

National ICT Australia Technical Reports  
ISSN 1833-9646

Report Number PA006075

**Reporting Leaders and Followers  
Among Trajectories of Moving Point Objects**

Mattias Andersson  
Joachim Gudmundsson  
Patrick Laube  
Thomas Wolle

# Reporting Leaders and Followers Among Trajectories of Moving Point Objects

Mattias Andersson<sup>1</sup>, Joachim Gudmundsson<sup>2</sup>, Patrick Laube<sup>3\*</sup>, and Thomas Wolle<sup>2</sup>

<sup>1</sup> Department of Computer Science, Lund University, Lund, Sweden  
mattias@cs.lth.se

<sup>2</sup> National ICT Australia Ltd\*\*, Locked Bag 9013, Alexandria NSW 1435, Australia  
{joachim.gudmundsson, thomas.wolle}@nicta.com.au

<sup>3</sup> School of Geography and Environmental Science, The University of Auckland,  
Private Bag 92019, Auckland, New Zealand  
p.laube@auckland.ac.nz

**Abstract.** Widespread availability of location aware devices (such as GPS receivers) promotes capture of detailed movement trajectories of people, animals, vehicles and other moving objects, opening new options for a better understanding of the processes involved. In this paper we investigate spatio-temporal movement patterns in large tracking data sets. We present a natural definition of the pattern ‘one object is leading others’, which is based on behavioural patterns discussed in the behavioural ecology literature. Such leadership patterns can be characterised by a minimum time length for which they have to exist and by a minimum number of entities involved in the pattern. Furthermore, we distinguish two granularities (discrete and continuous) of the time axis for which patterns can start and end. For all variants of these leadership patterns, we describe algorithms for their detection, given the trajectories of a group of moving entities. A theoretical analysis as well as experiments show that these algorithms efficiently report leadership patterns. Complementary to those results, we also propose an approximation algorithm for one variant of the pattern and show the hardness/a lower bound on the time complexity for the continuous version of the problem.

**Keywords:** moving point objects, trajectories, movement patterns, leadership, spatio-temporal data structures, computational geometry

## 1 Introduction

Movement is the spatio-temporal process par excellence. Technological advances of location-aware devices, surveillance systems, and electronic transaction networks produce more and more opportunities to trace moving individuals. Consequently, an eclectic set of disciplines including geography [18], data base research [24], animal behaviour research [27], surveillance and security analysis [50, 48], transport analysis [35], and market research [51] shows an increasing interest in movement patterns of various entities moving in various spaces over various times scales.

At the same time traditional geographic analysis suffers from the legacy of cartography’s static perception of the world and is thus generally not suited for the analysis of individual movement trajectories [8, 53], sometimes referred to as geospatial lifelines especially in a GIScience context [44]. Many authors have therefore recently proposed to use geographical (and thus) spatio-temporal data mining as a promising alternative to overcome this methodological shortcoming [15, 46].

As can be seen from the pattern terminology, this research is largely inspired by movement patterns observed in gregarious animals, such as flocking sheep or schooling fish. It follows a strategy to link the proposed patterns as close as possible to observable patterns. The proposed pattern

---

\* Current work is funded by the Swiss National Science Foundation, grant no. PBZH2-110315.

\*\* National ICT Australia is funded by the Australian Government’s Department of Communications, Information Technology, and the Arts and the Australian Research Council through Backing Australia’s Ability and the ICT Research Centre of Excellence programs.

definitions are based on behavioural patterns discussed in the behavioural ecology literature and used for the modelling of realistic movement patterns of agent-based virtual life forms [36, 28].

Even though the leadership patterns in this paper are motivated and investigated with respect to animal behaviour research, their definitions are held generic and are thus applicable to arbitrary types of entities moving in a 2D-space. In general the input is a set of  $n$  moving point objects  $e_1, \dots, e_n$  whose locations are known at  $\tau$  consecutive time-steps  $t_1, \dots, t_\tau$ , that is, the trajectory of each object is a polygonal line that can self-intersect. For brevity, we will call moving point objects *entities* from now on. We assume that the movement of an entity from its position at a time-step  $t_j$  to its position at the next time-step  $t_{j+1}$  is described by the straight-line segment between the two coordinates, and that the entity moves along the segment with constant velocity.

The paper is organised as follows. This introductory section links previous research on movement patterns similar to our leadership with the latest related research. Section 2 defines our notion of leadership and features definitions and preliminaries. In Section 3 and 4, we present exact algorithms for the detection of leadership. Section 5 describes an approximation algorithm and Section 6 a lower bound. Then, in Section 7 we present experimental results, and we conclude the paper with future research and final remarks in Section 8.

## 1.1 Movement Patterns Revisited

Precursory to this research Laube and colleagues proposed the REMO framework (RElative MOtion) which defines similar behaviour in groups of entities [37–39]. They defined a collection of movement patterns based on similar movement properties such as speed, acceleration, or movement direction. Laube et al. [40] extended the framework by not only including movement properties, but also location itself. They defined several movement patterns, including flock (co-ordinately moving close together), trend-setter (anticipating a move of others), leadership (spatially leading a move of others), convergence (converging towards a spot) and encounter (meeting at a spot), and gave algorithms to compute them efficiently. Later Gudmundsson et al. [22] considered the same problems and extended the algorithmic results by primarily focusing on approximation algorithms – “Any exact values of  $m$  and  $r$  hardly have a special significance – 20 caribou meeting in a circle with radius 50 meters form as interesting a pattern as 19 caribou meeting in a circle with radius 51 meters.”

Benkert et al. [4] and Gudmundsson and van Kreveld [21] only recently revisited the flock pattern and gave a more generic definition that bases purely on the geometric arrangement of the moving entities and thus excludes the need of an analytical space as with the initial definition of the patterns [37, 40]. The present paper shall achieve a similar redefinition of the leadership pattern and its flavours and shall provide a generic, geometric definition for them as well as efficient algorithms for their detection. The aim is to define movement patterns that are sound and allow for the patterns to be detected and reported efficiently.

The model used in the REMO framework considers each time-step separately, that is, given  $m \in \mathcal{N}$  and  $r > 0$  a flock is defined by at least  $m$  entities within a circular region of radius  $r$  and moving in the same direction at some point in time. Benkert et al. [4] argued that this is not enough for many practical applications, e.g. a group of animals may need to stay together for days or even weeks before it is defined as a flock. They proposed the following definition of a flock:

**Definition 1.**  $(m, k, r)$ -flock<sub>A</sub> - Given a set of  $n$  trajectories where each trajectory consists of  $\tau$  line segments, a flock in a time interval  $I = [t_i, t_j]$ , where  $j - i + 1 \geq k$ , consists of at least  $m$  entities such that for every point in time within  $I$  there is a disk of radius  $r$  that contains all the  $m$  entities. Note that  $m, k \in \mathcal{N}$  and  $r > 0$  are given constants.

We will use a similar model when defining the leadership patterns, see Section 2. Using this model, Gudmundsson and van Kreveld [21] recently showed that computing the longest duration flock and the largest subset flock is NP-hard to approximate within a factor of  $\tau^{1-\varepsilon}$  and  $n^{1-\varepsilon}$  respectively. In the same model, Benkert et al. [4] described an efficient approximation algorithm for reporting and detecting flocks, where they let the size of the region deviate slightly from

what is specified. Approximating the size of the circular region with a factor of  $\Delta > 1$  means that a disk with radius between  $r$  and  $\Delta r$  that contains at least  $m$  objects may or may not be reported as a flock while a region with a radius of at most  $r$  that contains at least  $m$  entities will always be reported. Their main approach is a  $(2 + \varepsilon)$ -approximation with running time  $T(n) = O(\tau n k^2 (\log n + 1/\varepsilon^{2k-1}))$ . Note that even though the dependency on the number of entities (namely  $n$ ) is small, the dependency on the duration of the flock pattern (namely  $k$ ) is exponential. Al-Naymat et al. [2] handle the problem of considering many entities and long-duration patterns by using a preprocessing step where the number of dimensions (i.e. time-steps) is reduced by random projection.

## 1.2 Related Work

There is ample research on moving object databases (MOD) [58, 65, 66, 26]. Whereas most database research on MOD focuses on data structures, indexing and efficient querying techniques for moving objects [24, 1, 17, 25], only recently the potential of data mining for movement patterns has been acknowledged [55, 32, 33]. For example, Du Mouza and Rigaux propose mobility patterns that describe sequences of moves in a discrete 2D-space [12]. In a GIScience context, activity related movement patterns have been researched, often with respect to improving location-based services (LBS). Dykes and Mountain search episodes expressing distinctive characteristics of movement, including absolute speed, direction, sinuosity and measurements of their variations [14]. Smyth presents a data mining algorithm that assigns predefined activities to segments of trajectories by analysing some measurable motion descriptors, such as speed, heading, acceleration [59]. Frihida et al. mine origin-destination data for predefined activity patterns in a transport demand modelling study [19]. A common approach in database research is to take an existing spatial query type and then study its generalisations to spatio-temporal data. An example of this is the recent work on continuous  $k$ -nearest neighbour querying over mobile data [47, 67]. The focus within data mining research is to design techniques to discover new patterns in large repositories of spatio-temporal data. For example, Mamoulis et al. [43] mine periodic patterns moving between objects and Ishikawa et al. [29] mine spatio-temporal patterns in the form of Markov transition probabilities. More recently Verhein and Chawla [62] used association rule mining for patterns such as, sinks, sources, stationary region, and thoroughfares.

Spatio-temporal proximity of entities is a reasonable first premise for many situations that assume interactions between individuals. One obvious analytical toolset to uncover proximity patterns in individual trajectories is clustering. Even though the spatio-temporal nature of movement data adds additional complexity to clustering procedures, there have been some successful approaches for clustering trajectories [57, 10, 42]. However, spatio-temporal co-presence does not explicitly include the idea of interactions within individuals. Relations such as ‘leading’, ‘following’ or ‘setting a trend’ cannot be investigated by pure clustering alone.

When analysing sets of individual trajectories, interestingly, little research has investigated relations like ‘one entity leads others’. Exceptions are the work by Shirabe [56], and also the work by Börner and Penumarthy who investigate the formation and diffusion of social groups in a virtual 3D-environment [6]. Even though they show an interest in participants forming groups and participants following a leader, their major interest is group coherence, expressed as proximity in space. Their groups may be focused (entities are close to each other) or unfocused (entities are spread out), and may interchange between those stages. The depiction of the influence of positions of entities on the diffusion of others bases in this approach purely on map-based visualisation of entity arrangements. Another exception investigating leadership patterns is the REMO approach [37–40], discussed earlier.

Along with humans, also animals interact socially to gain from coordination of their behaviour [34, 9]. For foraging animals coordination behaviour is likely to reflect a trade-off between the risk of starvation and predation [52, 45]. Rands et al. [52] illustrated the spontaneous emergence of leaders and followers using a simulation model reproducing the decision process of a pair of foraging animals, balancing their energetic states. The pair produced spontaneous temporary leaders and followers following a simple rule of thumb: ‘I should forage if either my reserves have

fallen below a certain threshold value, or my partner chooses to forage.’ In a very similar way, so-called local enhancement describes the behaviour ‘where the behaviour of one animal draws the attention of a second animal to a particular stimulus in a local environment’ [36]. For example, individual guppies (*Poecilia reticulata*) observe physical and behavioural characteristics of potential leaders in order to decide whom they shall follow for shoaling [36]. Leadership behaviour has also been studied in gray wolves (*Canis lupus*) [49]. In this context ‘leading’ has been quantified as a set of behavioural patterns observable in the field, involving ‘scent-marking’, ‘frequency and time in the lead while packs were travelling or hunting’ (frontal leadership), ‘initiation of pack behaviour’; and ‘non-frontal leadership’ (occurred when a wolf not in the lead broke ranks and led the pack in a new direction or activity).

The idea and the term of leadership have been used in several different contexts in the field of animal behaviour research, see Dumont et al. [13] for an overview. In general, one can distinguish two different readings:

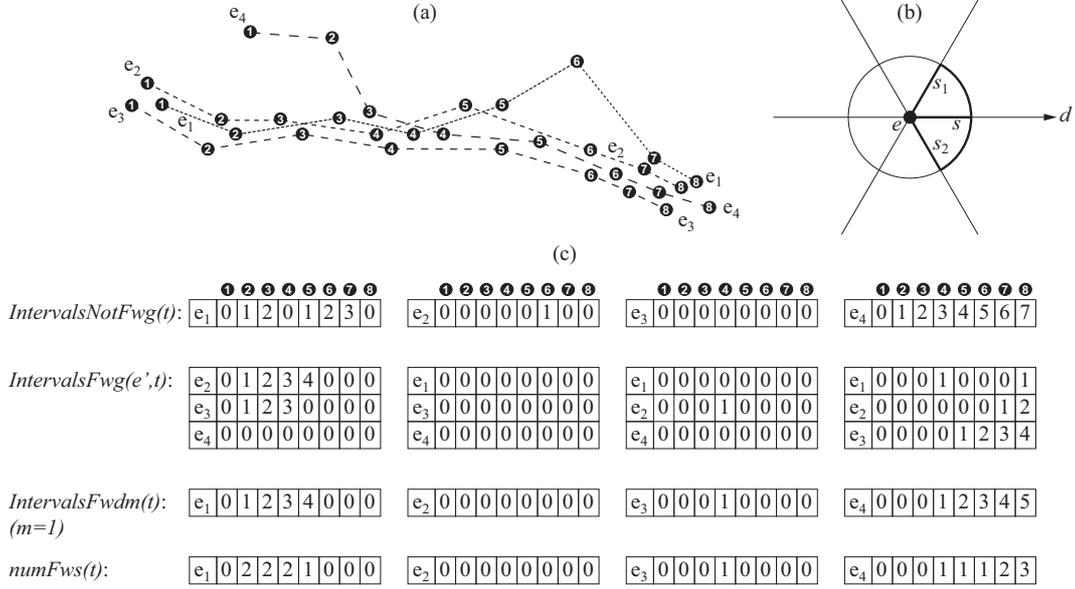
1. (i) ‘the event or process of one entity initiating a group movement (e.g. [49, 13, 7, 41])’ Leading in this sense is an active behaviour, referring to individuals that consistently initiate displacement of the group it belongs to. For example, Dumont et al. [13] found that in a group of 15 grazing heifers the same individual was reported to lead the group to new feeding places in 48% of all group movements.
2. (ii) ‘the event or process of one entity in front, leading a group movement (e.g. [7, 23])’ Intrinsic walking speeds may cause faster individuals (‘speeders’) to become leaders in the sense of individuals at the front, with respect to the direction of propagation [23]. Gueron and Levin model the spatial constellation ‘in front of’ as a function of the relative position with respect to the averaged position of its neighbours within a given range. Even though it has been found for grazing animals that leaders may guide a group being in front or chasing from behind, animals in front are considered to be more relevant to determine where the group will graze [13].

The use of the geometrical arrangement of moving entities has furthermore a long tradition for realistically modelling group behaviour, be it in animal behaviour science [23] or in the animation industry [54]. Most prominent is the flocking model implemented in NetLogo [61, 64], which mimics the flocking of birds [63]. The moving agents dynamically coordinate their movement basing on rules on alignment (turning in order to adopt direction of nearby agents), separation (turning to avoid getting too close to nearby agents) and cohesion (move towards other nearby agents). This model explicitly excludes the idea of an individual leading the others, but involves identical agents, each following the same set of rules. The basic model includes a maximal distance of vision  $r$  and 360 degree field of view. However, the IN-CONE primitive allows using a cone of vision, a most interesting concept with respect to the investigation of further structure in flocking entities that can, for example, be seen in V-shaped flocks as with migrating geese. Such front priority is also often used for agent-based models of schooling fish, where only individuals in front are candidates as interacting neighbours [28]. Inada’s and Kawachi’s model uses a wide-angled cone of perception, directed in movement direction and thus omitting a blind region behind the fish (Figure 1, Figure 2 in [28]). Just as NetLogo’s flocking model, their fish school model includes reaction fields that are limited by a given radius. Three nested reaction fields determine the agent’s reaction on neighbours: an inner-most repulsive-orientation field (push back to close entities), a middle parallel-orientation field (adopt neighbour’s direction), and an out-most attractive-orientation field (move towards neighbour). Jadbabaie et al. give a theoretical explanation for the spontaneous coordination of agents despite the absence of a centralised coordination and just following a simple nearest neighbour rule [30]. However, in an extension they also investigate the influence of a leader in their system.

