# A Trajectory K-Anonymity Model Based on Point Density and Partition

Wanshu Yu*
University of Glasgow
Glasgow, UK
South China University of
Technology, School of Computer
Science and Engineering
Guangzhou, China
2817055y@student.gla.ac.uk

Haonan Shi*
Case Western Reserve University
Cleveland, Ohio, USA
South China University of
Technology, School of Computer
Science and Engineering
Guangzhou, China
hxs896@case.edu

Hongyun Xu†
South China University of
Technology, School of Computer
Science and Engineering
Guangzhou, China
hongyun@scut.edu.cn

## ABSTRACT

As people's daily life becomes increasingly inseparable from various mobile electronic devices, relevant service application platforms and network operators can collect numerous individual information easily. When releasing these data for scientific research or commercial purposes, users' privacy will be in danger, especially in the publication of spatiotemporal trajectory datasets. Therefore, to avoid the leakage of users' privacy, it is necessary to anonymize the data before they are released. However, more than simply removing the unique identifiers of individuals is needed to protect the trajectory privacy, because some attackers may infer the identity of users by the connection with other databases. Much work has been devoted to merging multiple trajectories to avoid re-identification, but these solutions always require sacrificing data quality to achieve the anonymity requirement. In order to provide sufficient privacy protection for users' trajectory datasets, this paper develops a study on trajectory privacy against re-identification attacks, proposing a trajectory K-anonymity model based on Point Density and Partition (KPDP). Our approach improves the existing trajectory generalization anonymization techniques regarding trajectory set partition preprocessing and trajectory clustering algorithms. It successfully resists re-identification attacks and reduces the data utility loss of the k-anonymized dataset. A series of experiments on a real-world dataset show that the proposed model has significant advantages in terms of higher data utility and shorter algorithm execution time than other existing techniques.

## KEYWORDS

Trajectory dataset; Privacy protection; Re-identification attack; Trajectory clustering

## 1 INTRODUCTION

With the rapid development of mobile devices and communication technologies, location-based information online service platforms are widely used, which are highly relevant to people's daily lives and bring convenience. For example, after a user opens a navigation software, the application automatically sends location-based queries to the server, pulling map query results regarding the current area, such as nearby restaurants, car parks, shopping centres,

and banks. When a user accesses such a Location-based Information Service (LBS) application, the network operator of the mobile device can extensively record data about their movement trajectory [24], i.e. the sequence of location coordinates that the user passed over some time. Releasing the collected information to the public not only facilitates the research work of scientific organizations but also plays a vital role in the transparency of authorities such as operators and governments. However, the publication of data can be exploited by malicious attackers, resulting in the disclosure of user privacy.

Due to the rapid development of high-capacity storage and data analysis techniques, it is possible for attackers to distinguish the trajectory travelled by an individual from publicly released trajectory datasets and to obtain more sensitive privacy information by integrating their location and route with other databases [31]. Hence it is generally the correspondence between a user's spatiotemporal trajectory and the individual identity. However, it is not sufficient to simply erase the direct unique identifiers from the database to resist attacks. This is because once an attacker combines the quasi-identifiers with known background knowledge, it is possible to deduce the correspondence, thus causing danger to user privacy, property, security and reputation. Therefore, how to scientifically encrypt datasets has become an important issue in data release and privacy protection today.

In order to guarantee the privacy of users despite the public release of trajectory data, it is necessary to employ various techniques to process the trajectory data before releasing it. Many scholars have worked on the issue of trajectory privacy attacks and protection, proposing various techniques to achieve privacy protection in LBS, such as generalization, obfuscation and fuzzing. Nevertheless, although these existing techniques can protect user privacy from being attacked or exposed under certain circumstances, the corresponding algorithms are usually of high time and space complexity. Moreover, due to the specificity of trajectory shape distribution and the sensitivity of location information, the privacy protection processing will lead to the loss of information to a large extent, thus reducing the utility of the data.

To address the problem above, we propose a privacy protection methodology for user trajectories adopting machine learning techniques to prevent revealing the private information of LBS users on the one hand and to retain the features and accuracy of the original trajectories as far as possible on the other hand, so as to reduce the loss of information after data processing. Specifically, it is required

---

*Both authors contributed equally to this research
†Corresponding author

that trajectories from different users in the released dataset are indistinguishable from each other. As a result, trajectories in the original dataset typically need to be replaced with the generalized trajectory for several users. The process of thus replacing a specific value with a more general and imprecise value is called generalization [1]. The higher the level of generalization, the higher the extent of privacy protection, but the lower the data utility of the published trajectories and the higher the loss of information after generalization. In order to balance the degree of privacy protection and the generalization information loss, we preprocess trajectories by segmenting them according to the point density and generalize them based on the idea of DBSCAN cluster algorithm [11], achieving the resistance of the released dataset to re-identification attacks and preserving the distribution features of the trajectories in the best manner possible. To the best of our knowledge, our paper proposes such a partition preprocessing mechanism for the first time. Our main contributions are summarised as follows.

