



**HAL**  
open science

## An Efficient Spatio-Temporal Mining Approach to Really Know Who Travels with Whom!

Nhat Hai Phan, Pascal Poncelet, Maguelonne Teisseire

### ► To cite this version:

Nhat Hai Phan, Pascal Poncelet, Maguelonne Teisseire. An Efficient Spatio-Temporal Mining Approach to Really Know Who Travels with Whom!. BDA 2012 - 28e journées Bases de Données Avancées, Oct 2012, Clermont-Ferrand, France. lirmm-00798222

**HAL Id: lirmm-00798222**

<https://hal-lirmm.ccsd.cnrs.fr/lirmm-00798222v1>

Submitted on 11 Dec 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# An Efficient Spatio-Temporal Mining Approach to Really Know Who Travels with Whom!

Phan Nhat Hai<sup>1,2</sup>, Pascal Poncelet<sup>1,2</sup>, and Maguelonne Teisseire<sup>1,2</sup>

<sup>1</sup> IRSTEA Montpellier, UMR TETIS - 34093 Montpellier, France

{nhat-hai.phan, maguelonne.teisseire}@teledetection.fr

<sup>2</sup> LIRMM CNRS Montpellier - 34090 Montpellier, France pascal.poncelet@lirmm.fr

**Abstract.** Recent improvements in positioning technology has led to a much wider availability of massive moving object data. A crucial task is to find the moving objects that travel together. Usually, they are called spatio-temporal patterns. Due to the emergence of many different kinds of spatio-temporal patterns in recent years, different approaches have been proposed to extract them. However, each approach only focuses on mining a specific kind of pattern. In addition to the fact that it is a painstaking task due to the large number of algorithms used to mine and manage patterns, it is also time consuming. Additionally, we have to execute these algorithms again whenever new data are added to the existing database. To address these issues, we first redefine spatio-temporal patterns in the itemset context. Secondly, we propose a unifying approach, named *GeT\_Move*, using a frequent closed itemset-based spatio-temporal pattern-mining algorithm to mine and manage different spatio-temporal patterns. *GeT\_Move* is implemented in two versions which are *GeT\_Move* and *Incremental GeT\_Move*. Experiments are performed on real and synthetic datasets and the experimental results show that our approaches are very effective and outperform existing algorithms in terms of efficiency.

**Keywords:** Spatio-temporal pattern, frequent closed itemset, trajectories

## 1 Introduction

Nowadays, many electronic devices are used for real world applications. Telemetry attached on wildlife, GPS installed in cars, sensor networks, and mobile phones have enabled the tracking of almost any kind of data and has led to an increasingly large amount of data that contain moving objects and numerical data. Therefore, analysis on such data to find interesting patterns is attracting increasing attention for applications such as movement pattern analysis, animal behavior study, route planning and vehicle control.

Early approaches designed to recover information from spatio-temporal datasets included ad-hoc queries aimed as answering queries concerning a single predicate range or nearest neighbour. For instance, "*finding all the moving objects inside area A between 10:00 am and 2:00 pm*" or "*how many cars were driven between Main Square and the Airport on Friday*" [7]. Spatial query extensions in GIS applications are able to run this type of query. However, these techniques are used to find the best solution by

exploring each spatial object at a specific time according to some metric distance measurement (usually Euclidean). As results, it is difficult to capture collective behaviour and correlations among the involved entities using this type of queries.

Recently, many spatio-temporal patterns have been proposed [1, 3, 4, 6, 11, 15, 16, 17, 18]. In this paper, we are interested in the querying of patterns which capture 'group' or 'common' behaviour among moving entities. This is particularly true to identify groups of moving objects for which a strong relationship and interaction exist within a defined spatial region during a given time duration. Some examples of these patterns are flocks [1, 2], moving clusters [4, 12], convoy queries [3, 10], closed swarms [6, 9], group patterns [15], periodic patterns [18], etc...

To extract these kinds of patterns, different algorithms have been proposed. Naturally, the computation is costly and time consuming because we need to execute different algorithms consecutively. However, if we had an algorithm which could extract different kinds of patterns, the computation costs will be significantly decreased and the process would be much less time consuming. Therefore, we need to develop an efficient unifying algorithm.

In some real world applications (e.g. cars), object locations are continuously reported by using Global Positioning System (GPS). Therefore, new data is always available. If we do not have an incremental algorithm, we need to execute again and again algorithms on the whole database including existing data and new data to extract patterns. This is of course, cost-prohibitive and time consuming. An incremental algorithm can indeed improve the process by combining the results extracted from the existing data and the new data to obtain the final results.

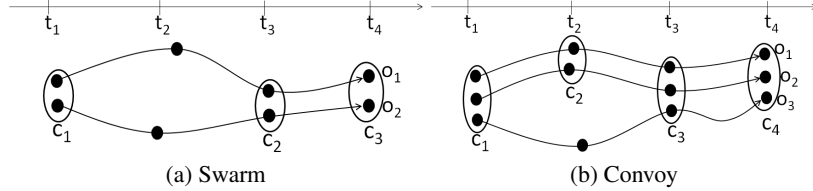
With the above issues in mind, we propose *GeT\_Move*: a unifying incremental spatio-temporal pattern-mining approach. Part of this approach is based on frequent closed itemset mining algorithm (FCI). The main idea of the algorithm and the main contributions of this paper are summarized below.

- We re-define the spatio-temporal patterns mining in the itemset context which enable us to effectively extract different kinds of spatio-temporal patterns.
- We present approaches, called *GeT\_Move* and *Incremental GeT\_Move*, which efficiently extract FCIs from which spatio-temporal patterns are retrieved.
- We present comprehensive experimental results over both real and synthetic databases. The results demonstrate that our techniques enable us to effectively extract different kinds of patterns. Furthermore, our approaches are more efficient compared to other algorithms in most of cases.

The remaining sections of the paper are organized as follows. Section 2 discusses preliminary definitions of the spatio-temporal patterns as well as the related work. The properties of these patterns are provided in an itemset context in Section 3. We introduce the *GeT\_Move* and *Incremental GeT\_Move* algorithms in Section 4. Experiments testing effectiveness and efficiency are shown in Section 5. Finally, we draw our conclusions in Section 6.

**Table 1.** An example of a Spatio-Temporal Database

Objects $O_{DB}$	Timesets $T_{DB}$	$x$	$y$
$o_1$	$t_1$	2.3	1.2
$o_2$	$t_1$	2.1	1
$o_1$	$t_2$	10.3	28.1
$o_2$	$t_2$	0.3	1.2

**Fig. 1.** An example of swarm and convoy where  $c_1, c_2, c_3, c_4$  are clusters which gather closed objects together at specific timestamps.

## 2 Spatio-Temporal Patterns

In this section we briefly propose an overview of the main spatio-temporal patterns. We thus define the different kinds of patterns and then we discuss the related work.

### 2.1 Preliminary Definitions

The problem of spatio-temporal patterns has been extensively addressed over the last years. Basically, spatio-temporal patterns are designed to group similar trajectories or objects which tend to move together during a time interval. So many different definitions can be proposed and today lots of patterns have been defined such as *flocks* [1, 2], *convoys* [3, 10], *swarms*, *closed swarms* [6, 9], *moving clusters* [4, 12], *group pattern* [15] and even *periodic patterns* [18].

In this paper, we focus on proposing a unifying approach to effectively and efficiently extract all these different kinds of patterns. First of all, we assume that we have a group of moving objects  $O_{DB} = \{o_1, o_2, \dots, o_z\}$ , a set of timestamps  $T_{DB} = \{t_1, t_2, \dots, t_n\}$  and at each timestamp  $t_i \in T_{DB}$ , spatial information<sup>3</sup>  $x, y$  for each object. For example, Table 1 illustrates an example of a spatio-temporal database. Usually, in spatio-temporal mining, we are interested in extracting a group of objects staying together during a period. Therefore, from now,  $O = \{o_{i_1}, o_{i_2}, \dots, o_{i_p}\} (O \subseteq O_{DB})$  stands for a group of objects,  $T = \{t_{a_1}, t_{a_2}, \dots, t_{a_m}\} (T \subseteq T_{DB})$  is the set of timestamps within which objects stay together. Let  $\varepsilon$  be a user-defined threshold standing for a minimum number of objects and  $min_t$  a minimum number of timestamps. Thus  $|O|$  (resp.  $|T|$ ) must be greater than or equal to  $\varepsilon$  (resp.  $min_t$ ). In the following, we formally define all the different kinds of patterns.

Informally, a *swarm* is a group of moving objects  $O$  containing at least  $\varepsilon$  individuals which are closed each other for at least  $min_t$  timestamps. Then a swarm can be formally defined as follows:

<sup>3</sup> Spatial information can be for instance GPS location.

