals of length τ_1 , *i.e.* of *meta-frames*. The procedure is applied recursively, until the product of the τ 's of each iteration equals the duration of the entire acquisition, $\tau_1\tau_2\tau_3 = T$. Finally, the partial solutions retrieved at each iteration are combined into the solution of the full problem at the last iteration. It is possible to prove that, when the conditions that guarantees the uniqueness of the solution (see Section 2.3.1 and Appendix D) are satisfied within each interval of length τ_1 , the solution obtained using the recursive approach coincides with the one obtained solving the problem over the entire temporal sequence. The reader is referred to Corollary 1 in Appendix D.

Note that the recursive approach offers two key advantages when compared to the classical sliding window one. First, it permits to evaluate the optimization problem for several interval interfaces at once, giving it a more global scope. Second, it allows postponing ambiguous choices at each iteration to later ones, effectively extending the temporal horizon over which these choices are made.

2.5 Making the algorithm robust against wrong or missing links

Dealing with real data, the sets of temporal and stereoscopic links are affected by noise, which results in missing links and fluctuations of the stereoscopic distances. In some particular situations, the optimization operated within finite intervals of time is not guaranteed to be equivalent to a truly global optimization over the entire temporal sequence.

We describe here two modifications of the algorithm that take into account such situations typical of experimental data: a way to postpone ambiguous choices to the next iteration, effectively extending the temporal horizon over which a choice is made, and a way to recognize and re-join fragments of the same trajectory.

2.5.1 Postponing ambiguous choices to the following iterations

The absence of some correct links may lead to one or more trajectories representing unreal objects. These trajectories are characterized by at least one long time gap during which the stereoscopic links are absent.

We detect these cases by using a threshold over the maximum acceptable number of consecutive frames of missing stereoscopic links, and we discard them. We then run the optimization algorithm. Next, we propagate the links over all the discarded 2D objects, creating new 2D paths. These are then passed to the optimization algorithm at the following iteration together with all

the matched 2D paths. In this way, we discard any ambiguous choice made locally within any interval at the current iteration, and postpone the decision to the following iteration, effectively extending the temporal window when necessary.

Such refined algorithm is applied at each iteration. At the last iteration, the trajectories lacking stereoscopic coherence are discarded. New 2D paths are obtained by propagating through the objects which belonged to the discarded trajectories, and they are added to the set of 2D paths. Finally, the set of all the 2D paths is passed to the optimization algorithm running for a second and last time over the full temporal sequence. This time, the trajectories lacking coherence will not be discarded.

2.5.2 Joining trajectory fragments

There are two reasons for the algorithm to output correct but fragmented trajectories. First, when a temporal link is missing. Second, when the modification described above breaks a wrong trajectory and reconstructs two correct fragments. In both cases, it is possible to re-join fragments of trajectories which are consecutive in time and stereoscopically coherent. We do so after each iteration by connecting fragmented 2D paths in each camera that are stereoscopically connected to the same full 2D path in another camera. Note that in our implementation that is designed for a system of three cameras, we actually match fragments of paths in one camera only with matched pairs of paths in the other two cameras.

2.6 Final quality check and 3D reconstruction

In the case of field experiments with freely moving animals, individuals happen to leave the field-of-view of one camera for long times during the recorded events. Furthermore, the noise present in real experimental images often results in errors of the segmentation routine, both miss-detections and over-detections. Because of these reasons, it is not possible to correctly track all the objects for the entire temporal sequence, especially when the size of the problem is very large and the image data is heavily corrupted with noise.

To ameliorate these issues, we discard those few trajectories lacking stereoscopic coherence for a considerably long time gap. As shown in the following Section 4, this amounts to roughly 4% of the final trajectories for an average-sized dataset (512 objects and 500 frames, see Table 1). We then cut those trajectories, in order to save the fragments which do satisfy the stereoscopic coherence. Such operation results in a minor trajectory fragmentation.

Finally, we are left with a set of matched triplets of 2D paths. Given a triplet of 2D paths, we can reconstruct the corresponding trajectory in 3D by applying standard stereometric formulas [27] to each triplets of synchronous 2D points belonging to the paths, as shown on the last row of Fig. 2.