- We investigate the shortcomings of existing trajectory privacy-preserving algorithms and propose a trajectory K-anonymity model based on Point Density and Partition (KPDP). The deficiencies of the existing models mainly stem from the irregularity of the shape distribution of real trajectories and the specific data structure, making it difficult to measure the similarity between trajectories, and thus unable to accurately cluster and generalize trajectories, resulting in a high information loss in the released dataset relative to the original dataset. Based on this situation, KPDP can segment trajectories based on point density before clustering them so that the length of trajectories is relatively balanced and the spatial distribution characteristics of the original trajectories are retained, yielding a lower generalization information loss than other models.
- To further enhance the utility of anonymized trajectory datasets and to achieve k-anonymity, this paper proposes an adaptive DBSCAN trajectory clustering algorithm. The algorithm measures the distance between trajectories using the loss from the alignment of trajectories and then clusters them based on sample density. However, due to the uncertainty of the number of samples in the clusters and the possible presence of noise from DBSCAN, direct adoption of its idea cannot guarantee k-anonymity. We consequently developed an adaptive DBSCAN trajectory clustering algorithm that can automatically adjust the values of parameters based on the number of trajectories and noise in each cluster and repeatedly call the core module to cluster. The main advantage of DBSCAN over other unsupervised machine learning-based algorithms is that it is not constrained by given values of parameters and can produce clustering results that better reflect the characteristics of the trajectory distribution, thus improving the data utility of the released dataset.
- We conducted extensive experiments based on a realistic trajectory dataset to evaluate the privacy-preserving effects of segmentation preprocessing mechanisms and trajectory clustering algorithms under different privacy metrics. The experiment results show that our approach performs better in terms of information loss and running time compared to other existing approaches.

The subsequent structure of this paper is organized as follows. Section 2 introduces and defines trajectory privacy attacks, privacy anonymity criteria, privacy-preserving methods, generalization hierarchy models, and trajectory alignment techniques. We then show an overview of KPDP in Section 3. Following this framework, we illustrate the rationale of the segmentation preprocessing mechanism and the design of the anonymization model in Section 4 and 5. The experiment results and evaluation are presented in Section 6. Finally, we conclude with an overview of our contributions in Section 7.

## 2 BACKGROUND AND RELATED WORKS

### 2.1 Attack Model

A trajectory privacy attack is the acquisition of a user's private information from a trajectory dataset by an attacker with background knowledge. In general, most studies assume that the background knowledge known to the attacker is part of the spatiotemporal points on the user's trajectory, and the privacy information the attacker attempts to disclose is the complete trajectory data of that victim. For a given anonymized dataset, Zhen Tu et al. [39] denote the set of users as $U = U_i$ and the corresponding set of trajectories as $T = T_i$, where $T_i$ denotes the spatiotemporal points of the trajectory of user $U_i$. A constant number of partial points sampled from the actual trajectories is considered the attacker's external background knowledge, denoted as $E = E_i$, where $E_i$ denotes the attacker's external observation of the user $U_i$. With any external information $E_i$, an attacker makes a successful re-identification attack if he can match only one trajectory, whose formulation is shown in Eq. (1).

$$C_i = \begin{cases} 1 & \left|T_j \mid T_j \cap E_j, T_j \in T\right| = 1, \\ 0 & otherwise \end{cases} \quad \sum_i C_i \geq 1 \quad (1)$$

where $C_i$ denotes whether the user $U_i$ is re-identifiable and $|*|$ denotes the size of the set $*$.

In addition, adversaries can also launch attacks based on more public information. Zhen Tu et al. [38] stated that an attacker could infer a victim's motivation and behaviour to visit a location by associating the Point of Interest (PoI) that the user passes on a map with the primary function of its corresponding location. Huaxin Li et al. [23] matched the locations shared by users on social networks with their real travel trajectories to enable external attackers to infer information such as their age, gender, and education. John Krumm [18] quantified the effectiveness of using different attack algorithms to recognize the location of subjects' homes and then identify them through a programmable web search engine. According to [40], the types of location privacy attacks explicitly include Single position attack [27], Multiple position attack [2, 13, 36] and Context linking attack [15, 25, 33]. Although there are many approaches to attacking user privacy, re-identification attack remains the most fundamental problem. This paper focuses on studying resistance to privacy issues caused by re-identification attack.

## 2.2 Privacy Model

The protection of individuals from re-identification attacks has been a topic of much discussion in recent years. The k-anonymity criterion is the most commonly used privacy-preserving metric to resist re-identification attacks for data publishing within the privacy and anonymity domain. K-anonymity is a concept introduced by Samarati and Sweeney in 1998. K-anonymity requires that each record stored in a published dataset should be indistinguishable from at least $k-1$ other records [29, 30, 35], i.e. it requires that the same quasi-identifier refers to at least multiple records, making it impossible for adversaries to connect records with other databases by quasi-identifiers and thus deduce user identity and more private information. Current k-anonymity implementations are mostly used to protect data anonymity for category and numerical attributes in general relational databases, including Generalization and suppression, Incognito, Top-down specialization, Clustering, and Multidimensional partitioning [12, 21, 22]. However, for such irregular geometric data structure as trajectory, it requires a specific processing method to achieve k-anonymity [3, 15, 32, 41].

In this paper, we need to ensure at least $k$ distinct trajectories in each cluster obtained from the original trajectory set and generalize them to identical anonymous records to form a trajectory dataset that conforms to k-anonymity.