**Definition 1** *Swarm* [6]. A pair  $(O, T)$  is a swarm if:

$$\left\{ \begin{array}{l} (1) : \forall t_{a_i} \in T, \exists c \text{ s.t. } O \subseteq c, c \text{ is a cluster.} \\ \text{There is at least one cluster containing} \\ \text{all the objects in } O \text{ at each timestamp in } T. \\ (2) : |O| \geq \varepsilon. \\ \text{There must be at least } \varepsilon \text{ objects.} \\ (3) : |T| \geq \min_t. \\ \text{There must be at least } \min_t \text{ timestamps.} \end{array} \right. \quad (1)$$

For example, as shown in Figure 1a, if we set  $\varepsilon = 2$  and  $\min_t = 2$ , we can find the following swarms  $(\{o_1, o_2\}, \{t_1, t_3\})$ ,  $(\{o_1, o_2\}, \{t_1, t_4\})$ ,  $(\{o_1, o_2\}, \{t_3, t_4\})$ ,  $(\{o_1, o_2\}, \{t_1, t_3, t_4\})$ . We can note that these swarms are in fact redundant since they can be grouped together in the following swarm  $(\{o_1, o_2\}, \{t_1, t_3, t_4\})$ .

To avoid this redundancy, Zhenhui Li et al. [6] propose the notion of *closed swarm* for grouping together both objects and time. A swarm  $(O, T)$  is *object-closed* if when fixing  $T$ ,  $O$  cannot be enlarged. Similarly, a swarm  $(O, T)$  is *time-closed* if when fixing  $O$ ,  $T$  cannot be enlarged. Finally, a swarm  $(O, T)$  is a closed swarm if it is both object-closed and time-closed and can be defined as follows:

**Definition 2** *Closed Swarm* [6]. A pair  $(O, T)$  is a closed swarm if:

$$\left\{ \begin{array}{l} (1) : (O, T) \text{ is a swarm.} \\ (2) : \nexists O' \text{ s.t. } (O', T) \text{ is a swarm and } O \subset O'. \\ (3) : \nexists T' \text{ s.t. } (O, T') \text{ is a swarm and } T \subset T'. \end{array} \right. \quad (2)$$

For instance, in the previous example,  $(\{o_1, o_2\}, \{t_1, t_3, t_4\})$  is a closed swarm.

A *convoy* is also a group of objects such that these objects are closed each other during at least  $\min_t$  time points. The main difference between convoy and closed swarm is that convoy lifetimes must be consecutive:

**Definition 3** *Convoy* [3]. A pair  $(O, T)$ , is a convoy if:

$$\left\{ \begin{array}{l} (1) : (O, T) \text{ is a swarm.} \\ (2) : \forall i, 1 \leq i < |T|, t_{a_i}, t_{a_{i+1}} \text{ are consecutive.} \end{array} \right. \quad (3)$$

For instance, on Figure 1b, with  $\varepsilon = 2$ ,  $\min_t = 2$  we have two convoys  $(\{o_1, o_2\}, \{t_1, t_2, t_3, t_4\})$  and  $(\{o_1, o_2, o_3\}, \{t_3, t_4\})$ .

Until now, we have considered that we have a group of objects that move close to each other for a long time interval. For instance, as shown in [21], moving clusters and different kinds of flocks virtually share essentially the same definition. Basically, the main difference is based on the clustering techniques used. Flocks usually consider a rigid definition of the radius while moving clusters and convoys apply a density-based clustering algorithm (e.g. DBScan [5]). Moving clusters can be seen as special cases of

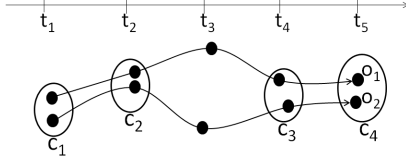


Fig. 2. A group pattern example.

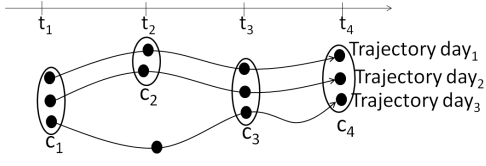


Fig. 3. A periodic pattern example.

convoys with the additional condition that they need to share some objects between two consecutive timestamps [21]. Therefore, in the following, for brevity and clarity sake we will mainly focus on convoy and density-based clustering algorithms.

According to the previous definitions, the main difference between convoys and swarms is about the consecutiveness and non-consecutiveness of clusters during a time interval. In [15], Hwang et al. propose a general pattern, called a *group pattern*, which essentially is a combination of both convoys and closed swarms. Basically, group pattern is a set of disjointed convoys which are generated by the same group of objects in different time intervals. By considering a convoy as a timepoint, a group pattern can be seen as a swarm of disjointed convoys. Additionally, group pattern cannot be enlarged in terms of objects and number of convoys. Therefore, group pattern is essentially a closed swarm of disjointed convoys. Formally, group pattern can be defined as follows:

**Definition 4** *Group Pattern* [15]. Given a set of objects  $O$ , a minimum weight threshold  $min_{wei}$ , a set of disjointed convoys  $T_S = \{s_1, s_2, \dots, s_n\}$ , a minimum number of convoys  $min_c$ .  $(O, T_S)$  is a group pattern if:

$$\begin{cases} (1) : (O, T_S) \text{ is a closed swarm with } \varepsilon, min_c. \\ (2) : \frac{\sum_{i=1}^{|T_S|} |s_i|}{|T_{DB}|} \geq min_{wei}. \end{cases} \quad (4)$$

Note that  $min_c$  is only applied for  $T_S$  (e.g.  $|T_S| \geq min_c$ ).

For instance, see Figure 2, with  $min_t = 2$  and  $\varepsilon = 2$  we have a set of convoys  $T_S = \{(\{o_1, o_2\}, \{t_1, t_2\}), (\{o_1, o_2\}, \{t_4, t_5\})\}$ . Additionally, with  $min_c = 1$  we have  $(\{o_1, o_2\}, T_S)$  is a closed swarm of convoys because  $|T_S| = 2 \geq min_c$ ,  $|O| \geq \varepsilon$  and  $(O, T_S)$  cannot be enlarged. Furthermore, with  $min_{wei} = 0.5$ ,  $(O, T_S)$  is a group pattern since  $\frac{|\{t_1, t_2\}| + |\{t_4, t_5\}|}{|T_{DB}|} = \frac{4}{5} \geq min_{wei}$ .

Previously, we overviewed patterns in which group objects move together during some time intervals. However, mining patterns from individual object movement is also interesting. In [18], N. Mamoulis et al. propose the notion of *periodic patterns* in which an object follows the same routes (approximately) over regular time intervals. For example, people wake up at the same time and generally follow the same route to their work everyday. Informally, given an object's trajectory including  $\mathcal{N}$  timepoints,  $\mathcal{T}_P$  which is the number of timestamps that a pattern may re-appear. An object's trajectory is decomposed into  $\lfloor \frac{\mathcal{N}}{\mathcal{T}_P} \rfloor$  sub-trajectories.  $\mathcal{T}_P$  is data-dependent and has no definite value. For example,  $\mathcal{T}_P$  can be set to 'a day' in traffic control applications since many vehicles have daily patterns, while annual animal migration patterns can be discovered by  $\mathcal{T}_P = \text{'a year'}$ . For instance, see Figure 3, an object's trajectory is decomposed into daily sub-trajectories.

Essentially, a periodic pattern is a closed swarm discovered from  $\lfloor \frac{N}{T_p} \rfloor$  sub-trajectories. For instance, in Figure 3, we have 3 daily sub-trajectories and from them we extract the two following periodic patterns  $\{c_1, c_2, c_3, c_4\}$  and  $\{c_1, c_3, c_4\}$ . The main difference in periodic pattern mining is the preprocessing data step while the definition is similar to that of a closed swarm. As we have provided the definition of a closed swarm, we will mainly focus on closed swarm mining below.

## 2.2 Related Work

As we mentioned before, many approaches have been proposed to extract patterns. The interested readers may refer to [14, 21] where short descriptions of the most efficient or interesting patterns and approaches are proposed. For instance, Gudmundsson and van Kreveld [1], Vieira et al. [2] define a flock pattern, in which the same set of objects stay together in a circular region with a predefined radius, Kalnis et al. [4] propose the notion of *moving clusters*, while Jeung et al. [3] define a convoy pattern.

Jeung et al. [3] adopt the DBScan algorithm [5] to find candidate convoy patterns. The authors propose three algorithms that incorporate trajectory simplification techniques in the first step. The distance measurements are performed on trajectory segments of as opposed to point based distance measurements. Another problem is related to the trajectory representation. Some trajectories may have missing timestamps or are measured at different time intervals. Therefore, the density measurements cannot be applied between trajectories with different timestamps. To address the problem of missing timestamps, the authors proposed to interpolate the trajectories by creating virtual time points and by applying density measurements on trajectory segments. Additionally, the convoy is defined as a candidate when it has at least  $k$  clusters during  $k$  consecutive timestamps.

Recently, Zhenhui Li et al. [6] propose the concept of swarm and closed swarm and the *ObjectGrowth* algorithm to extract closed swarm patterns. The *ObjectGrowth* method is a depth-first-search framework based on the objectset search space (i.e., the collection of all subsets of  $O_{DB}$ ). For the search space of  $O_{DB}$ , they perform depth-first search of all subsets of  $O_{DB}$  through a pre-order tree traversal. Even though, the search space remains still huge for enumerating the objectsets in  $O(2^{|O_{DB}|})$ . To speed up the search process, they propose two pruning rules. The first pruning rule, called *Apriori Pruning*, is used to stop traversal the subtree when we find further traversal that cannot satisfy  $min_t$ . The second pruning rule, called *Backward Pruning*, makes use of the closure property. It checks whether there is a superset of the current objectset, which has the same maximal corresponding timeset as that of the current one. If so, the traversal of the subtree under the current objectset is meaningless. After pruning the invalid candidates, the remaining ones may or may not be closed swarms. Then a *Forward Closure Checking* is used to determine whether a pattern is a closed swarm.