3 COMPLEXITY OF THE TRACKING PROBLEM

The global optimization requires comparing all the possible solutions and selecting the one that minimizes the cost function. This is a multidimensional assignment, and it is NP-hard. We solve it by finding the globally optimal solution using ILP techniques. We compute the cost of each possible triplet of linked 2D paths $\gamma \in \Gamma$, *i.e.* with at least one stereoscopic link. This implies that the number of variables to handle, H , corresponds to the number of possible trajectories, $|\Gamma|$. The parameter H strictly depends on the number P of 2D paths in the graph of each camera, obtained by propagation of the temporal links, and on the stereoscopic links between them. Therefore both H and P depend on the number of objects to be tracked, and they both grow exponentially with the temporal duration of the event. Let us analyze in detail such trend.

3.1 Complexity of the temporal graph of each camera space

The number of 2D paths, P , for a certain camera depends on the temporal length of the acquisition and on the connectivity of the graph on that camera. This dependence can be described by introducing a bifurcation coefficient $\alpha \geq 0$, which is an indirect measure of the number of occlusions per frame in each camera view. The higher is the number of occlusions between the detected objects, the higher is the average number of multiple links per

object in the corresponding graph, the higher is the value of α . We can predict P as a function of α , of the number of tracked objects N , and of the event duration T , as:

$$P = Ne^{\alpha T}. \quad (6)$$

For $\alpha = 0$, the number of paths is exactly equal to the number of real objects, hence $\alpha = 0$ corresponds to the ideal case of zero occlusions. Eq. 6 is confirmed by tests on synthetic data. In Fig. 5(a), the values of P as a function of T are plotted for a synthetic dataset of flocking birds ($N = 1024$, see Section 4 for details). For each T , we propagate only the correct temporal links and we measure $P(T)$ for several intervals lasting T frames. The mean value is plotted against T , and a linear fit is performed to retrieve the value of α . Typically, our experimental data (birds and insects) are characterized by $\alpha \in [0.001, 0.2]$.

3.2 Full computational complexity of the problem

The number of paths P is not by itself the computational bottleneck of the algorithm. What really matters is how many triplets built out of these P paths have a nonzero probability to be stereoscopically connected to each other, because this is what actually enters the global optimization problem. We call H the number of possible stereoscopic matches between the 2D paths across cameras. In the best case scenario, *i.e.* when there are no stereoscopic ambiguities, each one of the 2D paths