## 2 Definition of Leadership and Problem Statement

We consider  $n$  entities moving in the two dimensional plane during the time interval  $[t_1, t_\tau]$ , see Figure 1(a) for an example. The infinite set  $T_p$  of time-points is defined as  $T_p = \{t \mid t \in [t_1, t_\tau]\}$ ,

and the set  $T_s$  of time-steps is the set of discrete time-points given as input, i.e.  $T_s = \{t_1, \dots, t_\tau\}$ . We specify open and closed time intervals by  $(t_x, t_y)$  and  $[t_x, t_y]$ , respectively. A unit-time-interval is an open interval  $I$  between two consecutive time-steps, i.e.  $I = (t_{x-1}, t_x)$ , for a time-step  $t_x$  with  $x > 1$ .



**Fig. 1.** (a) A set of 4 entities moving from left to right over 7 unit-time-intervals, i.e. over 8 time-steps. (b) Illustrating the definition of the front-region as the disc-segment within bold lines. (c) The follow-arrays of the four entities, where we use the front-region as depicted in (b) and  $\alpha = \beta$ .

For describing our leadership patterns, we need a couple of parameters specifying these patterns. More specifically, we assume that we are given numbers  $m$  (specifying the size of a pattern, i.e. the minimum number of entities involved in a pattern),  $k$  (the minimum temporal length of a pattern), a radius  $r$  (spatial size of a pattern) and angles  $\alpha$  and  $\beta$  (also determining spatial characteristics of a pattern). We consider them as constants during the rest of the paper, i.e. we will not carry them along as parameters of functions or other notations.

At time-point  $t_x$ , an entity  $e_j$  is located at a position with coordinates  $xpos(e_j, t_x)$  and  $ypos(e_j, t_x)$ . As we do not have spatial information of an entity between two time-steps we make the following assumption for the remainder of this paper.

**Assumption 1** *We assume that all entities move between two consecutive time-steps with constant direction and constant velocity.*

The same assumption has been used in earlier papers [4]. It enables us to interpolate the positions of entities between time-steps. Even though we have no bound on the accuracy of this interpolation compared to the real positions of the entities, it appears to be a reasonable approach when tackling our leadership problems.

Suppose we are given an entity  $e_j$  at time-point  $t$  with  $t_{x-1} < t < t_x$  for  $t_x \in T_s$ . We say  $e_j$  is heading into *direction*  $d$  where  $d$  is an angle in  $[0, 2\pi)$  that is specified by the line segment  $e_j$  is moving along between time-steps  $t_{x-1}$  and  $t_x$ . (If  $e_j$  does not move between  $t_{x-1}$  and  $t_x$  then we define  $d$  to be the direction of the line segment  $e_j$  is moving along between the time-steps  $t_{x-2}$  and  $t_{x-1}$ . If no such time-steps exist, then we define  $d := 0$ .) The difference between two directions  $d_1$  and  $d_2$  is denoted by  $\|d_1 - d_2\|$ , which is an angle in  $[0, \pi]$ . We declare the direction of an entity at a time-step  $t_x$  to be *undefined*, because at time-steps an entity might change its direction.

However, the *direction* of an entity  $e_i$  at a time-step  $t_x$  with respect to  $(t_{x-1}, t_x)$  is the direction  $e_i$  is heading to at any time-point in  $(t_{x-1}, t_x)$ . Therefore, when considering time intervals with certain properties of entities depending on direction, we implicitly exclude time-steps from those intervals in the remainder of this paper.

Given an entity  $e$  and a time-point  $t \notin T_s$ , we define the *front-region* of  $e$  at time  $t$  in the following way. Consider the disk  $C$  with radius  $r$  centred at  $(xpos(e, t), ypos(e, t))$ . Furthermore, consider three line segments  $s_1$ ,  $s_2$  and  $s$  of length  $r$ , all having one end point at  $(xpos(e, t), ypos(e, t))$ . Segment  $s$  points in the direction  $d$  that  $e$  is heading to at time  $t$ , and segments  $s_1$  and  $s_2$  are the well defined segments forming angles of  $\frac{\alpha}{2}$  and  $-\frac{\alpha}{2}$  with  $s$ , respectively. The part of the disk  $C$  that contains  $s$  and is bounded by the segments  $s_1$  and  $s_2$  is the *front-region*, see Figures 1(b) and 2. We denote this wedge-shaped region by  $front(e)$  at time  $t$ . An entity  $e_j$  is said to be *in front* of an entity  $e_i$  at time  $t \notin T_s$  if and only if  $e_j \in front(e_i)$  at time  $t$ .

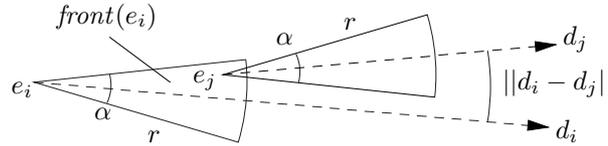
**Definition 2.** Let  $d_i$  and  $d_j$  be the directions of the entities  $e_i$  and  $e_j$  at time  $t \notin T_s$ , respectively. Entity  $e_i$  is said to follow entity  $e_j$  at time  $t$ , iff  $e_j \in front(e_i)$  at time  $t$  and  $\|d_i - d_j\| \leq \beta$ .

That is, entity  $e_i$  has to be in front of  $e_j$  and  $e_j$  and  $e_i$  are heading into directions that differ by at most  $\beta$ . An entity  $e_j$  is said to follow entity  $e_i$  at time  $[t_x, t_y]$  for time-points  $t_x, t_y$ , if and only if  $e_i$  follows  $e_j$  at time  $t$  for all time-points  $t \in [t_x, t_y] \setminus T_s$ .

**Definition 3.** An entity  $e_i$  is said to be a leader at time  $[t_x, t_y]$  for time-points  $t_x, t_y$ , if and only if  $e_i$  does not follow anyone at time  $[t_x, t_y]$ , and  $e_i$  is followed by sufficiently many entities at time  $[t_x, t_y]$ . If there is an entity that is a leader of at least  $m$  entities for at least  $k$  time units, we have a leadership pattern.

See Figure 2 for an example of some notations.

As mentioned above, the direction of an entity is not defined at a time-step. Therefore, in the remainder of this paper we exclude time-steps from any sets or intervals that specify any properties such as direction, being in front, following and leading. e.g. if  $e_i$  is following  $e_j$  at time  $I$ , then we mean  $I \setminus T_s$ .



**Fig. 2.** The front regions of  $e_i$  and  $e_j$  as wedges of edge length  $r$  and apex angle  $\alpha$ . Entity  $e_j$  is in front of  $e_i$ . The entities are heading into directions  $d_i$  and  $d_j$ , respectively. If  $\|d_i - d_j\| \leq \beta$  then  $e_i$  follows  $e_j$ .

*Example 1.* Consider the entities in Figure 1(a) where we use a front-region as depicted in Figure 1(b). We see that  $e_2$  is following  $e_1$  at time  $(t_1, t_5)$ ,  $e_1$  is not following any other entity at time  $(t_1, t_3)$  and  $(t_4, t_7)$  and hence  $e_1$  is a leader of  $e_2$  at time  $(t_1, t_3)$  and  $(t_4, t_5)$ .

In the remainder of this section, we consider two entities  $e_i$  and  $e_j$  and two consecutive time-steps  $t_{x-1}$  and  $t_x$ . The next lemma tells us that if we want to check whether an entity is following any other entity during the entire interval  $(t_{x-1}, t_x)$ , we only have to check this at the two end points with respect to  $(t_{x-1}, t_x)$ . The lemma is rather intuitive and can be proven with very much the same ideas as in the proof of Lemma 2 in [4].

**Lemma 1.** Let  $e_i$  and  $e_j$  be two entities, and let  $t_{x-1}$  and  $t_x$  be two consecutive time-steps. If  $e_i$  follows  $e_j$  at time-points  $t_y$  and  $t_z$  with  $t_{x-1} < t_y \leq t_z < t_x$  then under Assumption 1,  $e_i$  follows  $e_j$  at any time-point  $t \in [t_y, t_z]$ .

Note that the lemma is also true for  $t_{x-1} = t_y$  and  $t_z = t_x$ , however, the directions of the entities at these time-points are with respect to  $(t_{x-1}, t_x)$ . Therefore, the time that an entity  $e_i$  follows another entity  $e_j$  between two consecutive time-steps  $t_{x-1}$  and  $t_x$  is a single subinterval of  $[t_{x-1}, t_x]$ , and such an interval can be computed in a straightforward manner. If the directions of  $e_i$  and  $e_j$  differ by more than  $\beta$  then one entity cannot follow the other. Otherwise, we have to compute if and when  $e_j$  is in front of  $e_i$ . As  $e_j$ 's and also  $front(e_i)$ 's movement is described by straight lines, this can be done in constant time as stated in the next lemma.

**Lemma 2.** *Given two entities  $e_i$  and  $e_j$  and two time-steps  $t_{x-1}$  and  $t_x$ , we can compute in constant time the subinterval of  $[t_{x-1}, t_x]$  for which  $e_i \in front(e_j)$  and for which  $e_j$  follows  $e_i$ , under Assumption 1.*

In the remainder of this paper, we require that Assumption 1 holds true, and we regularly use Lemmas 1 and 2.

## 2.1 The Problems We Consider

A leadership pattern exists if there is an entity that is a leader of sufficiently many entities over a long enough series of time-steps or time-points. Such a pattern is characterised by two values  $m$  which is the size of the set of followers, and  $k$  which is the length of a pattern. As mentioned in related work [4, 22] specifying exactly which of the patterns should be reported is often a subject for discussion. For instance, a leadership pattern of length  $k + 1$  implies the existence of two leadership patterns of length  $k$ . However, the pattern of length  $k + 1$  might be more interesting to report from a practical point of view. Therefore, we consider the following problems where we assume that  $m$  and  $k$  are constants.

- LP-report-all: For each entity  $e$ , report all time intervals where  $e$  is a leader of at least  $m$  entities for at least  $k$  time units.
- LP-max-length: Compute the length of a longest leadership pattern of size at least  $m$ , i.e. compute the largest value  $k^{max}$  such that there is an entity  $e$  that is a leader of at least  $m$  entities for  $k^{max}$  time units.
- LP-max-size: Compute the size of a largest leadership pattern of length at least  $k$ , i.e. compute the largest value  $m^{max}$  such that there is an entity  $e$  that is a leader of  $m^{max}$  entities for at least  $k$  time units.

All these problems come in four different flavours which are combinations of the granularity of the time axis (discrete vs. continuous) and the consistency of the set of followers (varying vs. non-varying). More specifically, we consider each of the problems in a discrete case, where patterns (and follow behaviour) can only start and end at the discrete time-steps, and in the continuous case, where patterns can start and end at any time-points. As discussed above, the data for the continuous case relies on Assumption 1. The other variation concerns the set of followers. If there is a subset  $S$  of entities such that for each time-point of the duration of the pattern all entities in  $S$  follow the leader (there may be additional followers as well at some time-points), then we call this a non-varying (subset) leadership pattern. In contrast to this, if we allow the subset of followers to change from one unit-time-interval to the next during the duration of the pattern (some entities may drop out, other may join in), then we call such a pattern a varying (subset) leadership pattern, as long as always at least  $m$  entities are following at each unit-time interval of the pattern.

## 2.2 Our Contribution

Table 1 summarises the time complexity bounds of our algorithms for three considered problems for the four different models as stated by lemmas in this article. The space complexity is in all cases  $O(n\tau)$  which is linear in the size of the input. For the continuous case, we have a lower bound of  $\Omega(n^2\tau)$  on the time complexity in the model of [16], see Section 6. In Section 5, we show how a  $(1 + \varepsilon)$ -approximation of the discrete problems with varying followers can be computed in  $O(\frac{n\tau}{\varepsilon} \log^2 n)$  time.

model	discrete	continuous
non-varying	$O(n^2\tau)$ time (Lemma 4, 5, 6)	$O(n^2\tau \log n)$ time (Lemma 13, 14, 16)
varying	$O(n^2\tau)$ time (Lemma 7, 8, 9)	$O(n^2\tau \log n)$ time (Lemma 17, 18, 19)

**Table 1.** An overview over the running times of all leadership problems.

### 3 Exact Algorithms for the Discrete Case

In the discrete case, patterns can only start and end at time-steps. We first describe arrays storing information about the follow behaviour of the entities with respect to a fixed entity  $e_i$ . Later, these arrays will be used to solve our leadership problems.

#### 3.1 Getting Ready – Computing Follow-Arrays for an Entity $e_i$

For an entity  $e_i$  to determine whether it is a leader at the time  $(t_x, t_y)$ , we need to know whether  $e_i$  is not following any other entity and whether  $e_i$  is followed by sufficiently many entities at  $(t_x, t_y)$ . We consider  $e_i$  at this time as a potential leader, and we compute a couple of follow-arrays called ‘ $IntervalsNotFwg(t_x)$ ’, ‘ $IntervalsFwg(t_x, e_j)$ ’, ‘ $IntervalsFwd_m(t_x)$ ’ and ‘ $NumFws(t_x)$ ’. The first three arrays store the number of consecutive unit-time-intervals that there is a certain follow-behaviour. In contrast to this, the fourth arrays stores numbers of entities with a certain follow-behaviour.

***IntervalsNotFwg***: (short for ‘the number of unit-time-intervals  $e_i$  is not following at  $t_x$ ’) The array  $IntervalsNotFwg(t_x)$  is a one dimensional array storing nonnegative integers. Such an integer for time-step  $t_x$  specifies for how many past consecutive unit-time-intervals (the last one ending at  $t_x$ )  $e_i$  is not following any other entity. That is, if  $IntervalsNotFwg(t_x) = y$ , then  $e_i$  is not following any other entity during the time interval  $(t_{x-y}, t_x)$ . This interval covers exactly the  $y$  past consecutive unit-time-intervals. Formally  $IntervalsNotFwg(t_x)$  is defined as follows:

$$IntervalsNotFwg(t_x) = \begin{cases} 0 & \Leftrightarrow \exists j, j \neq i : e_i \text{ follows } e_j \text{ during} \\ & \text{the unit-time-interval } (t_{x-1}, t_x) \\ IntervalsNotFwg(t_{x-1}) + 1 & \text{otherwise} \end{cases}$$

To compute the values of the  $IntervalsNotFwg$ -array, we use two nested loops. The outer loop runs from  $t_x = t_2, \dots, t_\tau$  (we start at  $t_x = 2$  and set  $IntervalsNotFwg(t_1) := 0$ ). The inner loop ranges over  $e_j = e_1, \dots, e_n$  and  $e_j \neq e_i$ . After each round of the inner loop we update  $IntervalsNotFwg(t_x)$  according to whether we found an entity  $e_j$  such that  $e_i$  follows  $e_j$  at time  $(t_{x-1}, t_x)$ . According to Lemma 2 each such single test can be done in constant time.

***IntervalsFwg***: (short for ‘the number of unit-time-intervals  $e_i$  is followed by  $e_j$  at  $t_x$ ’) The array  $IntervalsFwg(t_x, e_j)$  is a  $(\tau \times n - 1)$  matrix storing nonnegative integers specifying for how many past consecutive unit-time-intervals (the last one ending at  $t_x$ )  $e_j$  is following  $e_i$  (with  $e_j \neq e_i$ ).

$$IntervalsFwg(t_x, e_j) = \begin{cases} IntervalsFwg(t_{x-1}, e_j) + 1 & \Leftrightarrow e_j \text{ follows } e_i \text{ at time } (t_{x-1}, t_x) \\ 0 & \text{otherwise} \end{cases}$$

Filling the  $IntervalsFwg$ -array with the right values can also be done with two nested loops, one outer loop for  $t_x = t_2, \dots, t_\tau$  and one inner loop for  $e_j = e_1, \dots, e_n$  and  $e_j \neq e_i$  – initially set

$IntervalsFwg(t_1, e_j) := 0$ . We test whether  $e_j$  follows  $e_i$  at the unit-time-interval  $(t_{x-1}, t_x)$ , and if so, we update  $IntervalsFwg(t_x, e_j)$ .