In [15], Hwang et al. propose two algorithms to mine group patterns, known as the *Apriori-like Group Pattern mining* algorithm and *Valid Group-Growth* algorithm. The former explores the Apriori property of valid group patterns and extends the Apriori algorithm [8] to mine valid group patterns. The latter is based on idea similar to the FP-growth algorithm [20]. Recently in [7], A. Calmeron proposes a frequent itemset-based approach for flock identification purposes.

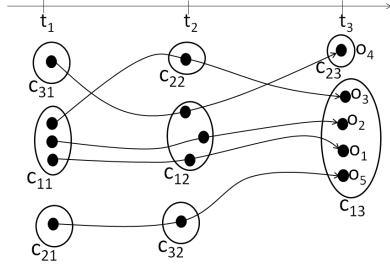


Fig. 4. An illustrative example.

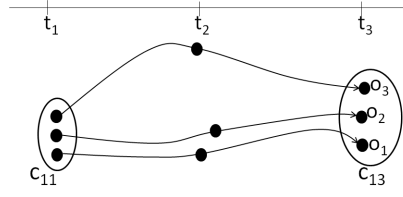


Fig. 5. A swarm from our example.

Table 2. Cluster Matrix.

$T_{DB}$		$t_1$			$t_2$			$t_3$	
Clusters	$C_{DB}$	$c_{11}$	$c_{21}$	$c_{31}$	$c_{12}$	$c_{22}$	$c_{32}$	$c_{13}$	$c_{23}$
$O_{DB}$	$o_1$	1			1			1	
	$o_2$	1			1			1	
	$o_3$	1				1		1	
	$o_4$			1	1				1
	$o_5$		1				1	1	

Even if these approaches are very efficient they suffer the problem that they only extract a specific kind of pattern. When considering a dataset, it is quite difficult, for the decision maker, to know in advance the kind of patterns embedded in the data. Therefore proposing an approach able to automatically extract all these different kinds of patterns can be very useful and this is the problem we address in this paper and that will be developed in the next sections.

### 3 Spatio-Temporal Patterns in Itemset Context

Extracting different kinds of patterns requires the use of several algorithms and to deal with this problem, we propose an unifying approach to extract and manage different kinds of patterns.

Basically, patterns are evolution of clusters over time. Therefore, to manage the evolution of clusters, we need to analyse the correlations between them. Furthermore, if clusters share some characteristics (e.g. share some objects), they could be a pattern. Consequently, if a cluster is considered as an item we will have a set of items (called itemset). The main problem essentially is to efficiently combine items (clusters) to find itemsets (a set of clusters) which share some characteristics or satisfy some properties to be considered as a pattern. To describe cluster evolution, spatio-temporal data is presented as a cluster matrix from which patterns can be extracted.

**Definition 5** *Cluster Matrix.* Assume that we have a set of clusters  $C_{DB} = \{C_1, C_2, \dots, C_n\}$  where  $C_i = \{c_{i_1 t_i}, c_{i_2 t_i}, \dots, c_{i_m t_i}\}$  is a set of clusters at timestamps  $t_i$ . A cluster matrix is thus a matrix of size  $|O_{DB}| \times |C_{DB}|$ . Each row represents an object and each column represents a cluster. The value of the cluster matrix cell,  $(o_i, c_j)$  is 1 (resp. empty) if  $o_i$  is in (resp. is not in) cluster  $c_j$ . A cluster (or item)  $c_j$  is a cluster formed after applying clustering techniques.



For instance, the data from Figure 4 is presented in a cluster matrix in Table 2. Object  $o_1$  belongs to the cluster  $c_{11}$  at timestamp  $t_1$ . For clarity reasons in the following,  $c_{ij}$  represents the cluster  $c_i$  at time  $t_j$ . Therefore, the matrix cell  $(o_1-c_{11})$  is 1, meanwhile the matrix cell  $(o_4-c_{11})$  is empty because object  $o_4$  does not belong to cluster  $c_{11}$ .

By presenting data in a cluster matrix, each object acts as a transaction while each cluster  $c_j$  stands for an item. Additionally, an itemset can be formed as  $\mathcal{Y} = \{c_{t_{a_1}}, c_{t_{a_2}}, \dots, c_{t_{a_p}}\}$  with life time  $T_{\mathcal{Y}} = \{t_{a_1}, t_{a_2}, \dots, t_{a_p}\}$  where  $t_{a_1} < t_{a_2} < \dots < t_{a_p}$ ,  $\forall a_i : t_{a_i} \in T_{DB}, c_{t_{a_i}} \in C_{a_i}$ . The support of the itemset  $\mathcal{Y}$ , denoted  $\sigma(\mathcal{Y})$ , is the number of common objects in every items belonging to  $\mathcal{Y}$ ,  $O(\mathcal{Y}) = \bigcap_{i=1}^p c_{t_{a_i}}$ . Additionally, the length of  $\mathcal{Y}$ , denoted  $|\mathcal{Y}|$ , is the number of items or timestamps ( $= |T_{\mathcal{Y}}|$ ).

For instance, in Table 2, for a support value of 2 we have:  $\mathcal{Y} = \{c_{11}, c_{12}\}$  verifying  $\sigma(\mathcal{Y}) = 2$ . Every items (resp. clusters) of  $\mathcal{Y}$ ,  $c_{11}$  and  $c_{12}$ , are in the transactions (resp. objects)  $o_1, o_2$ . The length of  $|\mathcal{Y}|$  is the number of items ( $= 2$ ).

Naturally, the number of clusters can be large; however, the maximum length of itemsets is  $|T_{DB}|$ . Because of the density-based clustering algorithm used, clusters at the same timestamp cannot be in the same itemsets.

Now, we will define some useful properties to extract the patterns presented in Section 2 from frequent itemsets as follows:

*Property 1. Swarm.* Given a frequent itemset  $\mathcal{Y} = \{c_{t_{a_1}}, c_{t_{a_2}}, \dots, c_{t_{a_p}}\}$ .  $(O(\mathcal{Y}), T_{\mathcal{Y}})$  is a swarm if and only if:

$$\begin{cases} (1) : \sigma(\mathcal{Y}) \geq \varepsilon \\ (2) : |\mathcal{Y}| \geq \min_t \end{cases} \quad (5)$$

*Proof.* After construction, we have  $\sigma(\mathcal{Y}) \geq \varepsilon$  and  $\sigma(\mathcal{Y}) = |O(\mathcal{Y})|$  then  $|O(\mathcal{Y})| \geq \varepsilon$ . Additionally, as  $|\mathcal{Y}| \geq \min_t$  and  $|\mathcal{Y}| = |T_{\mathcal{Y}}|$  then  $|T_{\mathcal{Y}}| \geq \min_t$ . Furthermore,  $\forall t_{a_j} \in T_{\mathcal{Y}}, O(\mathcal{Y}) \subseteq c_{t_{a_j}}$ , means that at every timestamp we have a cluster containing all objects in  $O(\mathcal{Y})$ . Consequently,  $(O(\mathcal{Y}), T_{\mathcal{Y}})$  is a swarm because it satisfies all the requirements of the *Definition 1*.

For instance, in Figure 5, for the frequent itemset  $\mathcal{Y} = \{c_{11}, c_{13}\}$  we have  $(O(\mathcal{Y}) = \{o_1, o_2, o_3\}, T_{\mathcal{Y}} = \{t_1, t_3\})$  which is a swarm with support threshold  $\varepsilon = 2$  and  $\min_t = 2$ . We can notice that  $\sigma(\mathcal{Y}) = 3 > \varepsilon$  and  $|\mathcal{Y}| = 2 \geq \min_t$ .