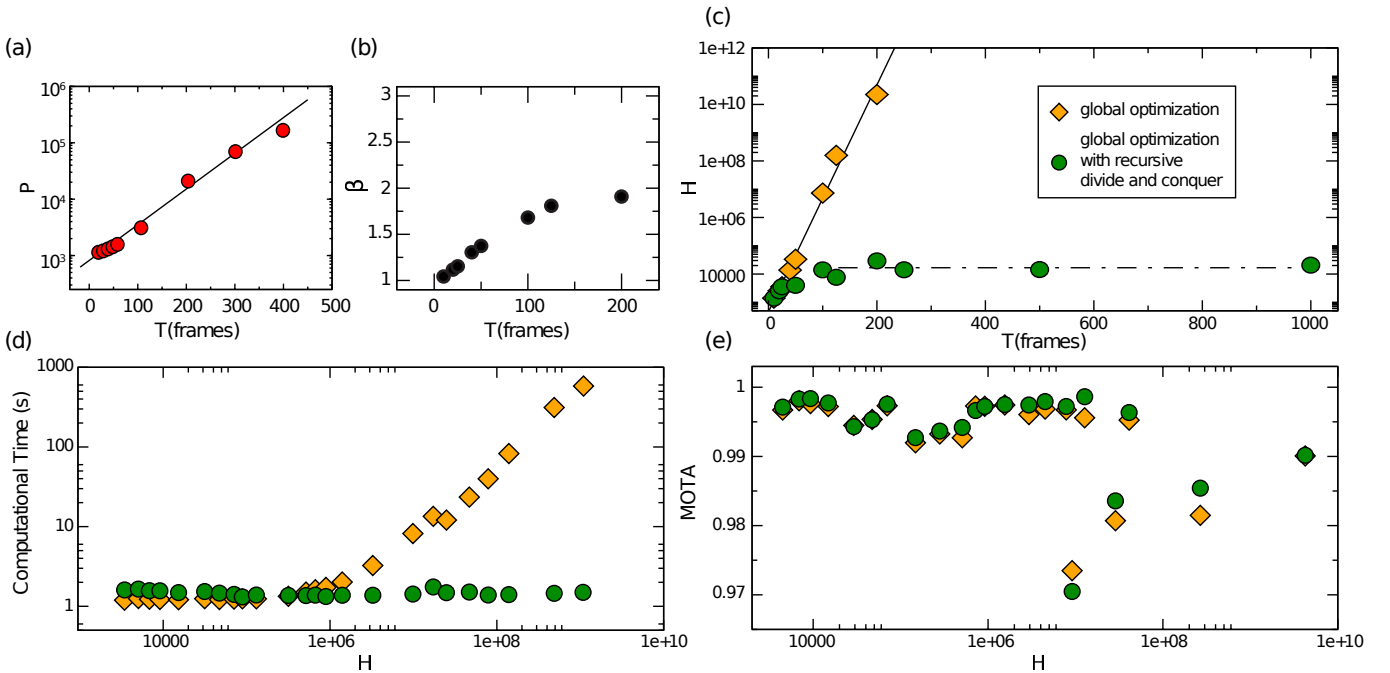


Fig. 5. The intrinsic complexity of the tracking problem is shown with synthetic data, with and without the recursive divide and conquer scheme is compared. In panel (a), the number P of possible paths for a synthetic dataset of 1024 objects is shown as a function of the considered temporal duration T . In panel (b), the dependence of β on T . In panel (c), the values of H as a function of T are plotted as yellow diamonds; these are compared to the values of H obtained when the recursive divide and conquer strategy is applied (choosing $\tau_1 = 25$ frames), and plotted as green circles. In panel (d), the computational time is plotted versus H for several runs with different complexities, with and without recursion as in panel (c). In panel (e), the MOTA (see Sec. 4) values obtained on several runs with different complexities are plotted versus H , with and without recursion as in panel (c).

belongs to only one trajectory (*i.e.*, $\gamma \in \Gamma$), and $H = P$. In the worst case scenario, each one of the 2D paths has at least one stereoscopic link with every other path in the image spaces of the other cameras and, for a three-camera system, $H = P^3$. We can therefore express H as a function of P in the following way,

$$H = P^\beta = (Ne^{\alpha T})^\beta, \quad (7)$$

where the parameter $\beta \in [1, 3]$ gives the measure of the degree of hybridization between 2D paths on different cameras, which is in turn a function of the optical density of the objects and of the number of occlusions. The longer the temporal acquisition, the higher the probability that one 2D path intersects another one; this is effect is large if there is high diffusion of the real 3D objects in the center of mass reference frame of the group. Therefore, we expect the exponent β to grow with the time duration of the event. This growth of $\beta(T)$ is indeed what we find (see Fig. 5(b)); although the saturation limit $\beta \sim 2$ is below the upper bound $\beta = 3$, the growth of β with time means that the exponential explosion of the computational complexity H is rather severe.

In Fig. 5(c) we report the computational complexity H as a function of the number of frames T for the same synthetic dataset (yellow diamonds). The exponential growth of H clearly shows that a multi-path branching algorithm by itself cannot solve the matching problem for long time intervals, however optimized it is and however large the memory resources are. Indeed, the very existence of optical occlusions, which is the reason to bifurcate paths in the first place, makes it impossible to reduce significantly the values of α and β .

3.3 Reducing the complexity via recursive divide and conquer

The modification through which we are able to drastically decrease the computational complexity of the problem is a recursive divide and conquer strategy.

The number of 2D paths created in each interval at the first iteration is $P_1 = Ne^{\alpha\tau_1}$. The number of possible matches between these paths at the first iteration is $H_1 = (Ne^{\alpha\tau_1})^{\beta_1}$, where $\tau_1 < T$ and $\beta_1 = \beta(\tau_1) < \beta(T)$. At the end of the first iteration, the algorithm chooses which 2D paths are kept in memory, and discards the other ones. The number of paths passed at the following iteration is of the same order of N ($N_1 \simeq N$ meta-objects are created in each interval). Therefore,

$$P_2 = Ne^{\alpha\tau_2} \quad \text{and} \quad H_2 = (Ne^{\alpha\tau_2})^{\beta_2}, \quad (8)$$

where $\tau_2 < \tau_1 < T$ and $\beta_1 < \beta_2 < \beta(T)$. At the n -th iteration,

$$P_n = Ne^{\alpha\tau_n} \quad \text{and} \quad H_n = (Ne^{\alpha\tau_n})^{\beta_n}, \quad (9)$$

where $\tau_n < \tau_{n-1} < \dots < \tau_1 < T$ and $\beta_1 < \beta_2 < \dots < \beta_n \leq \beta(T)$. The crucial point is that $\tau_n \beta_n \ll T \beta(T)$, because we can decide and tune τ_n to tame the exponential explosion of the computational complexity. Moreover,

such strategy allows us balancing the increase of β_i from iteration to iteration with a decrease of τ_i . As a result, we can handle a very large number of objects N for an arbitrarily long interval of time, T , regardless of the intrinsic complexity of the problem (expressed in terms of α and β).