**$IntervalsFws_m$ :** (short for ‘the number of unit-time-intervals  $e_i$  has at least  $m$  non-varying followers at  $t_x$ ’) The array  $IntervalsFwd_m(t_x)$  is a one-dimensional array storing integers specifying for how many consecutive past unit-time-intervals (the last one ending at  $t_x$ ) there are at least  $m$  entities following entity  $e_i$ .

$$IntervalsFwd_m(t_x) = \begin{cases} IntervalsFwd_m(t_{x-1}) + 1 & \Leftrightarrow \exists S \subseteq \{e_1, \dots, e_{i-1}, e_{i+1}, \dots, e_n\} : \\ & |S| \geq m \wedge \forall e_j \in S : \\ & e_j \text{ follows } e_i \text{ at time } (t_{x-1}, t_x) \\ 0 & \text{otherwise} \end{cases}$$

Given the array  $IntervalsFwg$ , computing the  $IntervalsFwd_m$ -array can be done by counting in each column of  $IntervalsFwg$  (if we imagine the array  $IntervalsFwg$  to be arranged to have  $\tau$  columns and  $n - 1$  rows) the number of entities following  $e_i$  at the current time-step.

**$NumFws$ :** (short for ‘the number of followers of  $e_i$  at  $t_x$ ’) Another array is  $NumFws(t_x)$  which is a one-dimensional array storing integers specifying how many entities are following entity  $e_i$  at time  $(t_{x-1}, t_x)$ .

$$NumFws(t_x) = \text{number of entities following } e_i \text{ at time } (t_{x-1}, t_x)$$

Again, counting in each row of  $IntervalsFwg$  the number of entities following  $e_i$  at the current time-step yields the corresponding value of the  $NumFws$  array.

From the above discussion on the corresponding arrays, we conclude with the following lemma.

**Lemma 3.** *The  $IntervalsNotFwg$ ,  $IntervalsFwg$ ,  $IntervalsFwd_m$  and  $NumFws$ -arrays for an entity  $e_i$  can be computed in  $O(n\tau)$  time and space.*

*Example 2.* Consider the entities in Figure 1(a) where we use a front-region as depicted in Figure 1(b). Figure 1(c) shows four columns (one for each entity) of follow-arrays. To fill the arrays  $IntervalsNotFwg$  and  $IntervalsFwg$ , we need the trajectories and the front-regions. Once that is done, the arrays  $IntervalsFwd_m$  and  $NumFws$  can be computed according to their definition.

### 3.2 Solving LP Problems with a Non-varying Subset of Followers

#### LP-report-all

In the discrete leadership version we assume that patterns can only start and end at time-steps  $T_s = \{t_1, \dots, t_\tau\}$ . To determine whether an entity  $e_i$  is a leader of a non-varying-subset of followers, we use the arrays  $IntervalsNotFwg$  and  $IntervalsFwg$ , and we combine their information to determine whether  $e_i$  is a leader. To that end we look for time-steps  $t_x$  such that  $IntervalsNotFwg(t_x) \geq k$ . For each such time-step  $t_x$ , we inspect the array  $IntervalsFwg(t_x, e_j)$  for  $j = 1, \dots, n$  and  $j \neq i$ , and we count the number of times that  $IntervalsFwg(t_x, e_j) \geq k$ . Let  $m(k)$  denote this number. Now we can report  $e_i$  as a leader for every time-step  $t_x$  for which  $m(k) \geq m$ . As we only need to traverse our arrays at most once, this can be done on  $O(n\tau)$  time.

*Example 3.* Let  $k = 1$  and  $m = 2$ . Looking at the follow-arrays of entity  $e_1$  in Figure 3, we see (shaded region) that  $e_1$  is not following anyone, but is followed by 2 entities, and this happens for at least  $k = 1$  unit-time-intervals at the time-steps 2 and 3. Hence, we would report two leadership-patterns with  $e_1$  as leader.

So far, we have seen that we can compute in  $O(n\tau)$  time and space at which time-steps an entity  $e_i$  is a leader. To find all leadership patterns amongst a set of entities we test any entity individually. As we only have to store one instance of each array at a time (i.e. we can reuse memory for different entities), we can conclude with the following lemma.

**Lemma 4.** *Reporting all non-varying-subset leadership patterns of size at least  $m$  and length at least  $k$ , amongst  $n$  trajectories over  $\tau$  time-steps can be done in  $O(n^2\tau)$  time and  $O(n\tau)$  space.*

<i>IntervalsNotFwg(t):</i>	<table style="border-collapse: collapse; text-align: center;"> <tr><td></td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td><td>8</td></tr> <tr><td>e<sub>1</sub></td><td>0</td><td>1</td><td>2</td><td>0</td><td>1</td><td>2</td><td>3</td><td>0</td></tr> </table>		1	2	3	4	5	6	7	8	e <sub>1</sub>	0	1	2	0	1	2	3	0	<table style="border-collapse: collapse; text-align: center;"> <tr><td></td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td><td>8</td></tr> <tr><td>e<sub>2</sub></td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td></tr> </table>		1	2	3	4	5	6	7	8	e <sub>2</sub>	0	0	0	0	0	1	0	0	<table style="border-collapse: collapse; text-align: center;"> <tr><td></td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td><td>8</td></tr> <tr><td>e<sub>3</sub></td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr> </table>		1	2	3	4	5	6	7	8	e <sub>3</sub>	0	0	0	0	0	0	0	0	<table style="border-collapse: collapse; text-align: center;"> <tr><td></td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td><td>8</td></tr> <tr><td>e<sub>4</sub></td><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td></tr> </table>		1	2	3	4	5	6	7	8	e <sub>4</sub>	0	1	2	3	4	5	6	7
	1	2	3	4	5	6	7	8																																																																				
e <sub>1</sub>	0	1	2	0	1	2	3	0																																																																				
	1	2	3	4	5	6	7	8																																																																				
e <sub>2</sub>	0	0	0	0	0	1	0	0																																																																				
	1	2	3	4	5	6	7	8																																																																				
e <sub>3</sub>	0	0	0	0	0	0	0	0																																																																				
	1	2	3	4	5	6	7	8																																																																				
e <sub>4</sub>	0	1	2	3	4	5	6	7																																																																				

<i>IntervalsFwg(e',t):</i>	<table style="border-collapse: collapse; text-align: center;"> <tr><td>e<sub>2</sub></td><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>e<sub>3</sub></td><td>0</td><td>1</td><td>2</td><td>3</td><td>0</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>e<sub>4</sub></td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr> </table>	e <sub>2</sub>	0	1	2	3	4	0	0	0	e <sub>3</sub>	0	1	2	3	0	0	0	0	e <sub>4</sub>	0	0	0	0	0	0	0	0	<table style="border-collapse: collapse; text-align: center;"> <tr><td>e<sub>1</sub></td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>e<sub>3</sub></td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>e<sub>4</sub></td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr> </table>	e <sub>1</sub>	0	0	0	0	0	0	0	0	e <sub>3</sub>	0	0	0	0	0	0	0	0	e <sub>4</sub>	0	0	0	0	0	0	0	0	<table style="border-collapse: collapse; text-align: center;"> <tr><td>e<sub>1</sub></td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>e<sub>2</sub></td><td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>e<sub>4</sub></td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr> </table>	e <sub>1</sub>	0	0	0	0	0	0	0	0	e <sub>2</sub>	0	0	0	1	0	0	0	0	e <sub>4</sub>	0	0	0	0	0	0	0	0	<table style="border-collapse: collapse; text-align: center;"> <tr><td>e<sub>1</sub></td><td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td><td>1</td></tr> <tr><td>e<sub>2</sub></td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>1</td><td>2</td></tr> <tr><td>e<sub>3</sub></td><td>0</td><td>0</td><td>0</td><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td></tr> </table>	e <sub>1</sub>	0	0	0	1	0	0	0	1	e <sub>2</sub>	0	0	0	0	0	0	1	2	e <sub>3</sub>	0	0	0	0	1	2	3	4
e <sub>2</sub>	0	1	2	3	4	0	0	0																																																																																																								
e <sub>3</sub>	0	1	2	3	0	0	0	0																																																																																																								
e <sub>4</sub>	0	0	0	0	0	0	0	0																																																																																																								
e <sub>1</sub>	0	0	0	0	0	0	0	0																																																																																																								
e <sub>3</sub>	0	0	0	0	0	0	0	0																																																																																																								
e <sub>4</sub>	0	0	0	0	0	0	0	0																																																																																																								
e <sub>1</sub>	0	0	0	0	0	0	0	0																																																																																																								
e <sub>2</sub>	0	0	0	1	0	0	0	0																																																																																																								
e <sub>4</sub>	0	0	0	0	0	0	0	0																																																																																																								
e <sub>1</sub>	0	0	0	1	0	0	0	1																																																																																																								
e <sub>2</sub>	0	0	0	0	0	0	1	2																																																																																																								
e <sub>3</sub>	0	0	0	0	1	2	3	4																																																																																																								

**Fig. 3.** Follow-arrays with highlighted entries to mark patterns with a non-varying subset of followers.

### LP-max-length

To compute the length of a longest pattern where  $e_i$  is the leader, we utilise a variable  $k^{max}$ . Initially  $k^{max} := 0$ , and at the end it will be equal to the length of a longest leadership pattern (for a specific  $m$ ). Now, for each  $t_x = t_1, \dots, t_\tau$  we check whether  $IntervalsNotFwg(t) > k^{max}$  and if so, we do the following. We inspect the column of the array  $IntervalsFwg$  corresponding to  $t$ . We traverse that column (i.e. we loop for  $j = 1, \dots, n, j \neq i$ ), and we count the number of entities  $e_j$  for which holds that  $IntervalsFwg(t_x, e_j) > k^{max}$ . Let this number be denoted by  $m(k^{max})$ . If  $m(k^{max}) \geq m$ , then we have at least  $m$  entities following  $e_i$  for more than  $k^{max}$  unit-time-intervals, and  $e_i$  is not following anyone during that time. Hence, we increase  $k^{max}$  by one and proceed with the next time-step  $t_{x+1}$ . As we only traverse the entire arrays once, it takes  $O(n\tau)$  time to compute the longest pattern where  $e_i$  is the leader.

The following concluding lemma might surprise, as the longest duration flock pattern is NP-hard to compute and cannot even be approximated within a factor of  $\tau^{1-\epsilon}$  [21].

**Lemma 5.** *The longest duration leadership pattern for a non-varying-subset of followers of size at least  $m$  can be computed in  $O(n^2\tau)$  time and  $O(n\tau)$  space.*

*Example 4.* Consider again Figure 3. For  $m = 1$ , the above described method would find entity  $e_4$  to be the leader (of one entity, namely  $e_3$ ) for four consecutive unit-time-intervals, which is the length of a longest pattern (for  $m = 1$ ).

### LP-max-size

It is also possible to compute the size of a largest non-varying-subset of followers that follows a leader for at least  $k$  unit-time-intervals. We utilise the arrays  $IntervalsNotFwg$  and  $IntervalsFwg$  and a variable  $m^{max}$ , initially set to 0. We update this variable whenever we find a larger set of followers. That is, for  $t_x := t_1, \dots, t_\tau$ , we test if both  $IntervalsNotFwg(t_x) \geq k$  and  $m(k) > m^{max}$ , and if so, we set  $m^{max} := m(k)$ , where  $m(k)$  is defined in the same way as  $m(k^{max})$  in the section above. Hence, we obtain the following lemma.

**Lemma 6.** *The size of a largest non-varying-subset of entities that follow a leader for at least  $k$  time-steps can be computed in  $O(n^2\tau)$  time and  $O(n\tau)$  space.*

*Example 5.* Consider again Figure 3, and let  $k = 1$ . The algorithm above computes  $m^{max} = 3$  as entity  $e_4$  is a leader of 3 entities for  $k = 1$  unit-time-interval at time-step 8.

### 3.3 Solving LP Problems with a Varying Subset of Followers

The variants of the problem of finding leadership patterns where the set of followers can change during the leadership pattern can be solved in a similar way as proposed in Section 3.2. To determine if an entity  $e_i$  is a leader of a varying-subset of followers, we use the follow-arrays  $IntervalsNotFwg(t_x)$ ,  $IntervalsFwd_m(t_x)$  and  $NumFws(t_x)$  as described in Section 3.1.

### LP-report-all

In the same flavour as described above, we can find out if  $e_i$  is a leader. We look for and report time-steps  $t_x$ , such that  $IntervalsNotFwg(t_x) \geq k$  and  $IntervalsFwd_m(t_x) \geq k$ . It is easy to see that reporting when  $e_i$  is a leader can be done in  $O(n\tau)$  time.

*Example 6.* Let  $m = 1$  and  $k = 2$ . Consider  $e_1$ 's follow-arrays in the upper half of Figure 4. Above method reports one time-steps (namely time-step 3) where  $e_1$  is a leader of at least  $m = 1$  entities for at least  $k = 2$  unit-time-intervals.

The complexity of finding all leadership patterns for  $n$  entities is summarised as follows.

**Lemma 7.** *Reporting all varying-subset leadership patterns amongst  $n$  trajectories over  $\tau$  time-steps can be done in  $O(n^2\tau)$  time and  $O(n\tau)$  space.*

	1 2 3 4 5 6 7 8	1 2 3 4 5 6 7 8	1 2 3 4 5 6 7 8	1 2 3 4 5 6 7 8
$IntervalsNotFwg(t)$ :	$e_1$ 0 1 2 0 1 2 3 0	$e_2$ 0 0 0 0 0 0 1 0 0	$e_3$ 0 0 0 0 0 0 0 0 0	$e_4$ 0 1 2 3 4 5 6 7
$IntervalsFwd_m(t)$ ( $m=1$ )	$e_1$ 0 1 2 3 4 0 0 0	$e_2$ 0 0 0 0 0 0 0 0	$e_3$ 0 0 0 1 0 0 0 0	$e_4$ 0 0 0 1 2 3 4 5
$IntervalsNotFwg(t)$ :	$e_1$ 0 1 2 0 1 2 3 0	$e_2$ 0 0 0 0 0 0 1 0 0	$e_3$ 0 0 0 0 0 0 0 0 0	$e_4$ 0 1 2 3 4 5 6 7
$numFws(t)$ :	$e_1$ 0 2 2 2 1 0 0 0	$e_2$ 0 0 0 0 0 0 0 0	$e_3$ 0 0 0 1 0 0 0 0	$e_4$ 0 0 0 1 1 1 2 3

**Fig. 4.** Follow-arrays with highlighted entries to mark patterns with a varying subset of followers.

### LP-max-length

For computing the longest duration flock, we use the arrays  $IntervalsNotFwg$  and  $IntervalsFwd_m$ , and we search for the largest  $k^{max}$  (initially  $k^{max} := 0$ ) such that there is a time-step  $t_x$  for which  $IntervalsNotFwg(t_x) \geq k^{max}$  and  $IntervalsFwd_m(t_x) \geq k^{max}$ . This can be done as follows. For  $t_x = t_1, \dots, t_\tau$ , we check if  $\min\{IntervalsNotFwg(t_x), IntervalsFwd_m(t_x)\} > k^{max}$ , and if so, we perform an update  $k^{max} := \min\{IntervalsNotFwg(t_x), IntervalsFwd_m(t_x)\}$  and proceed with the next time-step  $t_{x+1}$ .

**Lemma 8.** *The longest duration leadership pattern for a varying-subset of followers of size at least  $m$  can be computed in  $O(n^2\tau)$  time and  $O(n\tau)$  space.*

*Example 7.* Looking at the follow-arrays in the upper half of Figure 4, we see that  $e_4$  is a leader of at least  $m = 1$  entity for  $k^{max} = 5$  unit-time-intervals (starting at time-step 3 and ending at time-step 8).