## 2.3 Defense Techniques

Among diverse researches to achieve k-anonymity of trajectory data, generalization is one of the most dominant approaches. According to the different details of generalization techniques, such as the encoding and operation of the Domain Generalization Hierarchy (DGH) tree, there are three main types of generalization: full domain generalization, subtree generalization and cell level generalization [42]. Acar Tamersoy et al. [37] proposed a heuristic approach based on the concept of generalization to achieve k-anonymity. Sina Shaham et al. [31] used a heuristic and a variant k-means algorithm for trajectory clustering and anonymization. Marco Gramaglia et al. [14] used a k-normalization algorithm to address the efficiency problem of generalization during the anonymization of trajectory datasets.

In addition to generalization methods, many researchers have worked on resisting trajectory privacy attacks from multiple other perspectives. The authors of [8, 13, 28] provide special treatments for sensitive locations on maps to protect the semantic privacy of trajectories. Zhen Tu et al. [39] protect trajectories from re-identification and semantic attacks based on k-anonymity, l-diversity and t-confidentiality. There are also researches which generate stopping points and noise points to obfuscate the original trajectory set, demonstrating the effectiveness of resistance to virtual location-information [6, 10, 17]. Jiaxin Ding [9] prevents an attacker from identifying a specific user's trajectory by exchanging the user's ID at the intersection of the trajectory. Jae-Gil Lee et al. [20] sliced the trajectory into line segments and clustered the new set of segments based on a definition of the distance between the segments. A middleware structure and an algorithm for adjusting the resolution of location information along the spatial or temporal dimension was introduced by [15], which satisfies a specified anonymity constraint within a given region.

This paper proposes a trajectory k-anonymity approach to preserve privacy via generalization techniques with low loss of data utility and algorithm time complexity.

## 2.4 Generalization for Secure Data

Generalization techniques enable the goal of ensuring the privacy of published dataset without compromising data availability. Generalization and suppression[35] are used to provide privacy to individuals. Pre-defined generalization hierarchy[16] allows the construction of generalization hierarchies before data masking. Full domain generalization hierarchy[29] enables the mapping of attributes to a more general domain in the domain generalization hierarchy.

A DGH tree is a quantitative model of information loss for generalizing numerical or categorical attributes [17]. In the anonymity domain, the structure of DGH trees has several different ways of formation. The structure of a DGH tree for numerical attributes can either be predefined by the user based on the usage scenario [16] or built dynamically during the generalization process [4]. Category attributes often require the manual creation of hierarchical structures considering attribute characteristics and usage scenarios [34]. Due to the complexity of the practical situation, not all anonymization processes for category attributes apply to the DGH tree model. For the generalization of the trajectory dataset in this paper, the latitude and longitude of the trajectories are usually used to dynamically construct the DGH tree and generalize the trajectory set based on this model while calculating the caused information loss.

## 2.5 Trajectory Alignment

Dynamic Sequence Alignment (DSA) is a trajectory alignment algorithm derived from the sequence alignment method of proteins and DNA in biology [5, 19]. The algorithm uses a dynamic programming approach to obtain the loss matrix for the alignment of two trajectories recursively, and then from the backtracking of this loss matrix, find the strategy that can make the merging of these two trajectories produce the least information loss and subsequently obtain the merged trajectory and the minimum information loss value.

Based on this, the Progressive Sequence Alignment (PSA) algorithm is derived as a multi-sequence alignment algorithm capable of aligning a cluster of trajectories [7]. The PSA algorithm can sequentially perform DSA operations on the trajectories within a cluster to obtain a synthetic trajectory obtained by aligning all trajectories within the cluster, and the PSA algorithm is often used to generalize the clusters of trajectories formed by clustering because it first selects the longest trajectory from a group of trajectories as the base trajectory, and then sequentially selects the remaining trajectories within the group in order of trajectory length from longest to shortest to align and synthesize with the base trajectory based on DSA. The generalized trajectory generated after DSA process will become the new base trajectory for the subsequent DSA process until all trajectories in the group have been aligned with the base trajectory.

# 3 SYSTEM OVERVIEW

## 3.1 System Utility Measurement

In the KPDP framework, trajectory alignment is the key to performing trajectory anonymization, and information loss is incurred in trajectory alignment. In order to calculate the loss of KPDP in the process of anonymizing trajectories more accurately and efficiently, a new DGH tree is proposed in this paper. This DGH tree is a partially ordered tree structure, which is able to map the specific and generalized values of attribute $A$ for a certain attribute $A$. The root node of the DGH tree indicates the case with the highest degree of generalization. Our DGH tree is constructed by dividing a number of small intervals of equal length within the range of corresponding values taken and then using these small intervals as leaf nodes to construct a full binary tree. If the number of leaf nodes is not enough to fill the bottom level of the binary tree, some invalid points are added to fill it up. A simple illustration of a DGH tree with a 4-layer structure is shown in Figure 1. The leaf nodes numbered 12 can be generalized to the parent node 6 or the ancestor nodes 3 and 1. Specifically, the DGH trees of KPDP in this paper are two DGH trees formed by building latitude and longitude in the trajectory set, corresponding to the x-axis and y-axis coordinate systems on the map plane space, respectively.
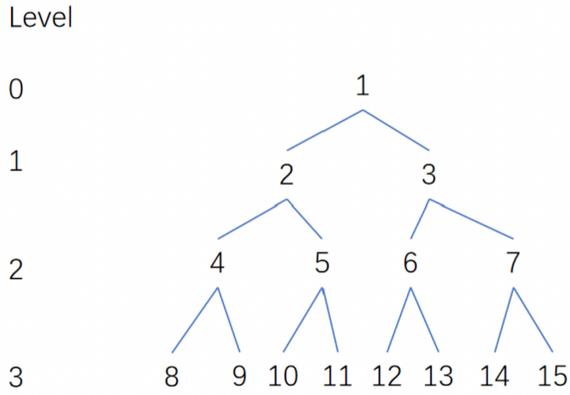