In Figure 5(d) we plot the computational time versus the problem complexity H , with and without the recursive divide and conquer strategy. Thanks to the recursive approach, the computational time is reduced by several orders of magnitude for large values of H . Note that, for small values of H , the non-recursive approach performs better and should be the preferred choice for small datasets. In Figure 5(e) we report a quality indicator (MOTA, see next section for its definition) for several runs on synthetic datasets with different complexities H , and comparing the results obtained with the global optimization with and without recursive scheme. The plot reveals that the two approaches perform similarly in terms of tracking accuracy.

4 VALIDATION WITH SYNTHETIC DATA

We validate our algorithm making use of synthetic datasets.

Synthetic data. We simulate 3D trajectories of flocking birds by adopting a model of self-propelled particles [30]. We use the positions projected in 2D planes directly, rather than generating realistic renderings from them. We do this since we are not interested in testing the performance of the segmentation routine and since it remains hard to predictably simulate the interaction of camera noise and very small objects. Instead, we simulate the errors of the segmentation routine adding white noise directly to the 2D coordinates of the projected objects, in terms of pixel displacements. We also simulate the formation of optical occlusions. For the details concerning the generation of the synthetic data, the reader is addressed to Appendix E. Note though that our simulation still preserves the correspondences between 3D trajectories, the set of 2D projected paths, and the perturbed paths, which can all be used as ground-truth data.

Quality parameters. Let G be the ground-truth set of trajectories and let N_G be the number of trajectories in G . The noisy 2D positions of the ground-truth trajectories are fed to our tracking algorithm, which outputs the set of trajectories O . The two sets G and O are compared, and the quality of the output is evaluated in terms of the following parameters:

- $MOTA$: Multiple Object Tracking Accuracy [31], *i.e.* the ratio of the number of correctly reconstructed 3D positions over the total number of 3D positions;
- G_{90} : the ratio of the number of ground-truth trajectories correctly reconstructed for at least 90%

TABLE 1

Summary of the synthetic datasets used to validate the new tracking software. For each dataset, we report its duration expressed in frames, the number of objects N_G of the ground-truth set G , the number of output trajectories N_O , the value of the parameter ξ , and the values of the quality parameters $MOTA$ and G_{90} .

Synthetic dataset	Duration (frames)	N_G	N_O	ξ	$MOTA$	G_{90}
S_1	125	256	256	0.19	0.999	1
S_2	125	512	512	0.24	0.9989	0.998
S_3	125	1024	1024	0.27	0.9970	0.987
S_4	250	256	257	0.19	0.9975	0.996
S_5	250	512	513	0.24	0.9958	0.988
S_6	250	1024	1030	0.27	0.9931	0.967
S_7	500	256	257	0.19	0.9989	0.988
S_8	500	512	517	0.24	0.9944	0.961
S_9	500	1024	1033	0.27	0.9895	0.928
S_{10}	1000	256	257	0.19	0.9991	0.984
S_{11}	1000	512	526	0.24	0.9873	0.902
S_{12}	1000	1024	1060	0.27	0.9784	0.869

frames over the entire event, over N_G . For example, given an event lasting 100 frames, G_{90} represents the percentage of ground-truth trajectories correctly reconstructed for 90 frames or more).

In the best case scenario – *i.e.* all the ground-truth trajectories are correctly reconstructed – $MOTA=1$ and $G_{90}=1$.

Results on synthetic datasets. Results for several synthetic datasets are shown in Table 1. The quality parameter $MOTA$ is always greater than 0.97, and greater than 0.99 for 9 datasets over 12. The percentage of correctly reconstructed trajectories is greater than 0.786. This percentage grows rapidly, as soon as we consider the trajectories which are reconstructed correctly for more than the 90% of the total duration: $G_{90} \geq 0.869$.

5 TESTS ON EXPERIMENTAL FIELD DATA

We also tested our algorithm using our experimental data of flocking birds and swarming insects acquired on the field, as well as using public benchmark datasets.

Testing the algorithm with our data, we analyzed 12 events of starling flocks [18], and 7 events of swarming midges [32], [33], as summarized in Table 2. Fig. 1 shows the reconstructed trajectories for four events, the bird flocks labelled E_3 and E_6 , and the midge swarms labelled E_{14} and E_{15} – see Table 2. The original video sequences of these four experimental events (in slow-motion, $0.15\times$ slower than the original speed), together with the reconstructed trajectories, are included as Supplemental Material. Clearly, ground-truth trajectories are

not available in the case of experimental data, and – due to the size of our datasets – manual inspection is not feasible, except for a limited number of ambiguous cases. The quality of the reconstructed trajectories is assessed in this case only in terms of trajectory fragmentation. In Table 2, we report the percentage of trajectories longer than the 90% of the duration of the acquisition. The majority of the reconstructed trajectories are of full-length, and trajectory fragmentation is negligible. Such high-quality data have been used to perform the analysis presented in [18] and [32], [33].