### LP-max-size

If we would like to compute the size of a largest varying set of followers that follow  $e_i$  for at least  $k$  time-steps, we cannot use the array  $IntervalsFwd_m$  directly as this array contains information only for one specific  $m$ . However, an easy way is to use binary search on  $m$  and recompute the  $IntervalsFwd_m$  array for each value of  $m$ . This adds an additional  $\log n$  factor to the running time.

We propose a slightly different approach. By spending linear preprocessing time, we can compute the minima of any substring of a sequence of numbers in  $O(1)$  time. For more information on this Range Minimum Query (RMQ), see e.g. [3]. Now, we use the array  $NumFws$  and we look for at least  $k$  consecutive unit-time-intervals such that the minimum number of followers in the array  $NumFws$  during that time is as large as possible and  $e_i$  can be a leader. That is, we are

looking for  $k$  consecutive unit-time-intervals such that  $e_i$  does not follow any other entity and for the largest minimum (to be referred to as  $m^{max}$ ) over all numbers of followers corresponding to those  $k$  consecutive unit-time-intervals. All minima can be computed in  $O(\tau)$  time [3], hence,  $m^{max}$  can be computed in linear time.

**Lemma 9.** *The size of a largest varying-subset of entities that follow a leader for at least  $k$  time-steps can be computed in  $O(n^2\tau)$  time and  $O(n\tau)$  space.*

*Example 8.* Consider the lower half of Figure 4 and let  $k = 2$ . Above algorithm computes  $m^{max} = 2$  at time-step 3 for entity  $e_1$  and at time-steps 8 for entity  $e_4$ .

## 4 Exact Algorithms for the Continuous Case

In contrast to the discrete version of the leadership pattern, where a pattern can only start or end at the given discrete time-steps, in the continuous version of the problem a pattern can start and end at any point in time. As we do not have spatial information of the entities between two consecutive time-steps we use Assumption 1 to tackle the continuous version in this section. The main ideas are similar to the discrete case, but instead of using arrays storing single numbers to represent follow-behaviour we will use sets of time intervals. First, we describe how to compute them for a fixed entity  $e_i$  and then we define two operations on (sets of) intervals. Later, these intervals and operations are used to solve our leadership problems.

### 4.1 Getting Ready – Follow-Intervals for an Entity $e_i$

**Computing Follow-Intervals:** A first step is to compute a set  $SetNotFwg$  of notfollowing-intervals representing when a fixed entity  $e_i$  is not following any other entity  $e_j$ . An interval  $I = (t_{x_a}, t_{y_a}) \in SetNotFwg$  with  $t_{x_a} \leq t_{y_a}$  means that entity  $e_i$  is not following any other entity during the whole time interval  $I$ . Because entities move on a straight line between two consecutive time-steps, cf. Assumption 1,  $e_i$  can be involved in at most two events that change its follow-behaviour (i.e. the events of beginning or ending to follow) for each entity between two consecutive time-steps. That is why the set  $SetNotFwg$  contains  $O(n\tau)$  intervals. We can compute this set with two nested loops one over all time-steps, another over all entities. By Lemma 2, this can be done in  $O(n\tau)$  time in total.

We also need information about which entities follow  $e_i$ . This information is again stored in a set  $SetFwd$  of intervals. (There is no need to store which entity belongs to which interval, as we do not want to output the names of the followers.) An interval  $I = (t_{x_a}, t_{y_a}) \in SetFwd$  with  $t_{x_a} < t_{y_a}$  means that  $e_i$  is followed by an entity, say  $e_j$ , during the whole time interval  $I$ . Also this set contains at most  $O(n\tau)$  intervals, as an entity can change its follow behaviour with respect to  $e_i$  at most twice between two consecutive time-steps. We can compute this set with two nested loops one over all time-steps, another over all entities. By Lemma 2, this can be done in  $O(n\tau)$  time in total.

Both sets of intervals can be computed in  $O(n\tau)$  time. For the subsequent methods, however, we need the start- and end points of the intervals in non-decreasing order and that the intervals are maximal. Obtaining the sets such that the start- and end points are sorted can be done in  $O(n\tau \log n)$  time in the following way. For each set we use two nested loops. The outer loop fixes an entity and the inner loop ranges over the time-steps. In that way it is easy to compute the intervals as maximal intervals. Whenever we compute start- or end points of an interval we can put them into  $\tau - 1$  buckets, namely one for each unit-time-interval, i.e. pair of consecutive time-steps. As we can have at most  $O(n)$  start- or end points in each bucket, we can sort all of them in all buckets in  $O(n\tau \log n)$  time. Combining the sorted sequences of each bucket results in a sorted sequence of all start- and end points.

**Lemma 10.** *In  $O(n\tau \log n)$  time and  $O(n\tau)$  space, the sets  $SetNotFwg$  and  $SetFwd$  for an entity  $e_i$  can be computed such that all the start- and end points of the intervals in each set are output in non-decreasing order.*

Next, we define operations that take and return a set of intervals as input and output. We also briefly describe how to compute these operations, if the set of intervals is given along with the start- and end points in sorted order.

**Combining Intervals:** We call the first operation under consideration *interval-combination* denoted as  $ic_x(S)$ , where  $S$  is a set of intervals of  $\mathbb{R}$ . The operation outputs a set of non-intersecting intervals. Every point in  $\mathbb{R}$  that is contained in at least  $x$  intervals of the input-set  $S$  will be in an interval of the output-set. Also, for every point in an interval of the output-set, there are at least  $x$  intervals in the input-set that all contain that point, see Figure 5. Note that  $ic_1(S)$  is the union of all intervals in  $S$  and  $ic_{|S|}(S)$  is the intersection of all intervals in  $S$ .

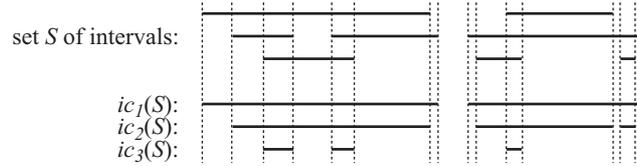
**Definition 4.** Given a set  $S$  of intervals and an integer  $x$ , the operation  $ic_x$  is defined as

$$ic_x(S) = \{I_1, \dots, I_y\}$$

where  $t \in I_z$ , for  $I_z \in \{I_1, \dots, I_y\}$  iff there exist at least  $x$  intervals  $J_1, J_2, \dots, J_x \in S$ , such that  $t \in J_1 \cap J_2 \cap \dots \cap J_x$ .

Let  $S$  be a set of intervals where the start- and end points are given in sorted order. The operation  $ic_x(S)$  can be computed by scanning or sweeping over the sorted start- and end points and keeping track of how many intervals are currently ‘active’.

**Lemma 11.** Suppose  $S$  is a set of intervals. If the start- and end points of the intervals in  $S$  are given in non-decreasing order, then we can compute  $ic_x(S)$  in  $O(|S|)$  time.



**Fig. 5.** The set  $S$  of intervals on the real line and the results after applying the  $ic_x$  operation for  $x \in \{1, 2, 3\}$ . Note that  $ic_x(S) = \emptyset$  for all  $x \geq 4$  as the intersection of any 4 intervals in  $S$  is empty.

**Clipping Intervals:** We also define another operation, which modifies single intervals. For an interval  $I = \{t_{x_a}, t_{x_b}\}$ , we cut or clip a part of length  $k$  at the beginning of  $I$ . If the resulting interval  $I'$  is non-empty, then that interval  $I'$  is the result of the operation.

$$clip_k(\{t_{x_a}, t_{x_b}\}) := \begin{cases} \{t_{x_a} + k, t_{x_b}\} & \Leftrightarrow t_{x_a} + k \leq t_{x_b} \\ \emptyset & \text{otherwise} \end{cases}$$

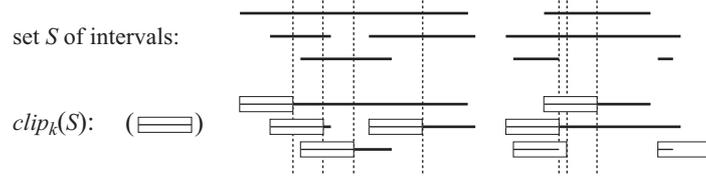
This operation can also be applied to all intervals of an entire set (cf. Figure 6), such that the order of the start- and end points of the intervals remains stable.

**Lemma 12.** Let be given a set  $S$  of intervals, where the start- and end points of the intervals in  $S$  are given in non-decreasing order. We can compute  $S' := \{I' \mid I' = clip_k(I), I \in S\}$  and output the start- and end points of all intervals in  $S'$  in non-decreasing order in  $O(|S|)$  time.

## 4.2 Solving LP Problems with a Non-Varying Subset of Followers

### LP-report-all

We first look at the non-varying-subset version. In the previous section we have seen that we can compute the interval-set  $SetFwd$  in  $O(n\tau \log n)$  time, where an interval in this set means that an



**Fig. 6.** The set  $S$  of intervals on the real line and the results after applying the  $clip_k$  operation. For the  $clip_k$  operation, the length of the interval in parentheses determines  $k$ .

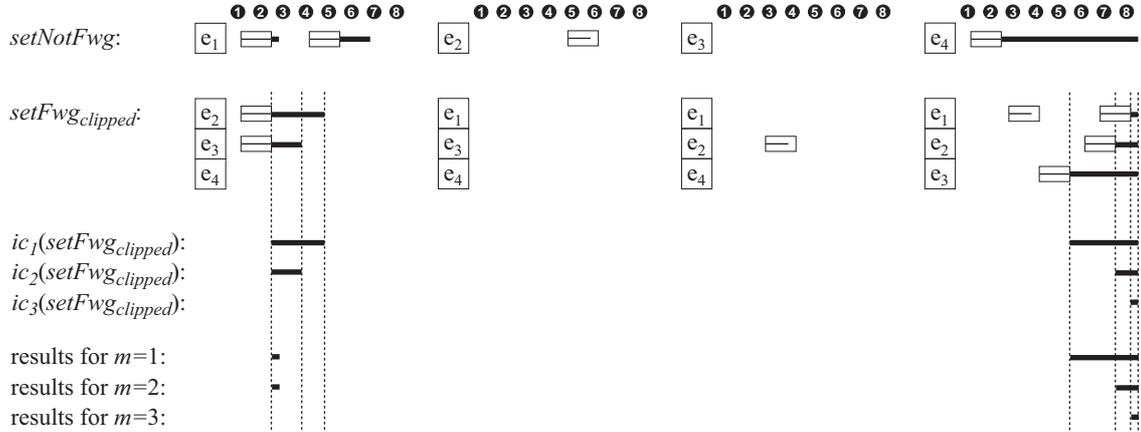
entity follows  $e_i$  for the time of the interval. Now, we are going to modify the intervals in the set  $SetFwd$ . For each interval  $I = \{t_{x_a}, t_{x_b}\} \in SetFwd$ , we apply the operation  $clip_k$  to obtain a set

$$SetFwd_{clipped} := \{I' \mid I' = clip_k(I), I \in SetFwd\}.$$

Note that  $SetFwd_{clipped}$  (see Figure 7) only contains intervals whose originals had length at least  $k$ . The meaning of an interval  $I' \in SetFwd_{clipped}$  is that there is an entity such that at each time-point  $t \in I'$  this entity has already followed  $e_i$  for at least  $k$  time units (which is not necessarily the same as  $k$  unit-time intervals). The set  $SetFwd_{clipped}$  can be computed in linear time with respect to the size of  $SetFwd$ , and this can be implemented such that the order of the (start- and end-points of the) intervals remains stable. We also clip the intervals of the set  $SetNotFwd$  to obtain a set

$$SetNotFwd_{clipped} := \{I' \mid I' = clip_k(I), I \in SetNotFwd\}$$

(see Figure 7). For each time-point in an interval in  $SetNotFwd_{clipped}$ , we have that  $e_i$  is not following any other entity for at least  $k$  time units.



**Fig. 7.** Illustration of clipping and combining the intervals, where the intervals represent the follow-behaviour of the entities in Figure 1. The *result*-intervals are shown for different values of  $m$ .

The next step is to compute yet another set  $S$  of intervals as an interim result using one of the operations introduced in Section 4.1,

$$S := ic_m(SetFwd_{clipped}).$$

For any time-point in an interval in  $S$  there are at least  $m$  entities following  $e_i$ , where each of those entities already followed  $e_i$  for at least  $k$  time units. Finally, we combine the information

represented by  $S$  and  $SetNotFwg_{clipped}$ . What we need is similar to a logical ‘and’ between intervals of those two sets, and this can be done by applying the  $ic_x$  again to obtain a set of *result*-intervals,

$$result := ic_2(S \cup SetNotFwg_{clipped}).$$

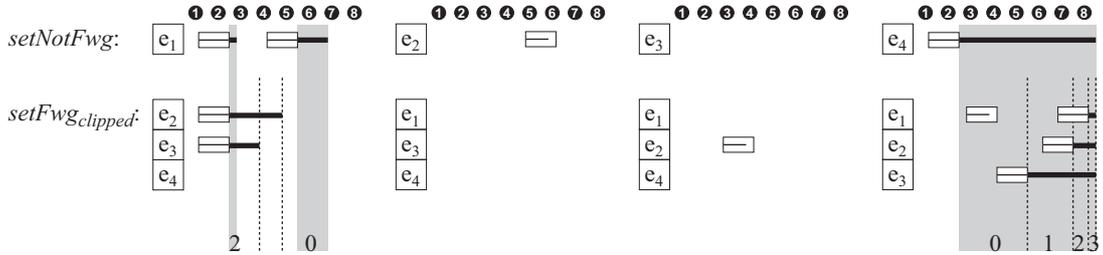
Note that the start- and end points of the set  $S \cup SetNotFwg_{clipped}$  can be sorted in linear time if the start- and end points of  $S$  and  $SetNotFwg_{clipped}$  are sorted. The set *result* contains all intervals for which  $e_i$  is a leader of at least  $m$  entities for at least  $k$  time units. If we would like to report the leadership patterns of all entities, we apply the above method to each entity. Hence, we can conclude with the following lemma.

**Lemma 13.** *Let be given  $n$  trajectories over  $\tau$  time-steps. Reporting all time intervals where there is an entity a leader of a non-varying subset of at least  $m$  entities for at least  $k$  time units, can be done in  $O(n^2\tau \log n)$  time and  $O(n\tau)$  space.*

### LP-max-size

We can use the sets  $SetNotFwg_{clipped}$  and  $SetFwd_{clipped}$ , where the intervals are given in non-decreasing order, to find the maximum  $m^{max}$  for which  $e_i$  is a leader of a non-varying set of  $m^{max}$  entities for at least  $k$  time units. To that end, we do not collapse the set  $SetFwd_{clipped}$  into a set  $S$  as described above, but we utilise a sweep-line approach over the intervals in  $SetNotFwg_{clipped}$  and  $SetFwd_{clipped}$ .

Sweeping means moving an imaginary vertical line over the horizontally arranged intervals, stopping at certain event points, and performing certain actions (see e.g. [11] for more details). In our case the event points are the start- and end points of the intervals. For any position of the sweep-line we say an interval  $I$  is *active*, if the sweep-line  $\ell$  intersects interval  $I$ . During the sweep, we keep track of the number of active intervals in  $SetFwd_{clipped}$ , where the intervals in  $SetNotFwg_{clipped}$  are used as a mask (see Figure 8). All this can be done in  $O(n\tau)$  time.



**Fig. 8.** Illustrating the sweep-line approach. The shaded region indicates how the  $SetNotFwg$  intervals are used as a mask. The numbers indicate the number of active  $SetFwd$  intervals.

**Lemma 14.** *Let be given  $n$  trajectories over  $\tau$  time-steps. Computing the maximum size of a non-varying subset of followers which follow a leader for at least  $k$  time units, can be done in  $O(n^2\tau \log n)$  time and  $O(n\tau)$  space.*

### LP-max-length

A method similar to the one presented above cannot be used directly, as the sets of intervals are computed for specified values of  $k$ . We could use binary search on  $k$ , however, this would add another  $\log \tau$  factor to the running time. The method described in this section also builds upon the sets  $SetNotFwg$  and  $SetFwd$  of intervals, introduced in Section 4.1. It finds the largest  $k^{max}$ , such

that there is a non-varying subset of at least  $m$  entities following entity  $e_i$  for  $k^{max}$  time-units. However, it can also be used to report patterns where  $e_i$  is a leader of at least  $m$  non-varying followers for at least  $k$  time-steps. We do this by sweeping over the sets of intervals, assuming they are given such that the start- and end points of the intervals are in non-decreasing order.