**Figure 1: Schematic diagram of the DGH tree structure used in the utility measurement of KPDP**

For KPDP, the information loss generated by this system mainly comes from the generalized information loss in the process of satisfying the k-anonymity criterion. The calculation of generalized information loss is based on the relationship between nodes on the DGH tree. The generalized information loss includes single-node generalized information loss as well as multi-node generalized information loss.

**Definition 1. Single-node generalized information loss:** The information loss incurred when generalizing a node to a parent or higher level node is calculated as shown in Eq. (2).

$$Loss_g(node_i, node_j) = log_2(LF(node_i)) - log_2(LF(node_j)) \quad (2)$$

Where $Loss_g(node_i, node_j)$ is the generalization information loss generated by generalizing $node_j$ to $node_i$, $LF(node_k)$ returns the

number of leaf nodes owned by $node_k$. The special case of leaf nodes being generalized to the root is called suppression[31], and in the suppression case, the generalization information loss is calculated as shown in Eq. (3).

$$Loss_g(node_i) = H \quad (3)$$

Where H denotes the height of the DGH tree.

**Definition 2. Multi-node generalized information loss:** Any two nodes on the DGH tree need to be generalized by finding the smallest subtree containing both nodes. The Lowest Cmmon Ancestor (LCA) of two nodes is the result of their generalization. The information loss caused by generalizing two nodes to their LCA nodes is calculated as shown in Eq. (4).

$$Loss_g(node_i, node_j, node_{LCA}) = Loss_g(node_{LCA}, node_i)$$
$$+Loss_g(node_{LCA}, node_j) \quad (4)$$

Since the trajectories input to KPDP system usually has irregular geometry, in order to cluster different trajectories to achieve the purpose of trajectory anonymization of KPDP, this paper uses PSA algorithm in order to cluster multiple trajectories in PSA. We need to calculate the trajectories with the smallest relative distance and the closest shape for clustering to achieve the purpose of trajectory anonymization. In order to complete the calculation process of PSA, this paper adopts the DSA algorithm to calculate the distance between trajectories and the information loss generated in the process of clustering trajectories and uses the information loss generated by trajectory alignment as a measure of the relative distance between trajectories in clustering. In DSA, the generalization information loss of generalizing two trajectory points and suppressing a certain trajectory point is calculated based on the DGH tree generalization model with the corresponding dimensional attributes. According to Eq. (3) and Eq. (4), for any two trajectories and, when DSA is performed on these two trajectories, the recursive equation of dynamic programming is shown in Eq. (5).

$$SAmatrix[i][j] =$$

$$min \begin{cases} SAmatrix[i-1][j-1] + (Loss_g(p_i.X, q_j.X, X_{LCA}) \\ +Loss_g(p_i.Y, q_j.Y, Y_{LCA})), \\ SAmatrix[i][j-1] + (Loss_g(q_j.X) + Loss_g(q_j.Y)), \\ SAmatrix[i-1][j] + (Loss_g(p_i.X) + Loss_g(p_i.Y)) \end{cases} \quad (5)$$

## 3.2 KPDP Workflow

KPDP is mainly composed of two parts, which are the Partition model and the Anonymization model, the trajectory dataset of multiple users is the input of KPDP, and the anonymized trajectory dataset is the output of KPDP. In this case, because the length difference of two trajectories close to each other is large, the information loss from DSA alignment is large, and thus the two trajectories cannot be grouped into one cluster in the clustering algorithm based on the distance of the trajectories, thus makes the clustering in Anonymization model less effective and generates a larger information loss. As shown in Figure 2, it can be found from Eq. (4) that in the process of aligning trajectory $tr_1$ with trajectory $tr_2$, $p_1$ and $q_1$, $p_2$ and $q_2$ are generalized to multiple nodes, while $q_3$, $q_4$, $q_5$ are generalized to the root node of DGH tree by a single node, and this process will produce excessive information loss. In this paper, we set up a Partition model to reduce the information loss

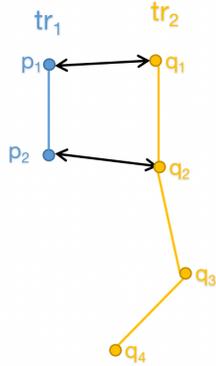of KPDP anonymization while ensuring the requirement of KPDP anonymization.