To the best of our knowledge, the only public benchmark datasets for 3D-tracking of animal groups are the thermal infrared videos (the raw image sequences with the corresponding sets of ground-truth trajectories) published by Z. Wu and coworkers [34]. We tested our tracking algorithm on the two datasets *Davis08-sparse* and *Davis08-dense*, and we evaluated the output trajectories using the quality parameters defined in [34]: the numbers of Mostly Tracked ($MT \geq 80\%$) trajectories, Mostly Lost ($ML \leq 20\%$) trajectories, track fragmentations (FM) and Identity Switches (IDS), as well as the $MOTA$. In Table 3, we report the quality of our output trajectories compared to the results on the same datasets published by Z. Wu *et al.* [34].

Our results exhibit better values of $MOTA$ and IDS on both datasets, *dense* and *sparse*, revealing that the trajectories are characterized by low values of false positives, identity switches and mismatches. In terms of MT, ML and FM, GReTA performs slightly worse on the *sparse* dataset than MHT and SDD-MHT; on the other hand, when applied to the *dense* dataset, its performance is comparable to the one of the other methods. This implies that MHT and SDD-MHT output a larger percentage of complete trajectories, which nevertheless are characterized by more identity switches and false positives – as revealed by the lower values of $MOTA$ and higher values of IDS. We were not surprised by this, as our tracking algorithm is intentionally designed to discard

TABLE 3

Comparison of the quality of the output trajectories retrieved using GReTA and the ones retrieved using the algorithms MHT and SDD-MHT, as published by Z. Wu *et al.* [34] (see Table IV therein) on the datasets labeled *Davis-08 sparse* and *Davis-08 dense*.

Dataset	Algorithm	MT (%)	ML (%)	FM (#)	IDS (#)	$MOTA$ (%)
<i>sparse</i>	MHT	96.6	0	105	97	64.1
	SDD-MHT	95.2	0	145	126	78.9
	GReTA	83.1	3.4	188	9	82.4
<i>dense</i>	MHT	71.9	2.5	274	355	-32.0
	SDD-MHT	61.1	3.0	454	444	44.9
	GReTA	78.8	2.9	335	8	80.3

TABLE 2

Summary of the field events analyzed with the new tracking software. For each event, we indicate the object type, the estimated number of objects, the duration (in frames and seconds), the acquisition frame-rate, and the percentage of reconstructed trajectories whose length is greater than 90% of the acquisition duration.

Experimental dataset	Object Type	Estimated # Objects	Duration (frames s)	Frame-rate (Hz)	% of trajectories with Length > 90%
<i>E1</i>	birds	179	440 5.50	80	87.0%
<i>E2</i>	birds	551	360 4.50	80	90.2%
<i>E3</i>	birds	365	128 1.60	80	78.6%
<i>E4</i>	birds	120	310 1.82	170	99.2%
<i>E5</i>	birds	50	1000 5.88	170	98.0%
<i>E6</i>	birds	482	761 4.48	170	84.3%
<i>E7</i>	birds	117	500 2.94	170	88.7%
<i>E8</i>	birds	110	661 3.89	170	97.2%
<i>E9</i>	birds	381	960 5.65	170	72.3%
<i>E10</i>	birds	168	300 1.76	170	81.2%
<i>E11</i>	birds	1270	300 1.76	170	87.6%
<i>E12</i>	birds	60	609 3.58	170	89.8%
<i>E13</i>	insects	37	2000 11.76	170	97.1%
<i>E14</i>	insects	332	465 2.73	170	80.2%
<i>E15</i>	insects	115	1000 5.88	170	85.6%
<i>E16</i>	insects	147	1000 5.88	170	84.6%
<i>E17</i>	insects	210	500 2.94	170	82.7%
<i>E18</i>	insects	124	1024 6.02	170	84.0%
<i>E19</i>	insects	633	250 1.47	170	82.3%

false positives, preferring short and correct trajectories to long but incorrect ones.

We believe that this benchmark proves the performance advantages, as well as the flexibility of GReTA to process very diverse experimental data.