During the sweep we keep track of the active intervals in *SetNotFwg*. Note that only one interval  $I \in \text{SetNotFwg}$  can be active at a time. By keeping a pointer  $p_1$  to  $I$ , we know for every time-point  $t$ , whether  $e_i$  is following any other entity. If  $e_i$  is following any other entity, then there is no interval in *SetNotFwg* active at time  $t$  (and  $p_1$  becomes a null-pointer). On the other hand, if there is an interval  $I \in \text{SetNotFwg}$  active at  $t$ , then we can compute for how long  $e_i$  is not following any other entity.

We also keep track of how many entities follow  $e_i$  and for how long. To this end, let  $t$  be a time-point during the sweep. Let  $A \subseteq \text{SetFwd}$  be the set of all intervals in *SetFwd* that are active at time  $t$ . We will not maintain  $A$ , but only a variable  $m'$  with  $m' = |A|$  (initially  $m' := 0$ ). Furthermore, we will maintain a pointer  $p_2$  to the interval in  $A$  with the  $m$ -th largest end point. If  $A$  contains less than  $m$  intervals, then  $p_2$  points to some interval. (As we will not use pointer  $p_2$  if  $A$  contains less than  $m$  intervals it is not important where  $p_2$  points to in that case.)

Before the sweep, we initialise  $k^{max} := 0$ , and after the sweep  $k^{max}$  will be the length of the longest leadership pattern with  $e_i$  as leader of a non-varying set of at least  $m$  entities. We also introduce an artificial interval which starts and ends before any other interval starts and we initialise  $p_2$  to point to that interval. This interval is introduced merely to have pointer  $p_2$  well initialised. The sweep does not take this interval into consideration. As mentioned above the event points are the start- and end points of the intervals, and if two event points have the same time, we process them one after the other, as if one was infinitesimally later than the other. By maintaining all invariants it is easy to see that for every position of the sweep-line with corresponding time  $t$ , we can check if there are at least  $m$  entities following  $e_i$ , i.e. if  $m' \geq m$ . In the case that there are at least  $m$  followers of  $e_i$ , we also can determine for how long in the future all these entities will follow  $e_i$ , by using the pointer  $p_2$ . Furthermore, we can check if  $e_i$  is following any other entity (by using  $p_1$ ), and if not for how long in the future  $e_i$  will not follow any other entity. Therefore, we can determine whether there is a leadership pattern, with  $e_i$  as leader of a non-varying set of at least  $m$  entities, and if there is such a pattern, we can also determine its length  $k'$ . If  $k' > k^{max}$  then we perform an update  $k^{max} := k'$ . More precisely, at each event point  $t$  we do the following, depending on the type of event point (there are four types).

- Case 1 ‘ $t$  is a start point of an interval  $I = \{t_{x_a}, t_{x_b}\} \in \text{SetNotFwg}$ ’:  
Make  $p_1$  point to  $I$ . If  $m' \geq m$  then let  $I' = \{t_{x_c}, t_{x_d}\}$  be the interval  $p_2$  is pointing to. Now we know that there are at least  $m$  followers until time  $t_{x_d}$  and that  $e_i$  is not following anyone until time  $t_{x_b}$ . Let  $k' := \min\{t_{x_b}, t_{x_d}\} - t$ . If  $k' > k^{max}$  then  $k^{max} := k'$ .
- Case 2 ‘ $t$  is an end point of an interval in *SetNotFwg*’:  
Make  $p_1$  a null-pointer.
- Case 3 ‘ $t$  is a start point of an interval  $I = \{t_{x_a}, t_{x_b}\} \in \text{SetFwd}$ ’:  
Increment  $m'$  as  $I$  is a new active interval (among those in *IntervalsFwg*).  
Subcase 3.1 ‘ $m' = m$ ’:  
Before encountering the current interval, there were less than  $m$  intervals of *SetFwd* active, and hence  $p_2$  points to a time-point that already has been passed by the sweep-line. To make  $p_2$  point to the interval in  $A$  with the  $m$ -th largest end point, we make a sub-sweep with another (imaginary) sweep-line starting at the time-point where  $p_2$  points to and go further until we find the first end point of an active interval  $I'$ . Make  $p_2$  point to that interval  $I' = \{t_{x_c}, t_{x_d}\}$ . We still have to check for a pattern and for its length, and hence, we do the following. We know that there are at least  $m$  followers until time  $t_{x_d}$ . If  $p_1$  is a nullpointer, stop this subcase, otherwise do the following. Let  $t_{x_e}$  be the end point of the interval  $p_1$  is pointing to. Let  $k' := \min\{t_{x_d}, t_{x_e}\} - t$ , as this is the time length  $e_i$  is not following anyone and  $e_i$  is followed by at least  $m$  entities. If  $k' > k^{max}$  then  $k^{max} := k'$ .  
Subcase 3.2 ‘ $m' > m$ ’:  
Let  $I' = \{t_{x_c}, t_{x_d}\}$  be the interval  $p_2$  is pointing to. Note that  $I'$  was the interval in  $A$  with the

$m$ -th largest end point just before  $I$  was encountered by the sweep-line. If  $t_{x_b} > t_{x_d}$  then we make again a sub-sweep starting with  $p_2$  at  $t_{x_d}$  until finding the next end point of an active interval. Make  $p_2$  point to that interval. This is done because after encountering  $I$ , interval  $I'$  becomes the interval with  $(m + 1)$ -th largest end point among the intervals in  $A$ . Therefore, we have to advance the pointer  $p_2$  by one active interval. If however,  $t_{x_b} < t_{x_d}$  then we do not have to modify  $p_2$ . Now, let  $I' = \{t_{x_c}, t_{x_d}\}$  be the (possibly new) interval  $p_2$  is pointing to. Also here we have to check for a pattern and for its length. This is done in a similar way as in the previous subcase.

Subcase 3.3 ‘ $m' < m$ ’:

As before and after encountering  $I$  there are less than  $m$  active intervals, we do not need to do anything in this case.

- Case 4 ‘ $t$  is an end point of an interval  $I \in \text{SetFwd}$ ’:  
Decrement  $m'$  as  $I$  was an active interval (among those in  $\text{SetFwd}$ ).

Apart from the subsweeps it is easy to see that at any event point we spend constant time. Whenever we perform a subsweep in Subcases 3.1 and 3.2 the pointer  $p_2$  only moves forward (perhaps) multiple steps. Hence, we can see that all subsweeps together take  $O(n\tau)$  time.

**Lemma 15.** *Let be given the sets  $\text{SetNotFwg}$  and  $\text{SetFwd}$  for an entity  $e_i$ , where the start- and end points of those intervals are given in non-decreasing order. Determining the largest  $k^{\max}$  such that  $e_i$  is a leader of a non-varying subset of at least  $m$  entities for  $k^{\max}$  time units, can be done in  $O(n\tau)$  time.*

By doing this sweep-line approach for each entity, we can compute the overall longest duration leadership pattern.

**Lemma 16.** *Let be given  $n$  trajectories over  $\tau$  time-steps. Computing the maximum length of a leadership pattern with a non-varying subset of followers of size at least  $m$ , can be done in  $O(n^2\tau \log n)$  time and  $O(n\tau)$  space.*

### 4.3 Solving LP Problems with a Varying Subset of Followers

#### LP-report-all

After considering the case for the non-varying subset in Section 4.2, the case for a varying subset is rather easy. Here, we do not require that all entities follow  $e_i$  for  $k$  time-units. Hence, with the terminology as used before we compute a set

$$S := ic_m(\text{SetFwd}).$$

For any time-point in an interval in  $S$  there are at least  $m$  entities following  $e_i$ . As  $e_i$  still has to be followed for at least  $k$  time-units, we clip all intervals in  $S$  to obtain

$$S' := \{I' \mid I' = \text{clip}_k(I), I \in S\}.$$

As before, our last step is to combine  $S'$  and  $\text{SetNotFwg}_{\text{clipped}}$  to obtain the set of *result*-intervals,

$$\text{result} := ic_2(S' \cup \text{SetNotFwg}_{\text{clipped}}).$$

**Lemma 17.** *Let be given  $n$  trajectories over  $\tau$  time-steps. Reporting all time intervals where there is an entity a leader of a varying subset of at least  $m$  entities for at least  $k$  time units, can be done in  $O(n^2\tau \log n)$  time and  $O(n\tau)$  space.*

**LP-max-size**

In this case, we can use the approach mentioned in Section 3.3, where we spend additional time for binary search on  $m$  to find  $m^{\max}$ . This additional  $\log n$  factor does not increase the overall running time of our method as we have a time bound of  $O(n\tau \log n)$  for computing the sets *IntervalsNotFwg* and *IntervalsFwg*, such that their start- and end points are sorted in non-decreasing order.

**Lemma 18.** *Let be given  $n$  trajectories over  $\tau$  time-steps. Computing the maximum size of a varying subset of followers which follow a leader for at least  $k$  time units, can be done  $O(n^2\tau \log n)$  time and  $O(n\tau)$  space.*

**LP-max-length**

To find the length of a longest duration leadership pattern of a varying set of at least  $m$  entities, we can use a similar approach as in Section 4.3. We also compute a set

$$S := ic_m(\text{SetFwd}),$$

such that for each time-point in an interval  $I \in S$ , we know that there are at least  $m$  followers of  $e_i$ . To combine this with the information when  $e_i$  is not following any other entity, we apply the operation  $ic_x$  once again to obtain

$$\text{result} := ic_2(S \cup \text{SetNotFwg}).$$

Now, for any interval in *result* it holds that  $e_i$  is not following any other entity, and also that  $e_i$  is followed by at least  $m$  entities. Searching for the length of the longest interval in *result* solves the problem at hand for entity  $e_i$ .

**Lemma 19.** *Let be given  $n$  trajectories over  $\tau$  time-steps. Computing the maximum length of a leadership pattern with a varying subset of followers of size at least  $m$ , can be done  $O(n^2\tau \log n)$  time and  $O(n\tau)$  space.*

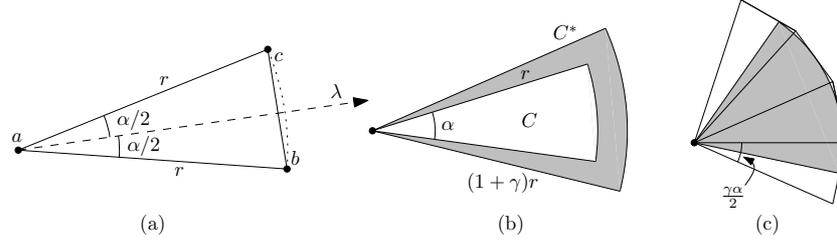
**5 Approximate leadership patterns**

In this section we will consider an approximate version of the leadership pattern, where we approximate the front region. For the discrete and varying cases of our leadership problems we have seen in Section 3.3 that the arrays *IntervalsNotFwg*, *IntervalsFwd<sub>m</sub>* and *NumFws* can be used to identify patterns for a fixed entity  $e_i$ . Recall, that these arrays are one dimensional and can be computed in  $O(n\tau)$  time and space (see Section 3.1 for more details). In this section, we will see how we can more efficiently compute approximate versions of them. The general idea is to approximate the front regions of the entities. Furthermore, we slightly modify the definition of leadership to allow for the approximation.

We say that an entity  $e$  does not follow anyone at time  $[t_x, t_y]$  if for each time-step  $t \in [t_x, t_y]$  the front-regions of  $e$  with respect to the adjacent unit time intervals are empty. An entity  $e$  is said to be followed by  $m$  entities at time  $[t_x, t_y]$  if for each time-step  $t \in [t_x, t_y]$  and each adjacent unit-time-interval  $I$  there are at least  $m$  entities with  $e$  in their front-regions with respect to  $I$ . The definition of leadership and the follows-arrays are now defined using the above definitions. Before we go into the details we will need some basic notations, definitions and data structures. Before we go into the details we will need some basic notations, definitions and data structures.

**5.1 Preliminaries**

Let  $T = (a, b, c)$  be an isosceles triangle with  $|ab| = |ac| = r$  and  $\angle cab = \alpha$ . The orientation  $\lambda$  of a triangle is the direction of its bisector through  $a$ , as shown in Fig. 9a. We say that  $T$  has apex at  $a$ , radius  $r$ , orientation  $\lambda$  and angle  $\alpha$ . Note that we can define a sector of a disk using the same parameters. The following lemma is well-known and can be obtained by using range trees [5].



**Fig. 9.** (a) Illustrating the definition of an isosceles triangle and a sector. (b) Illustrating the  $(1 + \gamma)$ -approximate region of  $C$ . (c) A sector with interior angle  $\alpha$  and radius  $r$  can be  $(1 + \gamma)$ -approximated with a small number of fixed orientation isosceles triangles.

**Lemma 20.** *Given a set  $S$  of  $n$  points in the plane and two positive real values  $\lambda < 2\pi$  and  $\alpha \leq \pi/2$  one can preprocess  $S$  in  $O(n \log^2 n)$  time such that range counting queries can be answered in  $O(\log^2 n)$  time, where the range is an isosceles triangle with orientation  $\lambda$  and angle  $\alpha$ .*

Consider two sectors  $C$  and  $C^*$  apex at  $a$  and orientation  $\lambda$ , where  $C$  has radius  $r$  and interior angle  $\alpha$  while  $C^*$  has radius  $(1 + \gamma)r$  and interior angle  $(1 + \gamma)\alpha$ , see Fig. 9b. A region  $R$  is said to be a  $(1 + \gamma)$ -approximation of  $C$  if and only if  $C^*$  contains  $R$  and,  $R$  contains  $C$ . If  $C$  is a front region of an entity  $e$  then  $R$  is said to be a  $(1 + \gamma)$ -approximate front region. The set of orientations that is a multiple of  $\frac{\gamma\alpha}{2}$  is said to be the set of fixed orientations. The following observation is straight-forward (see also Fig. 9c).

**Observation 1** *Given two positive real values  $\alpha \leq \pi/2$  and  $\gamma$ , every sector with interior angle  $\alpha$  can be  $(1 + \gamma)$ -approximated by the union of  $\frac{\gamma+2}{\gamma}$  isosceles triangles of fixed orientation and interior angle  $\frac{\gamma\alpha}{2}$ .*

Combining Lemma 20 and Observation 1 we obtain the following corollary that will be used repeatedly in the rest of this section.

**Corollary 1.** *Given a set  $S$  of  $n$  points in the plane and two positive real values  $\alpha \leq \pi/2$  and  $\gamma$  one can preprocess  $S$  in  $O(\frac{n}{\gamma\alpha} \log^2 n)$  time such that  $(1 + \gamma)$ -approximate range counting queries can be answered in  $O(\frac{1}{\gamma\alpha} \log^2 n)$  time, where the range is any sector with interior angle  $\alpha$ .*

In Section 3.3, it was shown how to extract discrete leadership patterns with varying followers from the *IntervalsNotFwg*, *IntervalsFwd<sub>m</sub>* and *NumFws*-arrays. Thus it suffices to show how the approximate versions of these arrays can be computed.

## 5.2 Approximating *IntervalsNotFwg*

The algorithms described in Sections 3.1 to compute the array *IntervalsNotFwg* for an entity  $e_i$  takes  $O(\tau n)$ . In this section it will be shown how this can be done much faster if we allow a  $(1 + \gamma)$ -approximation of the front regions instead of the exact front regions when deciding if an entity is following any other entity. The result will be stored in an array that we call a  $(1 + \varepsilon)$ -approximate *IntervalsNotFwg*-array.