Essentially, a closed swarm is a swarm which satisfies the *object-closed* and *time-closed* conditions therefore closed-swarm property is as follows:

*Property 2. Closed Swarm.* Given a frequent itemset  $\mathcal{Y} = \{c_{t_{a_1}}, c_{t_{a_2}}, \dots, c_{t_{a_p}}\}$ .  $(O(\mathcal{Y}), T_{\mathcal{Y}})$  is a closed swarm if and only if:

$$\begin{cases} (1) : (O(\mathcal{Y}), T_{\mathcal{Y}}) \text{ is a swarm.} \\ (2) : \nexists \mathcal{Y}' \text{ s.t. } O(\mathcal{Y}) \subset O(\mathcal{Y}'), T_{\mathcal{Y}'} = T_{\mathcal{Y}} \text{ and} \\ \quad (O(\mathcal{Y}'), T_{\mathcal{Y}}) \text{ is a swarm.} \\ (3) : \nexists \mathcal{Y}' \text{ s.t. } O(\mathcal{Y}') = O(\mathcal{Y}), T_{\mathcal{Y}} \subset T_{\mathcal{Y}'} \text{ and} \\ \quad (O(\mathcal{Y}), T_{\mathcal{Y}'}) \text{ is a swarm.} \end{cases} \quad (6)$$

*Proof.* After construction, we obtain  $(O(\mathcal{Y}), T_{\mathcal{Y}})$  which is a swarm. Additionally, if  $\nexists \mathcal{Y}'$  s.t.  $O(\mathcal{Y}) \subset O(\mathcal{Y}'), T_{\mathcal{Y}'} = T_{\mathcal{Y}}$  and  $(O(\mathcal{Y}'), T_{\mathcal{Y}})$  is a swarm then  $(O(\mathcal{Y}), T_{\mathcal{Y}})$  cannot be enlarged in terms of objects. Therefore, it satisfies the *object-closed* condition. Furthermore, if  $\nexists \mathcal{Y}'$  s.t.  $O(\mathcal{Y}') = O(\mathcal{Y}), T_{\mathcal{Y}} \subset T_{\mathcal{Y}'}$  and  $(O(\mathcal{Y}'), T_{\mathcal{Y}'})$  is a swarm then  $(O(\mathcal{Y}), T_{\mathcal{Y}})$  cannot be enlarged in terms of lifetime. Therefore, it satisfies the *time-closed* condition. Consequently,  $(O(\mathcal{Y}), T_{\mathcal{Y}})$  is a swarm and it satisfies *object-closed* and *time-closed* conditions and therefore  $(O(\mathcal{Y}), T_{\mathcal{Y}})$  is a closed swarm according to the *Definition 2*.

Due to space limitation, we do not provide the properties and proof for convoys, moving clusters which are basically extended by adding some conditions to *Property 1*. For instance, a convoy is a swarm which satisfies the consecutiveness in terms of time condition. For moving clusters [4], they need to share some objects between two timestamps (integrity proportion). Regarding to periodic patterns, the main difference in periodic pattern mining is the input data while the property is similar to *Property 2*. With a slightly modifying cluster matrix such as "each object  $o$  becomes a sub-trajectory", we can extract periodic patterns by applying *Property 2*.

Please remember that group pattern is a set of disjointed convoys which share the same objects, but in different time intervals. Therefore, the group pattern property is as follows:

*Property 3. Group Pattern.* Given a frequent itemset  $\mathcal{Y} = \{c_{t_{a_1}}, c_{t_{a_2}}, \dots, c_{t_{a_p}}\}$ , a minimum weight  $min_{wei}$ , a minimum number of convoys  $min_c$ , a set of consecutive time segments  $T_S = \{s_1, s_2, \dots, s_n\}$ .  $(O(\mathcal{Y}), T_S)$  is a group pattern if and only if:

$$\begin{cases} (1) : |T_S| \geq min_c. \\ (2) : \forall s_i, s_i \subseteq T_{\mathcal{Y}}, |s_i| \geq min_t. \\ (3) : \bigcap_{i=1}^n s_i = \emptyset, \bigcap_{i=1}^n O(s_i) = O(\mathcal{Y}). \\ (4) : \forall s \notin T_S, s \text{ is a convoy}, O(\mathcal{Y}) \not\subseteq O(s). \\ (5) : \frac{\sum_{i=1}^n |s_i|}{|T|} \geq min_{wei}. \end{cases} \quad (7)$$

*Proof.* If  $|T_S| \geq min_c$  then we know that at least  $min_c$  consecutive time intervals  $s_i$  in  $T_S$ . Furthermore, if  $\forall s_i, s_i \subseteq T_{\mathcal{Y}}$  then we have  $O(\mathcal{Y}) \subseteq O(s_i)$ . Additionally, if  $|s_i| \geq min_t$  then  $(O(\mathcal{Y}), s_i)$  is a convoy (*Definition 3*). Now,  $T_S$  actually is a set of convoys of  $O(\mathcal{Y})$  and if  $\bigcap_{i=1}^n s_i = \emptyset$  then  $T_S$  is a set of disjointed convoys. A little bit further, if  $\forall s \notin T_S, s$  is a convoy and  $O(\mathcal{Y}) \not\subseteq O(s)$  then  $\nexists T_{S'}$  s.t.  $T_S \subset T_{S'}$  and  $\bigcap_{i=1}^{|T_{S'}|} O(s_i) = O(\mathcal{Y})$ . Therefore,  $(O(\mathcal{Y}), T_S)$  cannot be enlarged in terms of *number of convoys*. Similarly, if  $\bigcap_{i=1}^n O(s_i) = O(\mathcal{Y})$  then  $(O(\mathcal{Y}), T_S)$  cannot be enlarged in terms of *objects*. Consequently,  $(O(\mathcal{Y}), T_S)$  is a closed swarm of disjointed convoys because  $|O(\mathcal{Y})| \geq \varepsilon, |T_S| \geq min_c$  and  $(O(\mathcal{Y}), T_S)$  cannot be enlarged (*Definition 2*). Finally, if  $(O(\mathcal{Y}), T_S)$  satisfies condition (5) then it is a valid group pattern due to *Definition 4*.

Above, we presented some useful properties to extract spatio-temporal patterns from itemsets. Now we will focus on the fact that from an itemset mining algorithm we are able to extract the set of all spatio-temporal patterns. We thus start the proof process by analyzing the swarm extracting problem. This first lemma shows that from a set of frequent itemsets we are able to extract all the swarms embedded in the database.

**Lemma 1.** Let  $FI = \{\mathcal{Y}_1, \mathcal{Y}_2, \dots, \mathcal{Y}_l\}$  be the frequent itemsets being mined from the cluster matrix with  $minsup = \varepsilon$ . All swarms  $(O, T)$  can be extracted from  $FI$ .

*Proof.* Let us assume that  $(O, T)$  is a swarm. Note,  $T = \{t_{a_1}, t_{a_2}, \dots, t_{a_m}\}$ . According to the *Definition 1* we know that  $|O| \geq \varepsilon$ . If  $(O, T)$  is a swarm then  $\forall t_{a_i} \in T, \exists c_{t_{a_i}}$  s.t.  $O \subseteq c_{t_{a_i}}$  therefore  $\bigcap_{i=1}^m c_{t_{a_i}} = O$ . Additionally, we know that  $\forall c_{t_{a_i}}, c_{t_{a_i}}$  is an item so  $\exists \mathcal{Y} = \bigcup_{i=1}^m c_{t_{a_i}}$  is an itemset and  $O(\mathcal{Y}) = \bigcap_{i=1}^m c_{t_{a_i}} = O, T_{\mathcal{Y}} = \bigcup_{i=1}^m t_{a_i} = T$ . Therefore,  $(O(\mathcal{Y}), T_{\mathcal{Y}})$  is a swarm. So,  $(O, T)$  is extracted from  $\mathcal{Y}$ . Furthermore,  $\sigma(\mathcal{Y}) = |O(\mathcal{Y})| = |O| \geq \varepsilon$  then  $\mathcal{Y}$  is a frequent itemset and  $\mathcal{Y} \in FI$ . Finally,  $\forall (O, T)$  s.t. if  $(O, T)$  is a swarm then  $\exists \mathcal{Y}$  s.t.  $\mathcal{Y} \in FI$  and  $(O, T)$  can be extracted from  $\mathcal{Y}$ , we can conclude that  $\forall$  a swarm  $(O, T)$ , it can be mined from  $FI$ .

We can consider that by adding constraints such as "consecutive lifetime", "time-closed", "object-closed", "integrity proportion" to swarms, we can retrieve convoys, closed swarms and moving clusters. Therefore, if *Swarm*, *CSwarm*, *Convoy*, *MCluster* respectively contain all swarms, closed-swarms, convoys and moving clusters then we have:  $CSwarm \subseteq Swarm$ ,  $Convoy \subseteq Swarm$  and  $MCluster \subseteq Swarm$ . By applying *Lemma 1*, we retrieve all swarms from frequent itemsets. Since, a set of closed swarms, a set of convoys and a set of moving clusters are subsets of swarms and they can therefore be completely extracted from frequent itemsets. Additionally, all periodic patterns also can be extracted because they are similar to closed swarms. Now, we will consider group patterns and we show that all of them can be directly extracted from the set of all frequent itemsets.

**Lemma 2.** Given  $FI = \{\mathcal{Y}_1, \mathcal{Y}_2, \dots, \mathcal{Y}_l\}$  contains all frequent itemsets mined from cluster matrix with  $minsup = \varepsilon$ . All group patterns  $(O, T_S)$  can be extracted from  $FI$ .