6 COMPARISON WITH PRIOR WORK

In order to situate our algorithm in the 3D tracking landscape, an estimate (when explicit information was not published) of the number of tracked objects N and of the temporal duration T (the average trajectory length is used in case of objects entering and leaving the field-of-view) is shown in Fig. 6 for a number of 3D tracking results published in the literature. The points scattered on the plot have been classified according to the field of investigation for which the respective algorithms have been developed: fluid dynamics experiments (■), biological experiments (●), and the experimental data presented in this paper and listed in Table 2 (△ and ▽ for birds and insects, respectively). The numbers next to the symbols correspond to the references to the papers from which T and N have been estimated – see Bibliography. Given that N and T are both valid – but qualitatively different – criteria of evaluation, to compare different methods according to N and T we can use a multi-objective optimization approach, the simplest of which is defining the Pareto frontier [35], [36] in the $\{T, N\}$ plane. We sketched with a dashed line the Pareto frontier for the plotted data-points in the 2D space of N and T . The best

tracking performance is given by the points closest to the frontier.

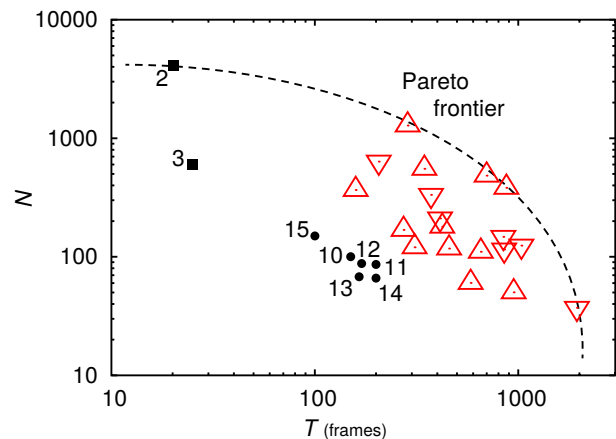


Fig. 6. Comparison of test-cases used for several tracking algorithms, quantified in terms of temporal duration T and estimated number of tracked objects N . The largest datasets processed by different tracking algorithms define the Pareto frontier in the two-dimensional space $\{T, N\}$. The points are classified according to the field of investigation for which the respective algorithms have been developed: fluid dynamics experiments (■), biological experiments (●), and the experimental data processed using GReTA and presented in this paper and listed in Table 2 (△ and ▽ for birds and insects, respectively). The numbers next to the symbols correspond to the references to the papers from which T and N have been estimated – see Bibliography.

The plot clearly shows that fluid dynamics tracking algorithms have been optimized to track large number of tracer particles for short times, and that previous tracking approaches restricted biological experiments to relatively small numbers of animals, even though for longer times. The data processed using GReTA mark the Pareto frontier, assessing the important step forward in terms of performance of the algorithm we proposed. This proved to be suitable for tracking large groups of objects for considerably long durations, without suffering from frequent optical occlusions.

7 CONCLUSIONS

We presented a novel Global and Recursive Tracking Algorithm (GReTA) in three dimensions able to reconstruct uninterrupted trajectories for large numbers of objects and long time intervals, even with frequent optical occlusions. This recursive divide and conquer algorithm is based on the idea of global optimization of the solution – global in space as well as in time. The applicability of a global optimization is limited by the computational complexity, which grows exponentially fast with the time duration of the sequence.

Here we achieve a dramatic reduction of the computational complexity by making use of a recursive divide and conquer strategy, which allows to first optimize the matches globally over shorter temporal intervals, and then iterate to cover the entire temporal sequence. In this way, the computational complexity is drastically reduced while preserving the global scope, permitting to track very large datasets (large in terms of number of objects and of duration of the video acquisition). We further proposed several adaptations making the algorithm robust against wrong or missing links.

We implemented the algorithm; we validated it making use of synthetic data with available ground-truth information; we tested it on new experimental field data of flocking birds and swarming insects; we compared its performance using public benchmark datasets. We showed that the algorithm is capable of reconstructing 3D-trajectories with negligible fragmentation, and that the quality of the trajectories is not affected by the recursive divide and conquer strategy. To the best of our knowledge, the results based on synthetic data and on the public datasets proved the superior performance of the proposed tracking approach compared to other existing methods.

We processed bird flock data, insect swarm data, and bats data, despite these systems being very different from each other: insects in a swarm fly in a very jerky manner and occlude frequently in the images, but for very short times; the flight of birds in a flock is highly coordinated, so that occlusions are typically very long-lasting, and can involve several birds at the time; bats exiting a cave continuously enter and leave the field-of-view. Because of this flexibility, we believe that the GReTA approach can be successfully applied to process the most diverse experimental data.