Process each entity  $e$  as follows. Consider every two consecutive time-steps  $t_{x-1}$  and  $t_x$  separately. Let  $R_{x-1}(e)$  and  $R_x(e)$  be the two front regions of  $e$  with respect to  $(t_{x-1}, t_x)$ . Perform two  $(1 + \gamma)$ -approximate range queries of  $R_{x-1}(e)$  and  $R_x(e)$ . If the number of reported points in both regions is zero then set *IntervalsNotFwg*( $t_x$ ) to be *IntervalsNotFwg*( $t_x - 1$ ) + 1 otherwise set it to be 0. If we set  $\gamma = \varepsilon$ , we get:

**Theorem 1.** *A  $(1 + \varepsilon)$ -approximate *IntervalsNotFwg*-array for an entity  $e_i$  can be computed in time  $O(\frac{\tau}{\varepsilon\alpha} \log^2 n)$  time, using  $O(\frac{\tau n}{\varepsilon\alpha} \log^2 n)$  preprocessing time and space.*

### 5.3 Approximating $IntervalsFwd_m$ and $NumFws$

Partition the set of entities into  $\sigma$  subsets  $S_1, \dots, S_\sigma$ , where an entity  $e$  belongs to  $S_j$  if  $d(e)$  is in the interval  $[(j-1) \cdot \theta, j \cdot \theta)$ , where  $\sigma = \frac{2\pi}{\theta}$  and  $d(e)$  is the direction of an entity  $e$  with respect to  $(t_{x-1}, t_x)$ . (The value of  $\theta$  will be below be set to  $\frac{\varepsilon\alpha}{6}$ ). For each set  $S_j$  apply Corollary 1 to obtain two data structures  $\mathcal{D}_{x-1}(S_j)$  and  $\mathcal{D}_x(S_j)$  at time-step  $t_{x-1}$  and  $t_x$  respectively. In total  $\tau\sigma$  data structures will be constructed.

Process every entity  $e$  as follows. As above we consider each pair of consecutive time-steps  $(t_{x-1}, t_x)$  independently. Initialise two counters  $f_{x-1} = f_x = 0$  and consider every set  $S_j$  that contains an entity  $e'$  such that  $\|d(e') - d(e)\| \leq \beta$ . Let  $R_j(e)$  be the sector with apex at  $e$ , radius  $r$ , orientation  $\lambda_j = (2j-1) \cdot \frac{\theta}{2}$  and interior angle  $(\alpha + \theta)$ , see Fig. 10a. Perform a  $(1+\gamma)$ -approximate counting query of  $R_j(e)$  in  $\mathcal{D}_{x-1}(S_j)$  and increase  $f_{x-1}$  with the number of reported points. The same is done with  $f_x$ , i.e., perform a  $(1+\gamma)$ -approximate counting query of  $R_j(e)$  in  $\mathcal{D}_x(S_j)$  and increase  $f_x$  with the number of reported points. Continue until all the sets have been processed.

Finally, if  $\min\{f_{x-1}, f_x\} \geq m$  then set  $IntervalsFwd_m(t_x) := IntervalsFwd_m(t_x - 1) + 1$  otherwise set it to be 0, and set  $NumFws(t_x) := \min\{f_{x-1}, f_x\}$ . When all the entities have been processed we are done.

**Lemma 21.** *The algorithm produces a  $((1+\gamma) \cdot (1+\delta))$ -approximate  $IntervalsFwd_m$ -array, where  $\delta < (1 + \frac{2}{1+\gamma}) \frac{\theta}{\alpha}$ .*

*Proof.* To prove the lemma we need to prove (1) all entities that follow  $e$  at time  $t_x$  with respect to  $(t_{x-1}, t_x)$  are counted, and (2) no entity at time  $t_x$  whose  $((1+\gamma)(1+\delta))$ -approximate front region does not include  $e$  is reported.

To prove (1) consider any entity  $e'$  that follows  $e$ . Assume that  $e' \in S_j$ . Since  $\|d(e') - d(e)\| \leq \beta$  the set  $S_j$  must be processed, and hence the query range  $R_j(e)$  will be considered. Since  $e'$  follows  $e$ , the sector with apex at  $e'$ , radius  $r$ , angle  $\alpha$  and orientation  $d(e')$  must contain  $e$ . As a result  $R_j(e)$  must contain  $e'$  at  $t_x$ , which completes (1).

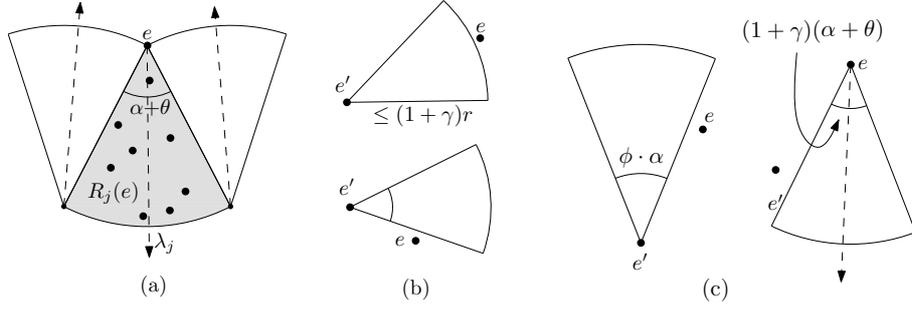
Consider an entity  $e'$  whose  $((1+\gamma)(1+\delta))$ -approximate front region does not contain  $e$ . It will be argued that  $e'$  will not be counted as a follower to  $e$  by the approximation algorithm. Assume that  $e'$  belongs to the set  $S_j$ . The algorithm may perform a  $(1+\gamma)$ -approximate counting query of  $R_j(e)$ . Recall that the radius of  $R_j(e)$  is  $r$  and the interior angle is  $\alpha + \theta$ . We will have two cases as shown in Fig. 10b. If the distance between  $e$  and  $e'$  is greater than  $(1+\gamma)r$  then we are done since  $R_j(e)$  cannot contain  $e'$ .

Otherwise, if the distance is less than  $(1+\gamma)r$  then we have to be more careful. The interior angle of the  $((1+\gamma)(1+\delta))$ -approximate front region of  $e'$  is at most  $(1+\gamma)(1+\delta)\alpha$  and the interior angle of the  $(1+\gamma)$ -approximation of  $R_j(e)$  is at most  $(1+\gamma)(\alpha + \theta)$ , as illustrated in Fig. 10c. Let  $\overline{d(e')} = d(e') + \pi$  and let  $\lambda_j(e)$  denote the orientation of  $R_j(e)$ . From the construction of  $S_j$  it holds that the angle between  $\overline{d(e')}$  and  $\lambda_j$  is bounded by  $\theta$ . If  $e$  lies outside the  $((1+\gamma)(1+\delta))$ -approximate front region of  $e'$  and  $e'$  lies within the  $(1+\gamma)$ -approximate region of  $R_j(e)$  then the following inequality must hold:

$$(1+\gamma) \left( \frac{\alpha + \theta}{2} \right) + \theta \geq (1+\gamma)(1+\delta) \frac{\alpha}{2}.$$

However, this inequality only holds if  $\delta \geq (1 + \frac{2}{1+\gamma}) \frac{\theta}{\alpha}$ . This completes (2), and hence also the lemma.  $\square$

**Theorem 2.** *Using  $O(\frac{\tau n}{\varepsilon\alpha} \log^2 n)$  preprocessing time and space  $(1+\varepsilon)$ -approximate  $IntervalsFwd_m$ - and  $NumFws$ -arrays for an entity  $e_i$  can be computed in time  $O(\frac{\tau}{\varepsilon\alpha} \log^2 n)$  time, for any  $0 < \varepsilon \leq 1$ .*



**Fig. 10.** (a) The region  $R_j(e)$  marked in grey. (b) Illustrating the two cases in the proof of Lemma 21. (c) If  $e$  is not included in the approximate front region of  $e'$  then  $e'$  cannot be included in the approximate region of  $R_j(e)$ .

*Proof.* The value of  $\gamma$  is decided by the approximation factor since  $(1 + \gamma)(1 + \delta) \leq (1 + \varepsilon)$ . We have:

$$\begin{aligned}
 (1 + \gamma)(1 + \delta) &\leq (1 + \gamma)\left(1 + \frac{2}{1 + \gamma}\frac{\theta}{\alpha}\right) \\
 &= 1 + \gamma + \frac{\theta}{\alpha}(3 + \gamma) \\
 &= 1 + \gamma + \frac{\varepsilon}{6}(3 + \gamma) \\
 &< 1 + \varepsilon,
 \end{aligned}$$

where the final inequality follows from setting  $\gamma < \varepsilon/3$ .

Since  $\gamma$  has been set we can now turn our attention to the running time. For each time-step and each set  $S_j$ ,  $1 \leq j \leq \sigma$ , the points are preprocessed as described in Corollary 1. Let  $n_j$  be the number of points in  $S_j$ , we have that the preprocessing time is:

$$\sum_{j=1}^{\sigma} \left( \frac{n_j}{\gamma\alpha} \log^2 n_j \right) = O\left( \frac{n}{\varepsilon\alpha} \log^2 n \right),$$

since  $n_1 + \dots + n_\sigma = n$ .

Next, consider an entity  $e$ . For each of the  $\tau$  time-steps the algorithm will perform  $\frac{2+\gamma}{\gamma}$  triangular range queries requiring a total of  $O\left(\frac{\tau}{\gamma\alpha} \log^2 n\right) = O\left(\frac{\tau}{\varepsilon\alpha} \log^2 n\right)$  time according to Observation 1 and Corollary 1.  $\square$

In Section 3.3 it was shown how to report all discrete leadership patterns with a varying set of followers by using the *IntervalsNotFwg*, *IntervalsFwd<sub>m</sub>* and *NumFws*-arrays. Also, computing the maximum length and the maximum size of a pattern is possible in this case. Thus, putting together Theorems 1 and 2 we conclude this section (assuming  $\alpha$  is a given constant):

**Corollary 2.** *Suppose we are given  $n$  trajectories over  $\tau$  time-steps. If the front-region is  $(1 + \varepsilon)$ -approximated, the discrete case with a varying set of followers of the problems LP-report-all, LP-max-length and LP-max-size can be solved in  $O\left(\frac{\tau n}{\varepsilon} \log^2 n\right)$  time and  $O\left(\tau n + \frac{n}{\varepsilon} \log^2 n\right)$  space.*

## 6 A Lower Bound for the Continuous Case

In this section we argue that it is likely that every algorithm for the continuous version of the leadership problem that detects leadership patterns between two consecutive time-steps in a set of  $n$  trajectories requires  $\Omega(n^2\tau)$  time in the worst case. We will present a transformation from the problem POINT-ON-3-LINES to a special case of our leadership problem.

*Problem:* POINT-ON-3-LINES

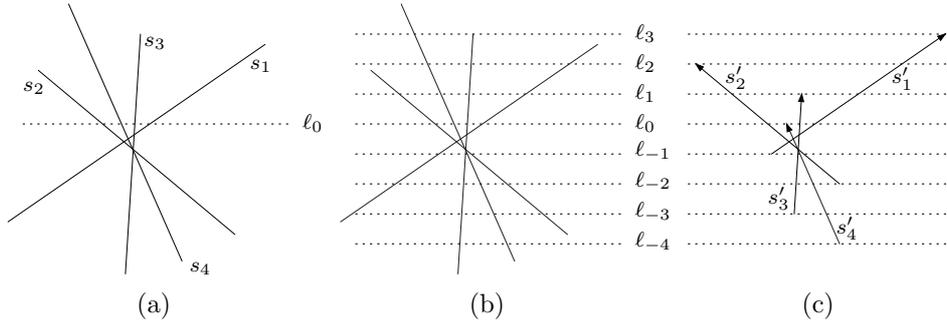
*Instance:* A set  $S$  of  $n$  lines in the plane.

*Question:* Is there a point in the plane that lies on at least three lines of  $S$ ?

The POINT-ON-3-LINES problem was proven to be 3-SUM-hard [20]. This means that it is at least as hard as 3-SUM for which no subquadratic time algorithm has been found yet, which is an indication for an inherent hardness of the problems. For a weak model of computation a lower bound of  $\Omega(n^2)$  for those problems exists [16].

*Remark 1.* Note that the intersection of two or more intervals can be a single number. Therefore it is possible (although not very likely from a practical point of view) that a leadership pattern in the continuous case exists only for one time-point. Such a pattern has length  $k = 0$ .

The transformation is as follows. Consider a set  $S = \{s_1, \dots, s_n\}$  of  $n$  lines in the plane. Let  $\ell_0$  be a horizontal line such that all intersections between two lines of  $S$  lie below  $\ell_0$  (see Figure 11(a)), and let  $\ell_{-1}$  be a line parallel to  $\ell_0$  such that all intersections between two lines of  $S$  lie above  $\ell_{-1}$ . Such two lines can be computed in  $O(n \log n)$  time. Without loss of generality, we assume that no line in  $S$  is parallel to  $\ell_0$ . Let  $\ell_{-n}, \dots, \ell_{-2}, \ell_1, \dots, \ell_{n-2}$  be  $2n - 2$  lines parallel to  $\ell_0$ , such that the distance between  $\ell_i$  and  $\ell_j$  equals  $|i - j|$  times the distance between  $\ell_{-1}$  and  $\ell_0$ , for  $i, j \in \{-n, \dots, n - 1\}$ , see Figure 11(b). For a line  $s_i \in S$ , let  $s'_i$  be the directed line segment that starts at the intersection of  $s_i$  and  $\ell_{-i}$  and ends at the intersection of  $s_i$  and  $\ell_{n-i}$ , as illustrated in Fig. 11(c). Let  $S'$  be the set of all directed line segments, i.e.  $S' = \{s'_i \mid s_i \in S\}$ .



**Fig. 11.** Illustrating a set  $S$  with  $n = 4$  lines  $s_1, \dots, s_4$ . (a)  $S$  and the horizontal line  $\ell_0$  (dotted). (b)  $S$  and the lines  $\ell_{-n}, \dots, \ell_{n-1}$  (dotted). (c) the set  $S'$ .

The set  $S'$  can be viewed as a set of  $n$  trajectories over two time-steps, where every directed segment  $s'_i \in S'$  corresponds to the trajectory of entity  $e_i$  starting at time  $t_{x-1}$  at the origin of  $s'_i$  and reaching the end point of  $s'_i$  at time  $t_x$ . For  $\alpha = 0$ ,  $\beta = \pi$  and  $r = \infty$ , an entity  $e_j$  follows an entity  $e_i$  if and only if  $s'_j$  and  $s'_i$  intersect and the entity  $e_i$  reaches the intersection point before  $e_j$ . Note that because of the construction, no two entities can be between  $\ell_{-1}$  and  $\ell_0$  at the same time and hence when an entity passes through an intersection of line segments in  $S'$ , no other entity passes through an intersection in  $S'$ .

**Lemma 22.** *There are at least three lines in  $S$  passing through the same point, if and only if there is a leadership pattern in  $S'$  with  $m = 2$ ,  $k = 0$ ,  $r = \infty$ ,  $\alpha = 0$  and  $\beta = \pi$ .*

*Proof.* Suppose three lines  $s_i, s_j$  and  $s_k$  intersect in a single point  $p$ . Consider the three corresponding entities  $e_i, e_j$  and  $e_k$  with trajectories  $s'_i, s'_j$  and  $s'_k$ . W.l.o.g. assume that  $e_i$  reaches the intersection point  $p$  first. As  $p$  is an intersection point between lines in  $S$  it is also an intersection

point between line segments in  $S'$ , and therefore when  $e_i$  passes through  $p$  no other entity is between  $\ell_{-1}$  and  $\ell_0$ . Thus, no other entity can intersect  $s'_i$  when  $e_i$  passes through  $p$ , and hence,  $e_i$  is not following another entity when it passes through  $p$ . Furthermore, because  $r = \infty$  and  $\alpha = 0$  the entity  $e_i$  is in the front-regions of  $e_j$  and  $e_k$ . It follows that we have a leadership pattern in  $S'$  with  $m = 2$  and  $k = 0$ .

Now, assume we have a leadership pattern in  $S'$  for  $m = 2$ ,  $k = 0$ ,  $r = \infty$ ,  $\alpha = 0$  and  $\beta = \pi$ . The leader  $e_i$  of that pattern will be followed by at least two entities, say  $e_j$  and  $e_k$ , at some point in time. This implies that  $s'_j$  and  $s'_k$  must intersect  $s'_i$  in the same point, and hence there are three lines in  $S$  passing through the same point.  $\square$

As the transformation described above can be done in  $O(n \log n)$  time, we can conclude with the following lemma.