*Proof.*  $\forall (O, T_S)$  is a valid group pattern, we have  $\exists T_S = \{s_1, s_2, \dots, s_n\}$  and  $T_S$  is a set of disjointed convoys of  $O$ . Therefore,  $(O, T_{s_i})$  is a convoy and  $\forall s_i \in T_S, \forall t \in T_{s_i}, \exists c_t$  s.t.  $O \subseteq c_t$ . Let us assume  $C_{s_i}$  is a set of clusters corresponding to  $s_i$ , we know that  $\exists \mathcal{Y}, \mathcal{Y}$  is an itemset,  $\mathcal{Y} = \bigcup_{i=1}^n C_{s_i}$  and  $O(\mathcal{Y}) = \bigcap_{i=1}^n O(C_{s_i}) = O$ . Additionally,  $(O, T_S)$  is a valid group pattern; therefore,  $|O| \geq \varepsilon$  so  $|O(\mathcal{Y})| \geq \varepsilon$ . Consequently,  $\mathcal{Y}$  is a frequent itemset and  $\mathcal{Y} \in FI$  because  $\mathcal{Y}$  is an itemset and  $\sigma(\mathcal{Y}) = |O(\mathcal{Y})| \geq \varepsilon$ . Consequently,  $\forall (O, T_S), \exists \mathcal{Y} \in FI$  s.t.  $(O, T_S)$  can be extracted from  $\mathcal{Y}$  and therefore all group patterns can be extracted from  $FI$ .

## 4 FCI-based Spatio-Temporal Pattern Mining Algorithm

In this section, we propose two approaches i.e., *GeT\_Move* and *Incremental GeT\_Move*, to efficiently extract patterns. The global process is illustrated in Figure 6.

In the first step, a clustering approach is applied at each timestamp to group objects into different clusters. For each timestamp  $t_a$ , we thus have a set of clusters  $C_a = \{c_{1t_a}, c_{2t_a}, \dots, c_{mt_a}\}$ , with  $1 \leq k \leq m, c_{kt_a} \subseteq O_{DB}$ . Spatio-temporal data can thus be converted to a cluster matrix  $CM$ .

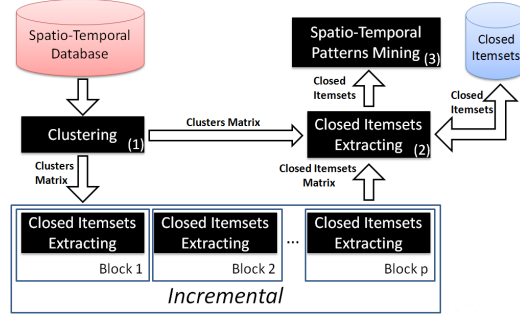


Fig. 6. The main process.

#### 4.1 GeT\_Move

After generating the cluster matrix  $CM$ , a FCI mining algorithm is applied on  $CM$  to extract all the FCIs. By scanning them and checking properties, we can obtain the patterns.

In this paper, we apply the LCM algorithm [19] to extract FCIs as it is known to be a very efficient algorithm. In LCM algorithm's process, we discard some useless candidate itemsets. In spatio-temporal patterns, items (resp. clusters) must belong to different timestamps and therefore items (resp. clusters) which form a FCI must be in different timestamps. In contrast, we are not able to extract patterns by combining items in the same timestamp. Consequently, FCIs which include more than 1 item in the same timestamp will be discarded.

Thanks to the above characteristic, we now have the maximum length of the FCIs which is the number of timestamps  $|T_{DB}|$ . Additionally, the LCM search space only depends on the number of objects (transactions)  $|O_{DB}|$  and the maximum length of itemsets  $|T_{DB}|$ . Consequently, by using LCM and by applying the above characteristic, *GeT\_Move* is not affected by the number of clusters and therefore the computing time can be greatly reduced.

The pseudo code of *GeT\_Move* is described in *Algorithm 1*. The core of *GeT\_Move* algorithm is based on the LCM algorithm which has been slightly modified by adding the pruning rule and by extracting patterns from FCIs. The initial value of FCI  $X$  is empty and then we start by putting item  $i$  into  $X$  (lines 2-3). By adding  $i$  into  $X$ , we have  $X[i]$  and if  $X[i]$  is a FCI then  $X[i]$  is used as a generator of a new FCI, call  $LCM\_Iter(X, \mathcal{T}(X), i(X))$  (lines 4-5). In  $LCM\_Iter$ , we first check properties of Section 3 (line 8) for FCI  $X$ . Next, for each transaction  $t \in \mathcal{T}(X)$ , we add all items  $j$ , which are larger than  $i(X)$  and satisfy the pruning rule, into occurrence sets  $\mathcal{J}[j]$  (lines 9-11). Next, for each  $j \in \mathcal{J}[j]$ , we check to see if  $\mathcal{J}[j]$  is a FCI, and if so, then we recall  $LCM\_Iter$  with the new generator (lines 12-14). Regarding to the  $PatternMining$  sub-function (lines 16-37), the algorithm basically checks properties of the itemset  $X$  to extract spatio-temporal patterns.

**Algorithm 1: GeT Move**


---

**Input** : Occurrence sets  $\mathcal{J}$ , int  $\varepsilon$ , int  $min_t$ , set of items  $C_{DB}$ , double  $\theta$ , int  $min_c$ , double  $min_{wei}$

```

1 begin
2    $X := I(\mathcal{T}(\emptyset));$  //The root
3   for  $i := 1$  to  $|C_{DB}|$  do
4     if  $|\mathcal{T}(X[i])| \geq \varepsilon$  and  $X[i]$  is closed then
5       LCM.Iter( $X[i], \mathcal{T}(X[i]), i$ );
6 LCM.Iter( $X, \mathcal{T}(X), i(X)$ )
7 begin
8   PatternMining( $X, min_t$ ); /* $X$  is a pattern?*/
9   foreach transaction  $t \in \mathcal{T}(X)$  do
10    foreach  $j \in t, j > i(X), j.time \notin time(X)$  do
11      insert  $j$  to  $\mathcal{J}[j]$ ;
12    foreach  $j \in \mathcal{J}[j]$  in the decreasing order do
13      if  $|\mathcal{T}(\mathcal{J}[j])| \geq \varepsilon$  and  $\mathcal{J}[j]$  is closed then
14        LCM.Iter( $\mathcal{J}[j], \mathcal{T}(\mathcal{J}[j]), j$ );
15      Delete  $\mathcal{J}[j]$ ;
16 PatternMining( $X, min_t$ )
17 begin
18   if  $|X| \geq min_t$  then
19     output  $X$ ; /*Closed Swarm*/
20      $gPattern := \emptyset; convoy := \emptyset; mc := \emptyset;$ 
21     for  $k := 1$  to  $|X| - 1$  do
22       if  $x_k.time = x_{(k+1)}.time - 1$  then
23          $convoy := convoy \cup x_k;$ 
24         if  $\frac{|\mathcal{T}(x_k) \cap \mathcal{T}(x_{k+1})|}{|\mathcal{T}(x_k) \cup \mathcal{T}(x_{k+1})|} \geq \theta$  then
25            $mc := mc \cup x_k;$ 
26         else
27           if  $|mc \cup x_k| \geq min_t$  then
28             output  $mc \cup x_k$ ; /*MovingCluster*/
29              $mc := \emptyset;$ 
30         else
31           if  $|convoy \cup x_k| \geq min_t$  and  $|\mathcal{T}(convoy \cup x_k)| = |\mathcal{T}(X)|$  then
32             output  $convoy \cup x_k$ ; /*Convoy*/
33              $gPattern := gPattern \cup (convoy \cup x_k);$ 
34           if  $|mc \cup x_k| \geq min_t$  then
35             output  $mc \cup x_k$ ; /*MovingCluster*/
36              $convoy := \emptyset; mc := \emptyset;$ 
37       if  $|gPattern| \geq min_c$  and  $size(gPattern) \geq min_{wei}$  then
38         output  $gPattern$ ; /*Group Pattern*/

```

39 Where:  $X$  is itemset,  $X[i] := X \cup i$ ,  $i(X)$  is the last item of  $X$ ,  $\mathcal{T}(X)$  is list of transactions that  $X$  belongs to,  $\mathcal{J}[j] := \mathcal{T}(X[j])$ ,  $j.time$  is time index of item  $j$ ,  $time(X)$  is a set of time indexes of  $X$ ,  $|\mathcal{T}(convoy)|$  is the number of transactions that the  $convoy$  belongs to,  $|gPattern|$  and  $size(gPattern)$  respectively are the number of convoys and the total length of the convoys in  $gPattern$ .

---

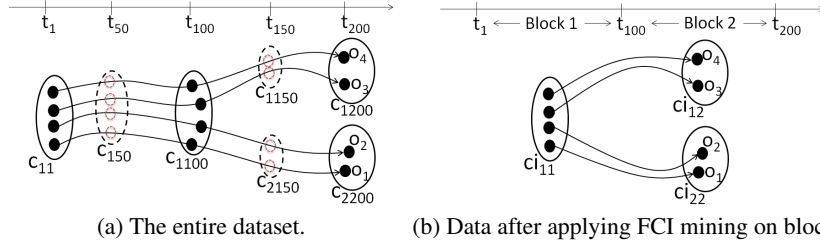


Fig. 7. A case study example. (b)- $ci_{11}$ ,  $ci_{12}$ ,  $ci_{22}$  are FCIs extracted from block 1 and block 2.