ACKNOWLEDGMENTS

This work was supported by grants IIT–Seed Artswarm, ERC–StG n.257126, and AFOSR–FA95501010250 (through the University of Maryland). F. Pellacini was partially supported by Intel. We acknowledge the advice of Carlo Lucibello on multi-objective optimization.

REFERENCES

- [1] R.J. Adrian, Particle-Imaging Techniques for Experimental Fluid Mechanics. *Annu. Rev. Fluid Mech.* **23**, 261–304 (1991).
- [2] B. Lüthi, A. Tsinober, and W. Kinzelbach, Lagrangian measurement of vorticity dynamics in turbulent flow. *J. Fluid Mech.* **528**, 87 (2005).
- [3] N.T. Ouellette, H. Xu, and E. Bodenschatz, A quantitative study of three-dimensional Lagrangian particle tracking algorithms. *Exp. Fluids* **40**, 301–313 (2006).
- [4] J.M. Cullen, E. Shaw, and H.A. Baldwin, Methods for measuring the three-dimensional structure of fish schools. *Animal Behaviour* **13**, 4, 534–536 (1965).
- [5] H.J. Dahmen and J. Zeil, Recording and reconstructing three-dimensional trajectories: a versatile method for the field biologist. *Proc. R. Soc. B* **222**, 107–113 (1984).
- [6] H. Pomeroy and F. Heppner, Structure of turning in airborne rock dove (*Columba livia*) flocks. *Auk* **109**, 256–267 (1992).
- [7] N.A. Malik, T. Dracos, and D.A. Papantoniou, Particle tracking velocimetry in three-dimensional flows. *Exp. Fluids* **15**, 279–294 (1993).
- [8] D.H. Doh, D.H. Kim, S.H. Choi, S.D. Hong, T. Saga, and T. Kobayashi, Single-frame (two-field image) 3-D PTV for high speed flows. *Exp. Fluids* **29**, 1 Suppl., S085–S098 (2000).
- [9] J. Willneff, 3D particle tracking velocimetry based on image and object space information. *Int. Arch. Photogrammetry and Remote Sensing and Spatial Inform. Sci.* **34**, 601 (2002).
- [10] Z. Wu, N.I. Hristov, T.L. Hedrick, T.H. Kunz, and M. Betke, Tracking a Large Number of Objects from Multiple Views. *Proc. 12th Int. Conf. Computer Vision (IEEE)*, 1546–1553 (2009).
- [11] D. Zou, Q. Zhao, H.S. Wu, and Y.Q. Chen, Reconstructing 3D motion trajectories of particle swarms by global correspondence selection. *Proc. 12th Int. Conf. Computer Vision (IEEE)*, 1578–1585 (2009).
- [12] H.S. Wu, Q. Zhao, D. Zou, and Y.Q. Chen, Acquiring 3D Motion Trajectories of Large Numbers of Swarming Animals. *Proc. Workshop 12th Int. Conf. Computer Vision (IEEE)*, 593–600 (2009).
- [13] H.S. Wu, Q. Zhao, D. Zou, and Y.Q. Chen, Automated 3D trajectory measuring of large numbers of moving particles. *Optics Express* **19**, 8 (2011).
- [14] Z. Wu, T.H. Kunz, and M. Betke, Efficient Track Linking Methods for Track Graphs Using Network-flow and Set-cover Techniques. *Proc. 24th Conf. Computer Vision and Pattern Recognition (IEEE)*, 1185–1192 (2011).
- [15] Y. Liu, H. Li, and Y.Q. Chen, Automatic Tracking of a Large Number of Moving Targets in 3D. *Proc. 12th European Conf. Computer Vision, Part IV*, 730–742 (2012).
- [16] R. Ardekani, A. Biyani, J.E. Dalton, J.B. Saltz, M.N. Arbeitman, J. Tower, S. Nuzhdin, and S. Tavaré, Three-dimensional tracking and behaviour monitoring of multiple fruit flies. *J. R. Soc. Interface* **10**, 78 (2012).
- [17] A. Cavagna, S.M. Duarte Queirós, I. Giardina, F. Stefanini, and M. Viale, Diffusion of individual birds in starling flocks. *Proc. R. Soc. B* **280**, 20122484 (2013).
- [18] A. Attanasi, A. Cavagna, L. Del Castello, I. Giardina, T.S. Grigera, A. Jelić, S. Melillo, L. Parisi, O. Pohl, E. Shen, and M. Viale, Information transfer and behavioural inertia in starling flocks. *Nature physics* **10**, 9, 691–696 (2014).
- [19] L. Biferale, E. Bodenschatz, M. Cencini, A.S. Lanotte, N.T. Ouellette, F. Toschi, and H. Xu, Lagrangian structure functions in turbulence: A quantitative comparison between experiment and direct numerical simulation. *Phys. Fluids* **20**, 065103 (2008).
- [20] A. D. Straw, K. Branson, T. R. Neumann and M. H. Dickinson, Multi-camera real-time three-dimensional tracking of multiple flying animals. *J. R. Soc. I*, (2010).