**Lemma 23.** *Finding continuous leadership patterns between two consecutive time-steps in a set of trajectories is 3-SUM-hard.*

## 7 Experimental Evaluation

This section is devoted to reporting the experimental results. The algorithms were implemented in Java<sup>4</sup> and all experiments were performed on a Linux operated PC with an Intel 3.6 GHz processor and 2 GB of main memory.

### 7.1 Input Data

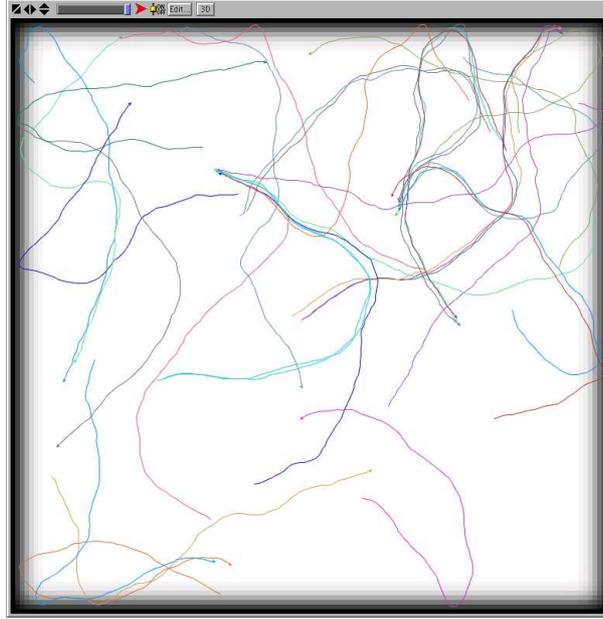
All input files were generated artificially with NetLogo [64]. More specifically, we modified NetLogo’s Flocking Model [63] such that entities do not wrap around the world-borders, but will be repulsed smoothly from walls, see Figure 12. Furthermore, we added some code for moderate random changes in an entity’s direction and saving the coordinates into a file. There are many parameters to modify the behaviour of the entities and thus also to modify how many flocks and leadership-patterns are created. However, we have no direct control over the exact number or length or size of patterns.

We generated files with variable number of entities (128-4096), two different sizes of the underlying universe  $U$  (i.e. coordinate space  $512 \times 512$  and  $1024 \times 1024$ ) and two different characteristics  $CH$  (i.e.  $CH = u$  and  $CH = c$ ) of the entity distribution.  $CH = u$  means that the parameters of the Flocking Model were chosen such that the entities are more uniformly distributed, i.e. only small clusters emerge. Flocks (and thus leadership patterns) still exist but their size and length are likely to be smaller than those of the other characteristic.  $CH = c$  means that the parameters of the Flocking Model were chosen such that the entities form few but rather large clusters, and hence, the flocks tend to contain more entities and have a longer duration. The number of time-steps is  $\tau = 1000$ .

### 7.2 Methods

We performed experiments with two variants of our algorithms for the discrete case. The first one is a straight forward implementation of the method described in Section 3. This method contains (among others) two nested loops ranging over all entities. The disadvantage from a practical point of view is that when looking for entities that might be in a front-region, then also entities that are too far away will be considered. Therefore, our second method tries to overcome this drawback by dividing the underlying plane into buckets (squares of side length  $r$ ). Now when looking for entities that might be in a front-region, only those entities will be considered that are in the nine neighbouring buckets (including the bucket at the centre). Note that for each of our

<sup>4</sup> Java was chosen because this increases the platform independence and it makes it easier to integrate the code into an existing larger framework.



**Fig. 12.** This screenshot of NetLogo’s modified Flocking model shows the trajectories of 32 entities in a universe with side lengths  $128 \times 128$ , run for 100 time-steps.

leadership problem, all methods always compute all arrays from scratch. Especially the arrays *IntervalsNotFwg* and *IntervalsFwg* could be used three times after computing them once. For an easier comparison however, we refrained from doing so.

### 7.3 Results

Tables 2 and 3 show the results of our algorithms for  $m = 10$ ,  $k = 20$ ,  $r = 20$ ,  $\alpha = \pi$  and  $\beta = \frac{\pi}{2}$ . From our point of view the running times and their asymptotic behaviour are much more interesting than for example the exact number of patterns found as we deal with artificial data. Nevertheless, in Table 2 we can see how many entities have been leaders (leaders), the number of leadership patterns found (report-all), the length of a longest duration pattern (max-length) and the number of entities in a pattern with most followers (max-size). Note that patterns with length  $> k$  will be reported multiple times as patterns of length  $k$ .

We observed that the vast majority of the running time is spent on computing the arrays *IntervalsNotFwg* and *IntervalsFwg* (which can be done in  $O(n^2\tau)$  time). Once these two arrays are computed, computing more arrays and/or extracting information to solve the leadership problem is very efficient (linear time). Therefore, our methods for the three different leadership problems result almost always in the same running times (they differ on average less than three percent), as they compute all arrays from scratch. Hence, Table 3 depicts the running times of our methods only for the report-all leadership problem.

### 7.4 Observations

**Non-Varying vs. Varying:** As we could expect running times for the patterns with a varying subset of followers are often higher, as one more array is computed for the ‘varying’ problems. However, this increase is very marginal compared to other influencing factors. We can also observe that the values for the ‘varying’ patterns are at least as big (sometimes slightly larger) as for the ‘non-varying’ patterns. This is because a non-varying pattern is also a varying pattern by definition.

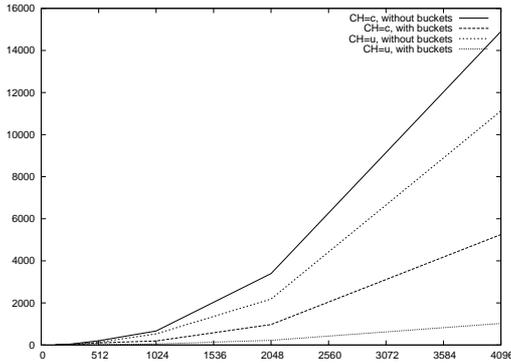
$n$	$U$	$CH$	non-varying				varying			
			leaders	report-all	max-length	max-size	leaders	report-all	max-length	max-size
128	$512^2$	$c$	2	89	37	23	2	161	41	23
256	$512^2$	$c$	10	211	56	66	12	389	71	71
512	$512^2$	$c$	11	329	44	200	13	520	68	238
1024	$512^2$	$c$	27	676	46	197	31	1142	66	208
2048	$512^2$	$c$	33	689	57	384	36	1022	72	419
4096	$512^2$	$c$	44	966	55	812	50	1541	78	959
128	$1024^2$	$c$	0	0	0	4	0	0	0	4
256	$1024^2$	$c$	0	0	2	8	0	0	2	8
512	$1024^2$	$c$	19	360	46	26	27	643	60	26
1024	$1024^2$	$c$	36	954	53	166	47	1591	73	183
2048	$1024^2$	$c$	80	1536	47	219	97	2833	101	224
4096	$1024^2$	$c$	98	2521	59	257	117	4299	78	324
128	$512^2$	$u$	0	0	0	3	0	0	0	4
256	$512^2$	$u$	0	0	9	7	0	0	13	7
512	$512^2$	$u$	1	7	25	10	6	54	41	13
1024	$512^2$	$u$	5	15	24	13	9	73	32	21
2048	$512^2$	$u$	8	40	34	15	29	187	40	25
4096	$512^2$	$u$	6	36	29	14	19	109	34	37
128	$1024^2$	$u$	0	0	0	3	0	0	0	3
256	$1024^2$	$u$	0	0	0	5	0	0	0	5
512	$1024^2$	$u$	0	0	0	7	0	0	0	7
1024	$1024^2$	$u$	1	1	20	10	1	2	20	10
2048	$1024^2$	$u$	6	26	25	11	20	152	40	17
4096	$1024^2$	$u$	16	87	42	24	49	279	42	24

Table 2. Resulting values of our methods.

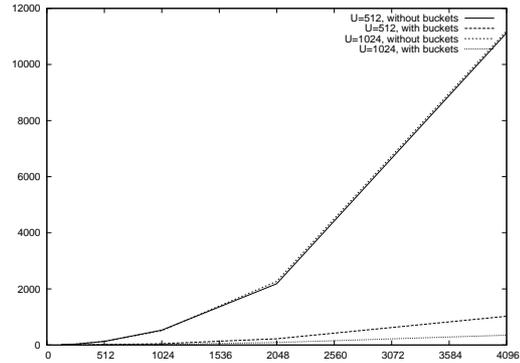
$n$	$U$	$CH$	without buckets		with buckets	
			non-varying	varying	non-varying	varying
128	$512^2$	$c$	9.44	9.56	2.35	2.86
256	$512^2$	$c$	40.91	41.89	12.58	14.69
512	$512^2$	$c$	203.53	212.88	102.74	119.15
1024	$512^2$	$c$	664.01	683.83	191.85	221.47
2048	$512^2$	$c$	3393.17	3457.26	972.34	1099.19
4096	$512^2$	$c$	14903.81	15046.69	5250.53	5651.59
128	$1024^2$	$c$	8.19	8.47	0.91	1.24
256	$1024^2$	$c$	32.79	33.82	2.83	3.63
512	$1024^2$	$c$	132.97	139.13	12.00	15.01
1024	$1024^2$	$c$	595.24	622.29	110.06	129.27
2048	$1024^2$	$c$	2809.12	2875.07	324.43	375.57
4096	$1024^2$	$c$	11143.28	11300.18	1477.70	1705.00
128	$512^2$	$u$	8.37	8.08	1.33	1.52
256	$512^2$	$u$	32.73	32.96	3.74	4.67
512	$512^2$	$u$	129.16	130.56	12.88	15.91
1024	$512^2$	$u$	529.79	523.12	47.54	54.52
2048	$512^2$	$u$	2184.42	2178.64	221.99	234.74
4096	$512^2$	$u$	11126.42	10978.71	1024.22	1037.39
128	$1024^2$	$u$	7.85	7.86	0.87	1.04
256	$1024^2$	$u$	31.45	32.79	2.28	3.01
512	$1024^2$	$u$	127.92	128.21	7.47	8.80
1024	$1024^2$	$u$	512.59	515.91	24.67	27.96
2048	$1024^2$	$u$	2268.77	2251.06	89.00	95.24
4096	$1024^2$	$u$	11201.63	11295.04	350.20	381.82

Table 3. Running times of our methods for the report-all problem. Reported times are in seconds.

**Without Buckets vs. With Buckets:** The approach to subdivide the space into buckets does not influence the reported values of our methods, however, it can have an impressive impact on the running times (see Figures 13 and 14). Depending on the input characteristics, we can observe speed-up factors between 2 and 32. The running time of the methods without ‘buckets’ is clearly quadratic in the number of entities. An asymptotic behaviour of the methods with ‘buckets’ is more difficult to identify, but note that also this method has a quadratic worst case running time.



**Fig. 13.** Running times depending on input size for non-varying report-all patterns for  $U = 512^2$ .



**Fig. 14.** Running times depending on input size for non-varying report-all patterns for  $CH = u$ .

**$CH = u$  vs.  $CH = c$ :** Almost always the input files with characteristic  $CH = c$  contain more patterns, longer patterns and larger patterns, which was expected as those files are much more likely to contain more and larger flocks. Hence, the input files with  $CH = u$  result in smaller running times (see also Figure 13). Interestingly these characteristics also indicate that the ‘bucket’ approach for speeding-up the computations has its limitations, because the speed-up factor of the ‘buckets’ method is strongly influenced by the characteristics. For  $CH = u$ , we observe speed-up factors around between 5 and 11 for the instances with  $U = 512^2$ , and between 7 and 32 for the instances with  $U = 1024^2$ . On the other hand, for  $CH = c$ , the speed-up factors are between 2 and 4 for the instances with  $U = 512^2$ , and between 5 and 11 for the instances with  $U = 1024^2$ . This can be explained by noting that the files with characteristic  $CH = c$  contain more and bigger flocks, and hence it is more likely that our algorithms encounter neighbouring buckets that are filled with more entities.

**$U = 512^2$  vs.  $U = 1024^2$ :** The difference between the universe with  $U = 512^2$  and  $U = 1024^2$  is that the former is much denser when filled with the same number of entities. As a result, in the larger universe ( $U = 1024^2$ ) less and smaller patterns exist. Also the running times are affected (see Figure 14). The methods with buckets run faster on instances with a larger universe, because we have more buckets and therefore, buckets are likely to contain less entities on average.

## 8 Discussion and Concluding Remarks

Movement patterns detect structure in large tracking data sets, and are thus key to a better understanding of the interactions amongst moving agents. In this paper we present the formal notion of a pattern called ‘leadership’, describing the event or process of one individual in front leading the movement of a group. Our approach is inspired by movement patterns documented in the animal behaviour and behavioural ecology literature.

The analysis of the interrelations of moving individuals has in the last five years attracted increasing attention, as a general reaction to the striking need for more powerful methods for surveillance and geospatial intelligence. Geographical information scientists are commissioned to develop methods that detect the expected and discover the unexpected from massive streams of disparate data, potentially originating from various sources [60]. Such methods need to be scalable, flexible and reliable.

Balancing the matching of formalised movement patterns (such as the presented leadership) with the inferring of unexpected space-time behaviours from visualised space-time paths, we argue that the former copes much better with increasing data sets. Whereas inferring from visualisation might be adequate for the analysis of individual events of interest [31], keeping track of hundreds of individuals cruising in the space-time aquarium is literally impossible [35]. By contrast, when detecting movement patterns such as flock or leadership, the number of entities  $n$  is just a performance factor but not an obfuscation factor.

Other authors recently used correlation analysis in order to discover leader and follower relationships amongst moving individuals. In Shirabe’s approach [56] such behaviour is detected by investigating pair-wise cross-correlation matrices with different time lags. Thus, lead-follow events can only be detected for pairs of individuals at a time. Our leadership pattern, in contrast, allows individuals to lead groups of followers. Since they operate on local-instantaneous events they can be detected in trajectories of variable lengths, as long as there is certain temporal overlap. Furthermore the approach in [56] has rather demanding constraints with respect to the analysed data set. It requires trajectories of equal length and strongly synchronous sampling. Even though we assumed the input data with the same characteristics, our algorithms for the continuous case can be easily applied to data without a synchronised sampling. The running times for sorting the sets of intervals for an entity would slightly increase, however, from  $O(n\tau \log n)$  to  $O(n\tau \log n\tau)$ . We argue that our leadership is thus a flexible concept, applicable to diverse data from various sources.

Movement patterns that are defined from the geometric arrangement of the involved entities (e.g. leadership), are more reliable than movement patterns that base on the intermediate step of an analysis matrix, as do the REMO patterns depend on an analysis matrix in Laube et al. [38]. The deterministic discretisation of the movement descriptors in eight cardinal direction classes introduces edge effects. An example shall illustrate such edge effects. Let  $22.5^\circ$  be one threshold of the discrete movement azimuth class ‘North’. Let furthermore the pattern under study be a flock pattern of four entities moving in the same direction at any time  $t$ . Why should a set of entities  $S_1$  with azimuths  $[21^\circ, 22^\circ, 22^\circ, 21^\circ]$  be a flock when another set  $S_2$  with azimuths  $[22^\circ, 23^\circ, 23^\circ, 22^\circ]$  is not? A definition requiring the entities to have a mean azimuth and some variance (e.g.  $\pm 22.5^\circ$ ) is a much more natural and thus reliable definition of flock. The definition of leadership in this paper follows for exactly the same reasons the road of using a geometrical arrangement instead of scanning a discretised matrix.

We provide a formal description of the pattern ‘leadership’ and subsequently algorithms for its efficient detection. Our experiments give indications which input-size can be processed within a reasonable amount of time. The resulting running times match the theoretical bounds, however for improved methods (with buckets) the running times strongly depend on the characteristics of the instance. When comparing the running times in this article with those reported in [4], we observe that the running times in the present work are much higher. This is because the used methods are different. The methods in [4] are faster but only report patterns of a specified length with a specified start- and end-time. The methods in this paper, however, are more flexible. Once the arrays *IntervalsNotFwg* and *IntervalsFwg* are computed we can very efficiently use them to report patterns of different lengths, and with different start- and end-times.