**Figure 2: Schematic diagram of the DGH tree structure used in the utility measurement of KPDP**

The specific workflow is shown in Figure 3. The trajectory dataset needs to be preprocessed by the Partition model first, which enables all trajectories to be processed in advance to keep the original geometric features of trajectories in the Anonymization model as much as possible, as well as to prevent excessive information loss in the process of Anonymization model. The partition model prevents the loss of information in the process of the Anonymization model. The processed datasets are transferred from the Partition model to the Anonymization model, which uses the PSA algorithm and the adaptive DBSCAN clustering algorithm proposed in this paper to complete the trajectory clustering, and finally outputs the anonymized trajectory datasets of the KPDP system.
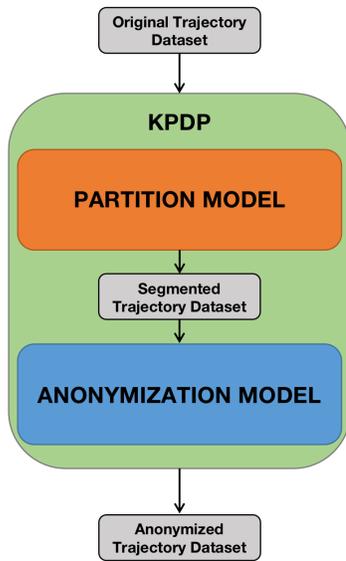


**Figure 3: KPDP Workflow**

# 4 PARTITION MODEL

Based on the workflow of KPDP in Section 3.2, this section focuses on the segmentation preprocessing of trajectories to reduce the generalization information loss of trajectories afterwards. This process refers to segmenting the trajectories based on point density before anonymizing the trajectory set so that the released dataset will retain the distribution features of the trajectories and reduce the generalization information loss in the alignment and clustering steps as much as possible. We illustrate the three main steps of the partition model - generating auxiliary points on the trajectory, then clustering the point set, and segmenting the trajectory based on the clustering distribution. The steps are interlocked to make the length of trajectories relatively average. The specific segmentation preprocess of the original trajectory set is schematically shown in Figure 4. First, for the original trajectory set in Figure
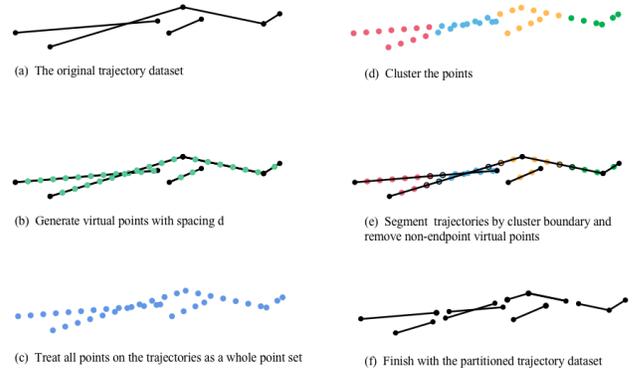


(a) The original trajectory dataset

(d) Cluster the points

(b) Generate virtual points with spacing d

(e) Segment trajectories by cluster boundary and remove non-endpoint virtual points

(c) Treat all points on the trajectories as a whole point set

(f) Finish with the partitioned trajectory dataset

**Figure 4: Schematic diagram of the preprocessing**

4(a), Figure 4(b) shows the generation of green auxiliary points on the trajectory with the same distance $d$. All the points in Figure 4(c), including the actual existing and virtual auxiliary points, are considered whole point sets for clustering. The clusters of points in Figure 4(d) are distinguished from each other by different colours, i.e., points of different colours belong to different clusters. These points are mapped back to the trajectory set in Figure 4(e), and the trajectories are partitioned at the neighbouring points belonging to different clusters according to the boundaries of the point clusters. The final result is shown in Figure 4(f), where the auxiliary points as trajectory endpoints after segmenting are kept in the trajectory set and form a new segmented dataset with other actual points.

We conducted extensive experiments to evaluate our segmentation model. The results demonstrate that, compared with the direct method of clustering and generalizing the trajectories, adding the preprocessing step can not only effectively reduce the overall generalization information loss but also speed up the running of the trajectory clustering algorithm.

## 4.1 Auxiliary Point Generation

Before segmenting the trajectory, relatively dense auxiliary points are added between adjacent points. These virtual auxiliary points are equally spaced, as shown in Figure 4(b). Since the actual points

on the trajectory are time-ordered, the primary purpose of setting auxiliary points is to make line segments of different lengths have the same effect on the density of points in their neighbourhoods so that the line segments represented by points can be more similar to the solid form of the line in space. auxiliary points are defined as follows.

**Definition 3. Auxiliary Point:** Points that do not exist in the trajectory dataset and are used to reflect the spatial distribution structure of the trajectory. For the line segment formed between two time sequence adjacent points, starting from the end of the previous time sequence, a auxiliary point is added for each fixed distance $d$ along the line segment.

The smaller the distance between the generated auxiliary points, the better the point set can reflect the distribution shape of trajectories in space. In contrast, if the spacing is too small, it will increase the amount of processed data and affect the operation efficiency.

## 4.2 Point Set Clustering

In order to make the length of trajectories relatively uniform, partitioning trajectories based on the difference in spatial trajectory density is our proposed solution. Regarding distribution, the density of trajectories in macroscopic space is reflected as the density difference of the points on the microscopic level. The point clustering algorithm can automatically gather the close points into clusters, reflecting the density distribution of points on the plane space.

For the trajectory dataset with auxiliary points, all coordinate points on the trajectory are regarded as the whole point set to be clustered. Meanwhile, the mapping relationship between each point and the cluster it belongs to is recorded for the subsequent segmentation operation.