- [21] S. Butail, N. Manoukis, M. Diallo, A.S. Yaro, A. Dao, S.F. Traore, J.M. Ribeiro, T. Lehmann, D.A. Paley, 3D tracking of mating events in wild swarms of the malaria mosquito *Anopheles gambiae*. *Engineering in Medicine and Biology Society, EMBC, 2011 Annual International Conference of the IEEE* **75**, 720–723 (2011).
- [22] S. Butail, N. Manoukis, M. Diallo, Ribeiro J. M., T. Lehmann and D. A. Paley, Reconstructing the flight kinematics of swarming and mating in wild mosquitoes. *J. R. Soc. I* **75**, 2624–2638 (2012).
- [23] A. Veeraghavan, M. Srinivasan, R. Chellappa, E. Baird and R. Lamont, Motion based correspondence for 3D tracking of multiple dim objects. *Proc. Int. Conf. Acoustics, Speech and Signal Processing (IEEE)* **2**, (2006).
- [24] Z. Wu, N.I. Hristov, T.H. Kunz and M. Betke, Tracking-reconstruction or reconstruction-tracking? Comparison of two multiple hypothesis tracking approaches to interpret 3D object motion from several camera views. *Motion and Video Computing, 2009. WMVC'09. Workshop on (IEEE)* 1–8, (2009).
- [25] S. Butail, D.A. Paley, 3D reconstruction of fish schooling kinematics from underwater video. *Robotics and Automation (ICRA), 2010 IEEE International Conference on* 2438–2443, (2010).
- [26] D.B. Reid, An algorithm for tracking multiple targets. *Automatic Control, IEEE Transactions on* **24**, 843–854 (1979)
- [27] R. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*, second ed. Cambridge, U.K.: Cambridge University Press, 2003.
- [28] M.L. Fisher, The Lagrangian relaxation method for solving integer programming problems. *Management science* **50**, 1861–1871 (2004).
- [29] CPLEX Optimization Incorporated, *Using the CPLEX Callable Library*, Incline Village, Nevada, 1994.
- [30] W. Bialek, A. Cavagna, I. Giardina, T. Mora, O. Pohl, E. Silvestri, M. Viale, and A. Walczak, Social interactions dominate speed control in poising natural flocks near criticality. *Proceedings of the National Academy of Sciences* , **111**, 20, 7212–7217 (2014).
- [31] B. Keni and S. Rainer, Evaluating multiple object tracking performance: the CLEAR MOT metrics. *EURASIP Journal on Image and Video Processing*, (2008).
- [32] A. Attanasi, A. Cavagna, L. Del Castello, I. Giardina, S. Melillo, L. Parisi, O. Pohl, B. Rossaro, E. Shen, E. Silvestri, and M. Viale, Collective behaviour without collective order in wild swarms of midges, *PLoS Computational Biology* **10**, 7, 1–15 (2014).
- [33] A. Attanasi, A. Cavagna, L. Del Castello , I. Giardina, S. Melillo, L. Parisi, O. Pohl, B. Rossaro, E. Shen, E. Silvestri, M. Viale, Finite-size scaling as a way to probe near-criticality in natural swarms. *Phys. Rev. Lett.* **113**, 238102, (2014).
- [34] Z. Wu, N. Fuller, D. Theriault, M. Betke, A Thermal Infrared Video Benchmark for Visual Analysis. *IEEE Computer Vision and Pattern Recognition Workshops (CVPRW)*, 201–208 (2014).
- [35] V. Pareto, *Cour d'economie politique*, Genève, Switzerland: Librairie Droz-Geneve, 1964 (first edition in 1896). See also http://en.wikipedia.org/wiki/Pareto_efficiency .
- [36] A. Messac, A. Ismail-Yahaya, and C.A. Mattson, The normalized normal constraint method for generating the Pareto frontier. *Struct. Multidisciplinary Opt.* **25**, 2, 86–98 (2003).