We also implemented the approximation algorithm and performed initial experiments. They show a better asymptotic behaviour of the approximation algorithm. However, the constant factors seem to be too large for practical purposes, because for our test-files the exact algorithms always outperformed the approximation algorithm.

One drawback of the given definition of leadership is that a leader has to be in the front region of all followers. For instance, for a very big flock of gnus this definition might not be applicable,

as some gnus at the end of the flock are too far away from the front-line to be able to see leading animals. Hence, one direction for future research could be the definition and analysis of cascading leaders or followers, where a cascading follower is a follower of a leader or a follower of another cascading follower.

For the many fields interested in movement, the overall challenge lies in relating movement patterns with the surrounding environment, in order to understand *where*, *when* and ultimately *why* the agents move the way they do. Conceptualising detectable movement patterns and the development of algorithms for their detection is a first important step towards this ambitious long-term goal. With its traditional spatial awareness, computational geometry can make immense contributions to the theoretical framework underlying movement analysis in geographical information science, behavioural ecology or surveillance and security analysis.

### Acknowledgements

The authors wish to thank Karin Schütz, AgResearch Ruakura, Hamilton, NZ for valuable comments on animal movement patterns and also Bojan Djordjevic for implementing the algorithms.

### References

1. P. K. Agarwal, L. Arge, and J. Erickson. Indexing moving points. *Journal of Computer and System Sciences*, 66(1):207–243, 2003.
2. G. Al-Naymat, S. Chawla, and J. Gudmundsson. Dimensionality reduction for long duration and complex spatio-temporal queries. accepted at ACM SAC, July 2007.
3. M. A. Bender and M. Farach-Colton. The LCA problem revisited. In *LATIN '00: Proceedings of the 4th Latin American Symposium on Theoretical Informatics*, volume 1776 of *Lecture Notes In Computer Science*, pages 88–94, London, UK, 2000. Springer-Verlag.
4. M. Benkert, J. Gudmundsson, F. Hübner, and T. Wolle. Reporting flock patterns. In *Proc. of the 14th Annual European Symposium on Algorithms (ESA 2006)*, volume 4168 of *Lecture Notes in Computer Science*, pages 660–671. Springer-Verlag, 2006.
5. J. L. Bentley. Multidimensional divide-and-conquer. *Communications of the ACM*, 23:214–229, 1980.
6. K. Börner and S. Penumarthy. Social diffusion patterns in three-dimensional virtual worlds. *Information Visualization*, 2(3):182–198, 2003.
7. S. M. C. Cavalcanti and F. F. Knowlton. Evaluation of physical and behavioral traits of llamas associated with aggressiveness toward sheep-threatening canids. *Applied Animal Behaviour Science*, 61(2):143–158, 1998.
8. N. R. Chrisman. Beyond the Snapshot: Changing the Approach to Change, Error, and Process. In M. J. Egenhofer and R. G. Golledge, editors, *Spatial and Temporal Reasoning in Geographic Information Systems*, pages 85–93. Oxford University Press, Oxford, UK, 1998.
9. L. Conradt and T. J. Roper. Group decision-making in animals. *Nature*, 421(6919):155–158, 2003.
10. M. D’Auria, M. Nanni, and D. Pedreschi. Time-focused density-based clustering of trajectories of moving objects. In *Proceedings of the Workshop on Mining Spatio-temporal Data (MSTD-2005)*, Porto, 2005.
11. M. de Berg, M. van Kreveld, M. Overmars, and O. Schwarzkopf. *Computational Geometry: algorithms and applications*. Springer, Berlin, 2nd edition, 2000.
12. C. Du Mouza and P. Rigaux. Mobility patterns. *GeoInformatica*, 9(4):297–319, 2005.
13. B. Dumont, A. Boissy, C. Achard, A. M. Sibbald, and H. W. Erhard. Consistency of animal order in spontaneous group movements allows the measurement of leadership in a group of grazing heifers. *Applied Animal Behaviour Science*, 95(1-2):55–66, 2005.
14. J. A. Dykes and D. M. Mountain. Seeking structure in records of spatio-temporal behaviour: visualization issues, efforts and application. *Computational Statistics and Data Analysis*, 43(4):581–603, 2003.
15. N. Eagle and A. Pentland. Reality mining: sensing complex social systems. *Personal and Ubiquitous Computing*, 10(4):255–268, 2006.
16. J. Erickson and R. Seidel. Better lower bounds on detecting affine and spherical degeneracies. *Discrete & Computational Geometry*, 13:41–57, 1995.
17. M. Erwig and R. H. Güting. Spatio-Temporal Data Types: An Approach to Modeling and Querying Moving Objects in Databases. *GeoInformatica*, 3(3):269–296, 1999.

18. A. U. Frank. Socio-Economic Units: Their Life and Motion. In A. U. Frank, J. Raper, and J. P. Cheylan, editors, *Life and motion of socio-economic units*, volume 8 of *GISDATA*, pages 21–34. Taylor & Francis, London, 2001.
19. A. Frihida, D. J. Marceau, and M. Theriault. Development of a Process-Based Model for Dynamic Interaction in Spatio-Temporal GIS. *GeoInformatica*, 82(3):211–235, 2004.
20. A. Gajentaan and M. H. Overmars.  $n^2$ -hard problems in computational geometry. Technical Report 1993-15, Department of Computer Science, Utrecht University, The Netherlands, 1993.
21. J. Gudmundsson and M. van Kreveld. Computing longest duration flocks in trajectory data. In *To appear in Proceedings 14th ACM Symposium on Advances in GIS*, 2006.
22. J. Gudmundsson, M. van Kreveld, and B. Speckmann. Efficient detection of motion patterns in spatio-temporal sets. *To appear in Geoinformatica*, 2006.
23. S. Gueron and S. A. Levin. Self-Organization of Front Patterns in Large Wildebeest Herds. *Journal of theoretical Biology*, 165(4):541–552, 1993.
24. R. Güting, M. H. Boehlen, M. Erwig, C. S. Jensen, N. Lorentzos, E. Nardelli, M. Schneider, and M. Vazirgiannis. A Foundation for Representing and Querying Moving Objects. *ACM Transactions on Database Systems (TODS)*, 2520(1):1–42, 2000.
25. R. Güting, M. H. Boehlen, M. Erwig, C. S. Jensen, N. Lorentzos, E. Nardelli, M. Schneider, and J. R. R. Viqueira. Spatio-temporal Models and Languages: An Approach Based on Data Types. In M. Koubarakis, T. Sellis, A. U. Frank, S. Grumbach, R. H. Gueting, C. S. Jensen, N. Lorentzos, Y. Manolopoulos, E. Nardelli, B. Pernici, H. J. Schek, M. Scholl, B. Theodoulidis, and N. Tryfona, editors, *Spatio-Temporal Databases: The CHOROCHRONOS Approach*, volume 2520 of *LNCS*, pages 117–176. Springer, Berlin, 2003.
26. R.H. Güting and M. Schneider. *Moving Objects Databases*. Morgan Kaufmann Publishers, 2005.
27. I. A. R. Hulbert. GPS and its Use in Animal Telemetry: The next five Years. In A. M. Sibbald and I. J. Gordon, editors, *Conference on Tracking Animals with GPS*, pages 51–60, Aberdeen, UK, 2001. Macaulay Insitute.
28. Y. Inada and K. Kawachi. Order and flexibility in the motion of fish schools. *Journal of theoretical Biology*, 214(3):371–387, 2002.
29. Y. Ishikawa, Y. Tsukamoto, and H. Kitagawa. Extracting mobility statistics from indexed spatio-temporal datasets. In *Proc. of the 2nd Workshop on Spatio-Temporal Database Management (STDBM)*, pages 9–16, 2004.
30. A. Jadbabaie, J. Lin, and A. S. Morse. Coordination of groups of mobile autonomous agents using nearest neighbor rules. *IEEE Transactions on Automatic Control*, 48(6):988–1001, 2003.
31. T. Kapler, R. Harper, and W. Wright. Correlating events with tracked movement in time and space: A geotime case study, 2005. Presented at the 2005 Intelligence Analysis Conference, Washington, DC.
32. G. Kollios, S. Sclaroff, and M. Betke. Motion mining: discovering spatio-temporal patterns in databases of human motion. In *Proceedings of the ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery*, 2001.
33. M. Koubarakis, Y. Theodoridis, and T. Sellis. Spatio-Temporal Databases in the Years Ahead. In M. Koubarakis, T. Sellis, A. U. Frank, S. Grumbach, R. H. Gueting, C. S. Jensen, N. Lorentzos, Y. Manolopoulos, E. Nardelli, B. Pernici, H. J. Schek, M. Scholl, B. Theodoulidis, and N. Tryfona, editors, *Spatio-Temporal Databases: The CHOROCHRONOS Approach*, volume 2520 of *LNCS*, pages 345–347. Springer, Berlin, 2003.
34. J. Krause and G. D. Ruxton. *Living in Groups*. Oxford Series in Ecology and Evolution. Oxford University Press, New York, NY, 2002.
35. M. P. Kwan. Interactive geovisualization of activity-travel patterns using three dimensional geographical information systems: a methodological exploration with a large data set. *Transportation Research Part C*, 8(1-6):185–203, 2000.
36. R. F. Lachlan, L. Crooks, and K. N. Laland. Who follows whom? Shoaling preferences and social learning of foraging information in guppies. *Animal Behaviour*, 56(1):181–190, 1998.
37. P. Laube and S. Imfeld. Analyzing relative motion within groups of trackable moving point objects. In M. J. Egenhofer and D. M. Mark, editors, *Geographic Information Science 2002*, volume 2478 of *LNCS*, pages 132–144, Berlin, 2002. Springer.
38. P. Laube, S. Imfeld, and R. Weibel. Discovering relative motion patterns in groups of moving point objects. *International Journal of Geographical Information Science*, 19(6):639–668, 2005.
39. P. Laube and R.S. Purves. An approach to evaluating motion pattern detection techniques in spatio-temporal data. *Computers, Environment and Urban Systems*, 30(3):347–374, 2006.
40. P. Laube, M. van Kreveld, and S. Imfeld. Finding REMO – detecting relative motion patterns in geospatial lifelines. In P. F. Fisher, editor, *Developments in Spatial Data Handling: Proceedings of the 11th International Symposium on Spatial Data Handling*, pages 201–214, Berlin, 2004. Springer.

41. J. B. Leca, N. Gunst, B. Thierry, and O. Petit. Distributed leadership in semifree-ranging white-faced capuchin monkeys. *Animal Behaviour*, 66:1045–1052, 2003.
42. Y. Li, J. Han, and J. Yang. Clustering moving objects. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, Seattle, WA, USA, 2004. ACM Press.
43. N. Mamoulis, H. Cao, G. Kollios, M. Hadjieleftheriou, Y. Tao, and D. W. Cheung. Mining, indexing, and querying historical spatiotemporal data. In *Proc. of the 10th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 236–245, 2004.
44. D. M. Mark. Geospatial Lifelines. In *Integrating Spatial and Temporal Databases*, volume 98471. Dagstuhl Seminars, 1998.
45. J.M. McNamara and A.I. Houston. The Value of fat reserves and the tradeoff between starvation and predation. *Acta Biotheoretica*, 38(1):37–61, 1990.
46. H. J. Miller and J. Han. Geographic data mining and knowledge discovery: An Overview. In H. J. Miller and J. Han, editors, *Geographic data mining and knowledge discovery*, pages 3–32. Taylor & Francis, London, UK, 2001.
47. K. Mouratidis, D. Papadias, and M. Hadjieleftheriou. Conceptual partitioning: an efficient method for continuous nearest neighbor monitoring. In *Proceedings of the 2005 ACM SIGMOD Conference on Management of Data*, pages 634–645, 2005.
48. R. T. Ng. Detecting outliers from large datasets. In H. J. Miller and J. Han, editors, *Geographic data mining and knowledge discovery*, pages 218–235. Taylor & Francis, London, UK, 2001.
49. R. O. Peterson, A. K. Jacobs, T. D. Drummer, L. D. Mech, and D. W. Smith. Leadership behavior in relation to dominance and reproductive status in gray wolves, *canis lupus*. *Canadian Journal of Zoology*, 80(8):1405–1412, 2002.
50. F. Porikli. Trajectory Distance Metric Using Hidden Markov Model based Representation. In *IEEE European Conference on Computer Vision, Workshop on PETS*, Prague, 2004.
51. Y. Qu, C. Wang, and X. S. Wang. Supporting fast search in time series for movement patterns in multiple scales. In *Seventh international conference on Information and knowledge management*, pages 251–258, Bethesda, Maryland, United States, 1998. ACM Press.
52. S. A. Rands, G. Cowlishaw, R. A. Pettifor, J. M. Rowcliffe, and R. A. Johnstone. Spontaneous emergence of leaders and followers in foraging pairs. *Nature*, 423(6938):432–434, 2003.
53. J. Raper. The Dimensions of GIScience, 2002. Keynote speech of GIScience 2002.
54. C.W. Reynolds. Flocks, herds and schools: A distributed behavioral model. In *Proceedings of the 14th annual conference on Computer graphics and interactive techniques*, volume 21, pages 25–34. ACM Press, 1987.
55. J. Roddick, K. Hornsby, and M. Spiliopoulou. An updated bibliography of temporal, spatial, and spatio-temporal data mining research. In *TSDM 2000*, volume 2007 of *LNCS*, pages 147–163, Berlin, 2001. Springer.
56. T. Shirabe. Correlation analysis of discrete motions. In *Proceedings of the Fourth International Conference on Geographic Information Science, GIScience 2006*, volume 4197 of *Lecture Notes In Computer Science*, pages 370–382, Berlin, 2006. Springer-Verlag.
57. G. Sinha and D. M. Mark. Measuring similarity between geospatial lifelines in studies of environmental health. *Journal of Geographical Systems*, 7(1):115–136, 2005.
58. A. P. Sistla, O. Wolfson, S. Chamberlain, and S. Dao. Modeling and Querying Moving Objects. In *13th International Conference on Data Engineering (ICDE13)*, 1997.
59. C. S. Smyth. Mining mobile trajectories. In H. J. Miller and J. Han, editors, *Geographic data mining and knowledge discovery*, pages 337–361. Taylor & Francis, London, UK, 2001.
60. James J. Thomas and Kristin A. Cook. A visual analytics agenda. *IEEE Computer Graphics and Applications*, 26(1):10–13, 2006.
61. S. Tisue and U. Wilensky. NetLogo: A Simple Environment for Modeling Complexity. In *International Conference on Complex Systems*, Boston, 2004.
62. F. Verhein and S. Chawla. Mining spatio-temporal association rules, sources, sinks, stationary regions and thoroughfares in object mobility databases. In *Proceedings of the 11th International Conference on Database Systems for Advanced Applications (DASFAA)*, 2006.
63. U. Wilensky. NetLogo Flocking model, 1998. <http://ccl.northwestern.edu/netlogo/models/Flocking>.
64. U. Wilensky. NetLogo (and NetLogo User Manual), 1999. <http://ccl.northwestern.edu/netlogo>.
65. O. Wolfson and E. Mena. Applications of Moving Objects Databases. In Y. Manolopoulos, A. Papadopoulos, and M. Vassilakopoulos, editors, *Spatial Databases: Technologies, Techniques and Trends*. Idea group Co., 2004.

66. O. Wolfson, B. Xu, S. Chamberlain, and L. Jiang. Moving Objects Databases: Issues and Solutions. In M. Rafanelli and M. Jarke, editors, *10th International Conference on Scientific and Statistical Database Management, Proceedings, Capri, Italy, July 1-3, 1998*, pages 111–131. IEEE Computer Society, 1998.
67. X. Xiong, M. F. Mokbel, and W. G. Aref. Sea-cnn: Scalable processing of continuous  $k$ -nearest neighbor queries in spatio-temporal databases. In *Proc. of the 21st International Conference on Data Engineering (ICDE 2005)*, pages 643–654, 2005.