We called the k-means point clustering method from the machine learning Sklearn library to divide the points into $k$ clusters based on the spatial Euclidean distance between them. The k-means algorithm is one of the most basic and widely used clustering algorithms that can divide data samples with different attribute values into a designated number of clusters and use the mean of all samples within each cluster as the representative points [26]. The main idea is to divide the data set into different classes by iteratively adjusting the clustering centres so that the mean error criterion function, which measures the clustering performance, is optimal, thus ensuring that the generated clustering results are compact within clusters and sparse from each other.

The effective operation of the clustering algorithm is generally based on the homogenization and standardization of the data feature variables. Since the attribute values used to calculate the Euclidean distance between trajectories only contain two dimensions, longitude and latitude, and there is no significant data disparity with uniform magnitude, the k-means algorithm can be directly applied to divide the point set on the plane space map into clusters.

## 4.3 Trajectory Segmentation

In this stage, we use the clustering boundaries generated by the k-means algorithm to segment the trajectories to reduce the disparity in the length of trajectories in the original dataset. Referring to [20], the sum of segmented trajectories is not necessarily the original trajectory but a characteristic reflection of its structure

distribution. Therefore, when trajectory clustering is performed later, the segments of a trajectory may belong to several different clusters and subsequently be generalized to different anonymous trajectories. However, the accuracy of the trajectory clustering will be relatively higher due to the reduced cost of information loss when aligning long and short trajectories later. In contrast, the overall trajectory clustering will lose more detailed features and incur higher generalization information loss during generalization. In the KPDP framework, after clustering the segmented trajectories, the length of trajectories within each cluster is relatively consistent, so the shape of the anonymous trajectories will be more reasonable.

The segmentation process of the trajectory set is described as follows. Iterate through each trajectory in the trajectory dataset containing auxiliary points and check whether the adjacent points on a trajectory belong to the same cluster. If they are not the same, the trajectory is segmented, and a new trajectory is generated. When the endpoint of a segmented trajectory is a auxiliary point, it will be added to the newly generated trajectory dataset as a real trajectory point, while other non-endpoint auxiliary points will be removed and will not be involved in the subsequent privacy-preserving processing. The pseudo-code for generating a segmented trajectory dataset is shown in Algorithm 1, whose input is the trajectory dataset containing auxiliary points.

---

**Algorithm 1:** Trajectory segmentation algorithm

**input** : Dataset $T$
**output**: Partitioned Dataset $T_{partitioned}$

Let $T_{partitioned}$ be an empty set that will store the new partitioned trajectory dataset;
**for** $tr$ **in** $T$ **do**
    Let $new_{tr}$ be an empty set that will store the new trajectory;
    Append $tr[0]$ as the first point to $new_{tr}$;
    **for** $p$ **in** $range(0, len(tr) - 1)$ **do**
        **if** *the point p and the adjacent point p+1 belong to the same cluster* **then**
            Append $tr[p+1]$ to $new_{tr}$
        **else**
            **if** *point p is not a real point* **then**
                Append $tr[p]$ to $new_{tr}$
            Append $new_{tr}$ to $T_{partitioned}$;
            Let $new_{tr}$ be an empty set that will store the new trajectory again;
            Append $tr[p+1]$ to $new_{tr}$;
    Append $new_{tr}$ to $T_{partitioned}$;
**return** $T_{partitioned}$

---

On the one hand, the maximum value of distance lost in segmenting the trajectory is $d$ because the spacing will not be smaller than the distance between the adjacent auxiliary points and the actual point or between two actual points when generating virtual auxiliary points along the trajectory direction before. In order to make the segmented trajectory closer to the original one, the parameter $d$ should be as small as possible without causing the algorithm to be overly complicated so that the loss due to segmentation can be

minimized when cutting the line segment between two adjacent points. On the other hand, the trajectory segmentation should not only ensure accuracy but also have simplicity, i.e., use as few points as possible to characterize the shape of the trajectory. The virtual auxiliary points that are not endpoints on the trajectory do not contribute significantly to the subsequent generalization process of the trajectory but rather increase the time complexity of the alignment algorithm, so they are discarded when generating the new segmented trajectory dataset.

## 5 ANONYMIZATION MODEL

In order to achieve the anonymity requirement, we introduce clustering algorithms that can gather data samples based on similarity. Clustering is an unsupervised learning method in the field of machine learning that is capable of discovering patterns implicit in a dataset. By clustering the preprocessed trajectory set wisely, it can produce a low information loss during generalization and anonymization, further maintaining the distribution characteristics and data utility of the original trajectory set. In this paper, two trajectory clustering algorithms are considered to construct the anonymization model, respectively, the iterative k'-means algorithm and the adaptive DBSCAN algorithm. Both use the alignment information loss obtained by DSA as the distance indicator between two trajectories and provide a design such that the number of trajectories contained in each cluster is no less than $k$, ensuring compliance with the privacy-preserving requirement of k-anonymity. Among them, the adaptive DBSCAN algorithm is the primary one that this paper focuses on as a method that can significantly improve the utility of the anonymous trajectory dataset and reduce the model running time, while the iterative k'-means algorithm is mainly used for comparison. These two algorithms run independently in the anonymization model of KPDP. After clustering, KPDP will apply PSA to generalize the trajectories of each cluster to derive the anonymous trajectory set for publication.