Table 3. Closed Itemset Matrix

Block $B$		$b_1$	$b_2$
Frequent Closed Itemsets $CI$		$ci_{11}$	$ci_{12}$ $ci_{22}$
$O_{DB}$	$o_1$	1	1
	$o_2$	1	1
	$o_3$	1	1
	$o_4$	1	1

## 4.2 Incremental GeT\_Move

Naturally, in real world applications (cars, animal migration), the objects tend to move together in short interval meanwhile their movements can be different in long interval. Therefore, the number of items (clusters) can be large and the length of FCIs can be long. Additionally, refer to [13, 19], the FCI mining algorithms search space are affected by the number of items and the length of itemsets. For instance, see Figure 7a, objects  $\{o_1, o_2, o_3, o_4\}$  move together during first 100 timestamps and after that  $o_1, o_2$  stay together while  $o_3, o_4$  move together in another direction. The problem here is that if we apply *GeT\_Move* on the whole dataset, the extraction of the itemsets can be very time consuming.

To deal with the issue, we propose the *Incremental GeT\_Move* algorithm. The main idea is to split the trajectories (resp. cluster matrix CM) into short intervals, called blocks. By applying FCI mining on each short interval, the data can then be compressed into local FCIs. Additionally, the length of itemsets and the number of items can be greatly reduced.

For instance, see Figure 7, if we consider  $[t_1, t_{100}]$  as a block and  $[t_{101}, t_{200}]$  as another block, the maximum length of itemsets in both blocks is 100 (instead of 200). Additionally, the original data can be greatly compressed (e.g. Figure 7b) and only 3 items remain:  $ci_{11}, ci_{12}, ci_{22}$ .

**Definition 6** *Block*. Given a set of timestamps  $T_{DB} = \{t_1, t_2, \dots, t_n\}$ , a cluster matrix  $CM$ .  $CM$  is vertically split into equivalent (in terms of intervals) smaller cluster matrices and each of them is a block  $b$ . Assume  $T_b$  is a set of timestamps of block  $b$ ,  $T_b = \{t_1, t_2, \dots, t_k\}$ , thus we have  $|T_b| = k \leq |T_{DB}|$ .

Assume that we obtain a set of blocks  $B = \{b_1, b_2, \dots, b_p\}$  with  $|T_{b_1}| = |T_{b_2}| = \dots = |T_{b_p}|$ ,  $\bigcup_{i=1}^p b_i = CM$  and  $\bigcap_{i=1}^p b_i = \emptyset$ . Given a set of FCI collections  $CI = \{CI_1, CI_2, \dots, CI_p\}$  where  $CI_i$  is mined from block  $b_i$ .  $CI$  is presented as a *closed itemset matrix* which is formed by horizontally connecting all local FCIs:  $CIM = \bigcup_{i=1}^p CI_i$ .

**Definition 7** *Closed Itemset Matrix (CIM)*. Closed itemset matrix is a cluster matrix with some differences as follows: 1) Timestamp  $t$  now becomes a block  $b$ . 2) Item  $c$  is a FCI  $ci$ .

For instance, see Table 3, we have two sets of FCIs  $CI_1 = \{ci_{11}\}$ ,  $CI_2 = \{ci_{12}, ci_{22}\}$  which are respectively extracted from blocks  $b_1, b_2$ . We have  $CIM$  which is created from  $CI_1, CI_2$  in Table 3.

Now, by applying FCI mining on closed itemset matrix  $CIM$ , we retrieve all FCIs from corresponding data. Note that items (in  $CIM$ ) which are in the same block cannot be in the same FCIs.

**Lemma 3.** Given a cluster matrix  $CM$  which is vertically split into a set of blocks  $B = \{b_1, b_2, \dots, b_p\}$  so that  $\forall \mathcal{Y}, \mathcal{Y}$  is a FCI and  $\mathcal{Y}$  is extracted from  $CM$  then  $\mathcal{Y}$  can be extracted from the closed itemset matrix  $CIM$ .

*Proof.* Let us assume that  $\forall b_i, \exists I_i$  is a set of items belonging to  $b_i$  and therefore we have  $\bigcap_{i=1}^{|B|} I_i = \emptyset$ . If  $\forall \mathcal{Y}, \mathcal{Y}$  is a FCI extracted from  $CM$  then  $\mathcal{Y}$  is formed as  $\mathcal{Y} = \{\gamma_1, \gamma_2, \dots, \gamma_p\}$  where  $\gamma_i$  is a set of items s.t.  $\gamma_i \subseteq I_i$ . Additionally,  $\mathcal{Y}$  is a FCI and  $O(\mathcal{Y}) = \bigcap_{i=1}^p O(\gamma_i)$  then  $\forall O(\gamma_i), O(\mathcal{Y}) \subseteq O(\gamma_i)$ . Furthermore, we have  $|O(\mathcal{Y})| \geq \varepsilon$ ; therefore,  $|O(\gamma_i)| \geq \varepsilon$  so  $\gamma_i$  is a frequent itemset. Assume that  $\exists \gamma_i, \gamma_i \notin CI_i$  then  $\exists \Psi, \Psi \in CI_i$  s.t.  $\gamma_i \subseteq \Psi$  and  $\sigma(\gamma_i) = \sigma(\Psi), O(\gamma_i) = O(\Psi)$ . Note that  $\Psi, \gamma_i$  are from  $b_i$ . Remember that  $O(\mathcal{Y}) = O(\gamma_1) \cap O(\gamma_2) \cap \dots \cap O(\gamma_i) \cap \dots \cap O(\gamma_p)$  and we have:  $\exists \mathcal{Y}'$  s.t.  $O(\mathcal{Y}') = O(\gamma_1) \cap O(\gamma_2) \cap \dots \cap O(\Psi) \cap \dots \cap O(\gamma_p)$ . Therefore,  $O(\mathcal{Y}') = O(\mathcal{Y})$  and  $\sigma(\mathcal{Y}') = \sigma(\mathcal{Y})$ . Additionally, we know that  $\gamma_i \subseteq \Psi$  so  $\mathcal{Y} \subseteq \mathcal{Y}'$ . Consequently, we obtain  $\mathcal{Y} \subseteq \mathcal{Y}'$  and  $\sigma(\mathcal{Y}) = \sigma(\mathcal{Y}')$ . Therefore,  $\mathcal{Y}$  is not a FCI. That violates the assumption and therefore we have: if  $\exists \gamma_i, \gamma_i \notin CI_i$  therefore  $\mathcal{Y}$  is not a FCI. Finally, we can conclude that  $\forall \mathcal{Y}, \mathcal{Y} = \{\gamma_1, \gamma_2, \dots, \gamma_p\}$  is a FCI extracted from  $CM$ ,  $\forall \gamma_i \in \mathcal{Y}$ ,  $\gamma_i$  must be belong to  $CI_i$  and  $\gamma_i$  is an item in closed itemset matrix  $CIM$ . Therefore,  $\mathcal{Y}$  can be retrieved by applying FCI mining on  $CIM$ .

By applying *Lemma 3*, we can obtain all the FCIs and from the itemsets, patterns can be extracted. Note that the Incremental GeT\_Move does not depend on the length restriction  $min_t$ . The reason is that  $min_t$  is only used in Spatio-Temporal Patterns Mining step. Whatever  $min_t$  ( $min_t \geq$  block size or  $min_t \leq$  block size), Incremental GeT\_Move can extract all the FCIs and therefore the final results are the same.

The pseudo code of *Incremental GeT\_Move* is described in *Algorithm 2*. The main different between the code of *Incremental GeT\_Move* and *GeT\_Move* is the *Update* sub-function. In this function, we step by step generate the closed itemsets matrix from blocks (line 14 and lines 22-26). Next, we apply *GeT\_Move* to extract patterns (line 5).

## 5 Experimental Results

A comprehensive performance study has been conducted on real datasets and synthetic datasets. All the algorithms are implemented in C++, and all the experiments are carried out on a 2.8GHz Intel Core i7 system with 4GB Memory. The system runs Ubuntu 11.10 and g++ version 4.6.1.

**Algorithm 2: Incremental GeT\_Move**


---

**Input** : Occurrence sets  $K$ , int  $\varepsilon$ , int  $min_t$ , double  $\theta$ , set of Occurrence sets (blocks)  $B$ , int  $min_c$ , double  $min_{wei}$

```

1 begin
2    $K := \emptyset; CI := \phi; int\ item\_total := 0;$ 
3   foreach  $b \in B$  do
4     |  $LCM(b, \varepsilon, I_b);$ 
5     |  $GeT\_Move(K, \varepsilon, min_t, CI, \theta, min_c, min_{wei});$ 
6   LCM(Occurrence sets  $\mathcal{J}$ , int  $\sigma_0$ , set of items  $C$ )
7   begin
8      $X := I(\mathcal{T}(\emptyset));$  //The root
9     for  $i := 1$  to  $|C|$  do
10    | if  $|\mathcal{T}(X[i])| \geq \varepsilon$  and  $X[i]$  is closed then
11    | |  $LCM\_Iter(X[i], \mathcal{T}(X[i]), i);$ 
12  LCM.Iter( $X, \mathcal{T}(X), i(X)$ )
13  begin
14    Update( $K, X, \mathcal{T}(X), item\_total + +$ );
15    foreach transaction  $t \in \mathcal{T}(X)$  do
16    | foreach  $j \in t, j > i(X), j.time \notin time(X)$  do
17    | | insert  $j$  to  $\mathcal{J}[j]$ ;
18    | foreach  $j, \mathcal{J}[j] \neq \phi$  in the decreasing order do
19    | | if  $|\mathcal{T}(\mathcal{J}[j])| \geq \varepsilon$  and  $\mathcal{J}[j]$  is closed then
20    | | |  $LCM\_Iter(\mathcal{J}[j], \mathcal{T}(\mathcal{J}[j]), j);$ 
21    | | Delete  $\mathcal{J}[j]$ ;
22  Update( $K, X, \mathcal{T}(X), item\_total$ )
23  begin
24    foreach  $t \in \mathcal{T}(X)$  do
25    | insert  $item\_total$  into  $K[t]$ ;
26     $CI := CI \cup item\_total;$ 

```

---

The implementation (source code) is available and also integrated in our online demonstration system<sup>4</sup>. As in [6], we only report the results on the following datasets<sup>5</sup>: *Swainsoni dataset* includes 43 objects evolving over time and 764 different timestamps. The interested readers may refer to our online demonstration system<sup>2</sup> for other experimental results. Additionally, similar to [6,3,10], we first use linear interpolation to fill in the missing data and then DBScan [5] ( $MinPts = 2, Eps = 0.001$ ) is applied to generate clusters at each timestamp.