### 5.1 Iterative K'-means algorithm

We borrowed the idea from [31] to perform k'-means clustering on trajectories (where the "'" is used to distinguish the "k" that has different meanings in k-means and k-anonymity) and ensure the number of trajectories within each cluster is at least $k$ by iteration. K'-means is a distance-based clustering algorithm. Its clustering similarity is calculated using the mean distance between objects within each cluster. The brief idea is to divide data objects into $k'$ clusters according to the input value of $k'$, making the similarity within each cluster higher and the similarity between different clusters lower. The iterative k'-means algorithm is used for comparison with the adaptive DBSCAN algorithm.

The basic k'-means algorithm works by first selecting any $k'$ objects from the dataset as the initial cluster centers and assigning the remaining objects to the most similar clusters (i.e., closest in the distance) to them based on their similarity (usually Euclidean distance). Then for each cluster, a new cluster center is calculated based on the mean value of the distances of all objects in the cluster. This process is repeated until the cluster centers no longer change or the standard measure of clustering performance converges.

Measuring the relative distance between trajectories is a major difficulty with an irregular data structure like spatiotemporal trajectories. The iterative k'-means algorithm in this paper uses the information loss generated by DSA of two trajectories to measure the relative distance of trajectories. In addition, when designing the trajectory clustering algorithm based on k'-means, many technical details need to be adjusted according to the characteristics of the trajectory data and the rationality of the processing method so that the iterative k'-means algorithm can effectively cluster and generalize the trajectories in the trajectory privacy protection model. Its workflow is described as follows: **(1)** Calculate the initial number of clusters based on the value of $k$ required for k-anonymity and the number of trajectories in the dataset. **(2)** A randomly selected trajectory from the trajectory set is used as the initial clustering center for each cluster. **(3)** Assigning all trajectories in the trajectory set to the cluster center that produces the least loss of alignment information with its DSA. **(4)** Apply the PSA algorithm to each cluster and generalize and merge the trajectories it contains to form a new cluster center. **(5)** Repeat steps (3) and (4) until the trajectories contained in each cluster no longer change, completing k'-means clustering of trajectories. **(6)** Dissolve the clusters containing less than $k$ trajectories and repeat the steps of k'-means clustering until all clusters conform to k-anonymity.

Compared with the basic k'-means algorithm, the iterative k'-means algorithm gets the centers of a cluster of trajectories by PSA, except for the initial clustering centers randomly selected from the set of trajectories. In addition, the ordinary k'-means algorithm determines whether to perform the next clustering iteration based on the change of the cluster centers. However, due to the specificity of the generalization trajectory, when the cluster assignment is no longer changed, it marks the iteration stop to reduce the algorithm complexity.

Theoretically, the iterative k'-means algorithm is random for selecting initial clustering centers. This may lead to a high overall generalization information loss by generalizing each cluster when the distribution of initial clustering centers is poor, reducing the data utility of the final trajectory set used for publication. The experimental performance of the iterative k'-means algorithm on real datasets will be discussed in Section 6.

### 5.2 Adaptive DBSCAN Algorithm

Inspired by the iterative k'-means algorithm, we propose the adaptive DBSCAN algorithm, which can capture the distribution characteristics among trajectories with more details. DBSCAN is a density-based spatial clustering of applications with noise, which measures the similarity between data samples in terms of density [11]. Compared with k'-means, DBSCAN does not need cluster centres to instruct clustering. Instead, it searches for high-density regions separated by low-density regions through the density connectivity of the samples. These separated high-density regions are the corresponding clusters to which the corresponding samples belong. In contrast to k'-means, which is unable to discover spherical clusters, DBSCAN can not only discover clusters of arbitrary shapes but also has special treatment for noisy samples to suppress the influence of abnormal data on clustering.

The basic DBSCAN algorithm requires two parameters to be entered before working: the neighbourhood radius *epsilon* and the minimum number of samples contained in the neighbourhood *minPts*. Once it starts running, the DBSCAN algorithm will traverse and label each sample in the dataset. First, for any sample that has not been labelled, find all samples whose relative distance to it is within *epsilon*. If the number of samples contained in the neighbourhood of the sample reaches the threshold indicator *minPts*, the sample and all samples in its neighbourhood will form a cluster, and the sample will be marked as visited. Then recursive processing is performed for the other samples in that cluster to extend the cluster by the same steps.

Conversely, if the number of samples contained in the neighbourhood of that sample is less than *minPts*, the sample is temporarily marked as noise. Once the recursion is over, the cluster has been sufficiently extended, i.e. all samples are marked as visited. The algorithm then proceeds to traverse the points in the dataset, and the points that have not been labelled are processed similarly. The basic DBSCAN algorithm outputs clusters from sample density expansion and possibly noisy samples that are still labelled as the noise at the end of the algorithm.

We conducted an intensive study on the utilization of DBSCAN ideas for the trajectory clustering algorithm and proposed an adaptive DBSCAN algorithm that meets the privacy preservation requirement. Similar to the iterative k'-means algorithm, the adaptive DBSCAN algorithm measures the relative distance of two trajectories by the information loss generated by DSA. In order to make the anonymous trajectory dataset generated by clustering and generalization fulfil the k-anonymity criterion, we assign $k$ as the value of *minPts* in the adaptive DBSCAN algorithm. This is because the parameter *minPts* is the threshold indicator of whether a trajectory is clustered with its neighbouring trajectories, so as long as the value of *minPts* is greater than $k$, it can ensure that the number of trajectories within each cluster is at least $k$, resulting in k-anonymity of the trajectory dataset.