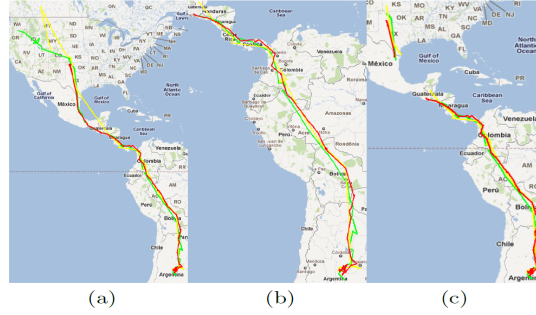
In the comparison, we employ *CMC*, *CuTS*<sup>6</sup> (convoy mining) and *ObjectGrowth* (closed swarm mining). Note that, in [6], *ObjectGrowth* outperforms *VG – Growth* [15] (a group patterns mining algorithm) in terms of performance and therefore we will only consider *ObjectGrowth* and not both. Note that, in the reported experiments,

<sup>4</sup> [www.lirmm.fr/~phan/index.jsp](http://www.lirmm.fr/~phan/index.jsp)

<sup>5</sup> <http://www.movebank.org>

<sup>6</sup> The source code of *CMC*, *CuTS*<sup>\*</sup> is available at [http://lsirpeople.epfl.ch/jeung/source\\_codes.htm](http://lsirpeople.epfl.ch/jeung/source_codes.htm)





**Fig. 8.** An example of patterns discovered from Swainsoni dataset. (a) One of discovered closed swarms, (b) One of discovered convoys, (c) One of discovered group patterns.

GeT\_Move and Incremental GeT\_Move extract closed swarms, convoys and group patterns while *CMC*, *CuTS\** only extract convoys and *ObjectGrowth* extract closed swarms.

### 5.1 Effectiveness

We proved that mining spatio-temporal patterns can be similarly mapped into itemsets mining issue. Therefore, in theoretical way, our approaches can provide the correct results. Experimentally, we do a further comparison, we first obtain the spatio-temporal patterns by applying *CMC*, *CuTS\**, *ObjectGrowth* as well as our approaches. To apply our algorithms, we split cluster matrix into blocks such as each block  $b$  contains 25 timestamps. Additionally, to retrieve all the spatio-temporal patterns, in the reported experiments, the default value of  $\varepsilon$  is set to 2 (two objects can form a pattern),  $min_t$  is 1. Note that the default values are the hardest conditions for examining the algorithms. Then in the following we mainly focus on different values of  $min_t$  in order to obtain different sets of convoys, closed swarms and group patterns. Note that for group patterns,  $min_c$  is 1 and  $min_{wei}$  is 0.

The results show that our proposed approaches obtain the same results compared to the traditional algorithms. An example of patterns is illustrated in Figure 8. For instance, see Figure 8a, a closed swarm is discovered within a FCI. Furthermore, from the itemset, a convoy and a group pattern are also extracted (i.e. Figure 8b, 8c).

### 5.2 Efficiency

To show the efficiency of our algorithms, we also generate larger synthetic datasets using Brinkhoff's network-based generator of moving objects<sup>7</sup> as in [6]. We generate 500 objects ( $|O_{DB}| = 500$ ) for  $10^4$  timestamps ( $|T_{DB}| = 10^4$ ) using the generator's default map with low moving speed. There are  $5 \times 10^6$  points in total. DBScan ( $MinPts = 3$ ,  $Eps = 300$ ) is applied to obtain clusters for each timestamp.

**Efficiency w.r.t.  $\varepsilon$ ,  $min_t$ .** Figure 9a, 10a show running time w.r.t.  $\varepsilon$ . It is clear that our approaches outperform other algorithms. *ObjectGrowth* is the lowest one and the

<sup>7</sup> <http://iapg.jade.hs.de/personen/brinkhoff/generator/>

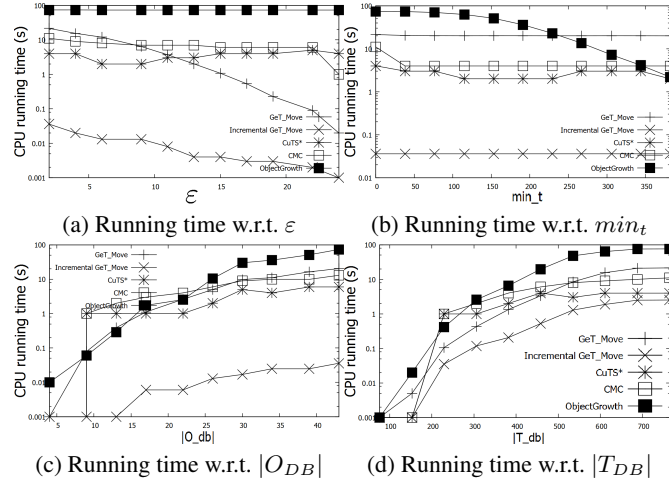


Fig. 9. Running time on Swainsoni Dataset.

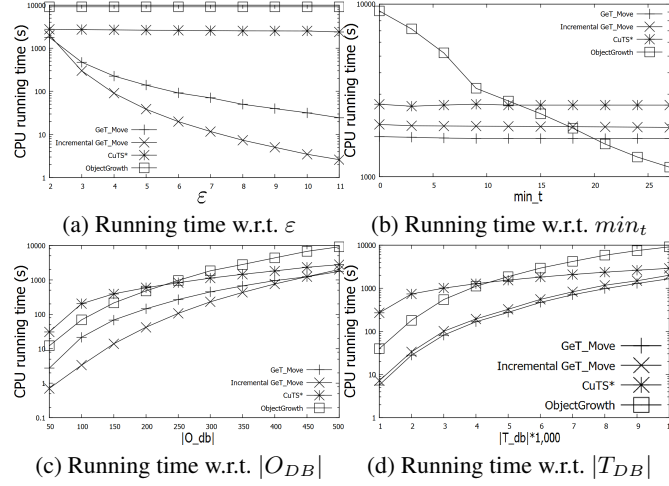


Fig. 10. Running time on Synthetic Dataset.

main reason is that with low  $min_t$  (default  $min_t = 1$ ), the *Apriori Pruning* rule (the most efficient pruning rule) is no longer effective. Therefore, the search space is greatly enlarged ( $2^{|O_{DB}|}$  in the worst case). Additionally, there is no pruning rule for  $\epsilon$  and therefore the change of  $\epsilon$  does not directly affect the running time of ObjectGrowth. A little bit further, Get\_Move is lower than Incremental Get\_Move. The main reason is that Get\_Move has to process with large number of items and long itemsets. While, thanks to blocks, the number of items is greatly reduced and itemsets are not long as the ones in Get\_Move.

Figure 9b, 10b show running time w.r.t.  $min_t$ . In almost all cases, our approaches outperform other algorithms. See Figure, 10b, with low  $min_t$ , our algorithm is much faster than the others. However, when  $min_t$  is higher ( $min_t > 20$  in Figure 10b) our

algorithms take more time than CuTS\* and ObjectGrowth. This is because with high value of  $min_t$ , the number of patterns is significantly reduced (Figure 11b, 12b) (i.e. no extracted convoy when  $min_t > 100$  (resp.  $min_t > 10$ ), Figure 11b (resp. Figure 12b)) and therefore CuTS\* and ObjectGrowth is faster. While, GeT\_Move and Incremental GeT\_Move have to work with FCIs.

**Efficiency w.r.t.  $|O_{DB}|, |T_{DB}|$ .** Figure 9c-d, Figure 10c-d show the running time when varying  $|O_{DB}|$  and  $|T_{DB}|$  respectively. In all figures, Incremental GeT\_Move outperforms other algorithms. However, with synthetic data (Figure 10d) and lowest values of  $\epsilon = 2$  and  $min_t = 1$ , GeT\_Move is a little bit faster than Incremental GeT\_Move.

**Scalability w.r.t.  $\epsilon$ .** We can notice that the running time of algorithms does not change significantly when varying  $min_t, |O_{DB}|, |T_{DB}|$  in synthetic data (Figures 10). However, they are quite different when varying  $\epsilon$  (default  $min_t = 1$ ). Therefore, we generate another large synthetic data to test the scalability of algorithms on  $\epsilon$ . The dataset includes 50,000 objects moving during 10,000 timestamps and it contains 500 million locations in total. The executions of CMC and CuTS\* stop due to a lack of memory capacity after processing 300 million locations. Additionally, ObjectGrowth can not provide the results after 1day running. The main reason is that with low  $min_t$  ( $= 1$ ), the search space is significant larger ( $\approx 2^{50,000}$ ). While, thanks to the LCM approach, our algorithms can provide the results within hours (Figure 13a).

**Efficiency w.r.t. Block-size.** To investigate the optimal value of block-size, we examine Incremental GeT\_Move by using the default values of  $\epsilon, min_t$  with different block-size values on real datasets (note: Buffalo<sup>3</sup>,  $|O_{DB}| = 165, |T_{DB}| = 3,000$ ) and synthetic dataset ( $|O_{DB}| = 500, |T_{DB}| = 1,000$ ). The optimal block -size range can be from 20 to 30 timestamps within which Incremental GeT\_Move obtains the best performance for all the datasets (Figure 13b). The main reason is that objects tend to move together in suitable short interval (from 20 to 30 timestamps). Therefore, by setting block-size in this range, the data is efficiently compressed into FCIs. Meanwhile, with larger block-size values, the objects' movements are quite different; therefore, the data compressing is not so efficient. Regarding to small block-size values (5-15), we have to face up to a large number of blocks so that the process is slowed down. In the previous experiments, block-size is set to 25.

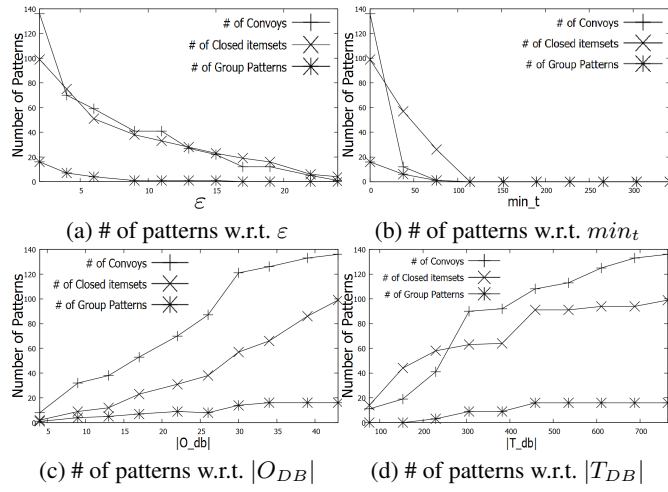
## 6 Conclusion and Discussion

In this paper, we propose unifying incremental approaches to automatically extract different kinds of spatio-temporal patterns by applying FCI mining techniques. Their effectiveness and efficiency have been evaluated by using real and synthetic datasets. Experiments show that our approaches outperform traditional ones.

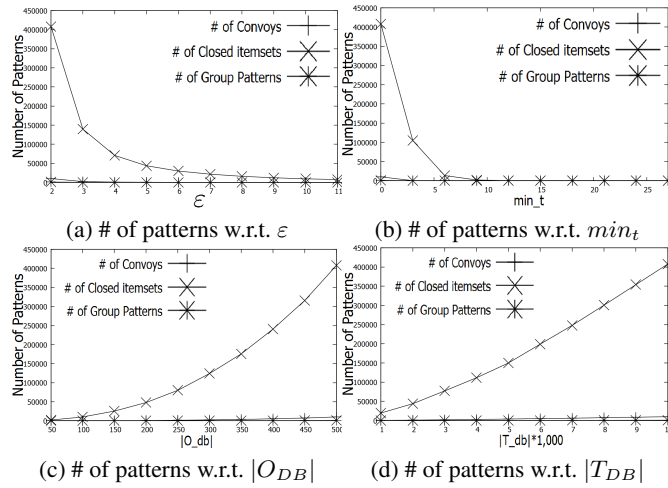
One next issue we plan to address is how to take into account the arrival of new objects which were not available for the first extraction. Now, we can store the result to improve the process when new object movements arrive. But, in this approach, we take the hypothesis is that the number of objects remains the same. However in some applications these objects could be different.

## References

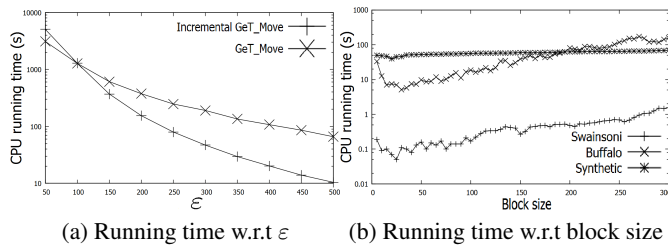
1. J. Gudmundsson, M. van Kreveld. *Computing longest duration flocks in trajectory data*. In: GIS 06, New York, NY, USA, pp.35-42.
2. MR. Vieira, P. Bakalov, VJ. Tsotras. *On-line Discovery of Flock Patterns in Spatio-Temporal Data*. In: GIS 09, New York, USA, pp.286-295.
3. H. Jeung, ML. Yiu, X. Zhou, CS. Jensen, HT. Shen. *Discovery of Convoys in Trajectory Databases*. PVLDB 2008, 1(1):1068-1080.
4. P. Kalnis, N. Mamoulis, S. Bakiras. *On Discovering Moving Clusters in Spatio-temporal Data*. In SSTD 2005, Angra dos Reis, Brazil, pages 364-381.
5. M. Ester, H.-P. Kriegel, J. Sander, X. Xu. *A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise*. KDD '96, Portland, pp. 226-231.
6. Z. Li, B. Ding, J. Han, R. Kays. *Swarm: Mining Relaxed Temporal Moving Object Clusters*. VLDB2010, Singapore, pp. 723-734.
7. A.O.C. Romero. *Mining moving flock patterns in large spatio-temporal datasets using a frequent pattern mining approach*. Master Thesis, University of Twente, faculty ITC, March 2011.
8. R. Agrawal and R. Srikant. *Fast algorithms for mining association rules*. VLDB'94, pp.487-499.
9. Z. Li, M. Ji, J.-G. Lee, L. Tang, Y. Yu, J. Han, and R. Kays. *Movemine: Mining moving object databases*. In SIGMOD 2010, Indianapolis, Indiana, pp.1203-1206.
10. H. Jeung, X. Zhou, H. T. Shen. *Convoy Queries in Spatio-Temporal Databases*. ICDE 2008, Cancun, Mexico, pp.1457-1459.
11. F. Verhein. *Mining Complex Spatio-Temporal Sequence Patterns*. SDM'09, Nevada, pp.605-616.
12. C.S. Jensen, D. Lin, and B.C. Ooi. *Continuous clustering of moving objects*. In KDE(2007), pp. 1161-1174. issn: 1041-4347.
13. C. Lucchese, S. Orlando, R. Perego. *DCI Closed: A Fast and Memory Efficient Algorithm to Mine Frequent Closed Itemsets*. ICDM FIMI 2004.
14. V. Bogorny and S. Shekhar. *Spatial and Spatio-Temporal Data Mining*. Tutorial on Spatial and Spatio-Temporal Data Mining, ICDM2010, Australia.
15. Y. Wang, E.-P. Lim, and S.-Y. Hwang. *Efficient Mining of Group Patterns from User Movement Data*. In DKE(2006), pp. 240-282.
16. H. Cao, N. Mamoulis, D.W. Cheung. *Discovery of Collocation Episodes in Spatiotemporal Data*. In ICDM'06, Hong Kong, pp.823-827, ISBN: 0-7695-2701-9.
17. J.-g. Lee, J. Han. *Trajectory Clustering: A Partition and Group Framework*. In SIGMOD 2007, pp. 593-604.
18. N. Mamoulis, H. Cao, G. Kollios, M. Hadjieleftheriou, Y. Tao, D. W. Cheung. *Mining, Indexing, and Querying Historical Spatiotemporal Data*. SIGKDD'04, pp.236-245.
19. T. Uno, M. Kiyomi, and H. Arimura. *LCM ver. 2: Efficient mining algorithms for frequent/closed/maximal itemsets*. ICDM FIMI 2004.
20. J. Han, H. Pei, Y. Yin. *Mining Frequent Patterns without Candidate Generation*. In SIGMOD'00, New York, USA.
21. J. Han, Z. Li, L. A. Tang. *Mining Moving Object, Trajectory and Traffic Data*. In DAS-FAA'10, Japan.



**Fig. 11.** # of patterns on Swainsoni Dataset. Note that # of FCIs is equal to # of closed swarms.



**Fig. 12.** # of patterns on Synthetic Dataset. Note that # of FCIs is equal to # of closed swarms.



**Fig. 13.** (a) Running time w.r.t  $\epsilon$  on large Synthetic Dataset, (b) Running time w.r.t block size.