As for the noisy samples that may exist in DBSCAN, we also handled them specifically in the trajectory privacy-preserving scenario. Analogous to the iterative k'-means algorithm, the adaptive DBSCAN algorithm repeatedly calls the core DBSCAN code until all clusters satisfy k-anonymity in order to make the noisy trajectories eventually satisfy the anonymity requirement as well. The noisy trajectories formed by DBSCAN each time will become the new input dataset for the next clustering, while another input parameter, the neighbourhood radius *epsilon*, will be enlarged appropriately to lower the judgment criterion of density connection between samples so that those noisy trajectories can be clustered more easily. The algorithm will not stop calling the DBSCAN core code until there are no more noisy samples in the dataset.

The pseudo-code of the adaptive DBSCAN algorithm is shown in Algorithm 2. It takes the trajectory dataset, the value of $k$ of the k-anonymity criterion and the neighbourhood radius parameter *epsilon* as inputs, and it outputs the clustered trajectory dataset. Its workflow is described as follows: **(1)** For trajectories that have not been labelled in the trajectory set, **(2)** if the number of trajectories with the information loss generated by DSA with that trajectory is less than the neighbourhood radius *epsilon* is greater than the threshold *minPts* (which takes the value of $k$), find all trajectories connected to that trajectory to form a cluster, and mark all trajectories in the cluster. **(3)** otherwise, mark the trajectory as noise, find the next unmarked trajectory, and repeat the previous step until all trajectories are marked. **(4)** For the set of trajectories that are still marked as noisy at this time, adaptively enlarge the value of the neighbourhood radius epsilon and repeat the above steps until all the generated clusters satisfy k-anonymity.

---

**Algorithm 2:** Adaptive DBSCAN algorithm

**input** : Dataset $T$, Anonymity Criterion $k$, Neighbor Radius *epsilon*

**output** : Trajectory Cluster Dataset $T_{clus_k}$

Let $T_{clus_k}$ be an empty set that will store the clusters with at least $k$ trajectories;

**while** *true* **do**

    $T, T_{clus_k} \leftarrow$ **TrajectoryDBSCANClustering**$(T, epsilon, k)$;

    *Append the cluster in $T_{clus}$ to $T_{clus_k}$*;

    **if** $|T| < 2 * k$ **then**

        *Cluster $T$'s remaining trajectories together and append the last cluster to $T_{clus_k}$*;

        **break**

    **if** $epsilon < top_{epsilon}$ **then**

        Increase the value of *epsilon adaptively*

    **else**

        Cluster $T$'s remaining trajectories together and append the last cluster to $T_{clus_k}$;

        **break**

**return** $T_{clus_k}$

---

In the loop of the algorithm, the value of neighbourhood radius *epsilon* will be changed adaptively based on the statistical distribution of the relative distance between trajectories. For example, in the first several rounds, the density of the trajectories is high, and the relative distances between trajectories will be concentrated at a low level. If we increase the neighbourhood radius *epsilon* by a small margin, we can efficiently cluster the trajectories in the area of high density. As looping times increase, the number of trajectories with low relative distances decreases, and the main distribution of distances among trajectories will tend to a higher value range. At this point, the neighbourhood radius *epsilon* must be raised by a greater magnitude to reduce invalid loops of the core function and improve the efficiency of the adaptive DBSCAN algorithm.

Theoretically, the time complexity of the adaptive DBSCAN algorithm will be an order of magnitude lower than the iterative k'-means algorithm due to no iteration of cluster centres required, which significantly reduces the time consumption in the trajectory clustering session. In addition, the algorithm can specialize noisy data to enhance the data utility of the generalized trajectories, and the stability of the trajectory cluster generation process is also an advantage it has. As for the logic of the algorithm, how to adaptively adjust the value of the neighbourhood radius parameter is the key point to improve the operation efficiency. Reducing the algorithm's complexity and optimizing the parameter values are unavoidable contradictions requiring a balanced algorithm design. Extensive

experiments on a real dataset will evaluate the performance of our proposed model.

## 6 EVALUATION

This section describes the experiments on trajectory privacy protection of the KPDP framework against re-identification attacks on real-world datasets. We mainly evaluate and analyze the effectiveness of the preprocessing step of the trajectory set and the performance of two trajectory clustering algorithms and explore the role of different values of parameters in the k-anonymity criterion on the experimental procedure and results. The experimental results reflect the superior performance of our method in all aspects.

### 6.1 Dataset Introduction

The trajectory dataset used in the experiment is from the Geolife project [45–47] and the T-Drive dataset [43, 44], which consists of GPS trajectories of mobile device users in the Beijing area, specifically including the longitude and latitude information and time series relationship of trajectory points. After obtaining the basic trajectory data, the original trajectory set used for trajectory privacy protection in this paper is all the trajectories intercepted in an area on the map of Beijing, China, corresponding to the latitude and longitude ranges of $116.300000 \sim 116.316000°E$ and $39.989500 \sim 40.000000°N$. The road network model composed of this trajectory set is shown in Figure 5, where The trajectory consists of longitude and latitude coordinate points collected after a certain time interval.
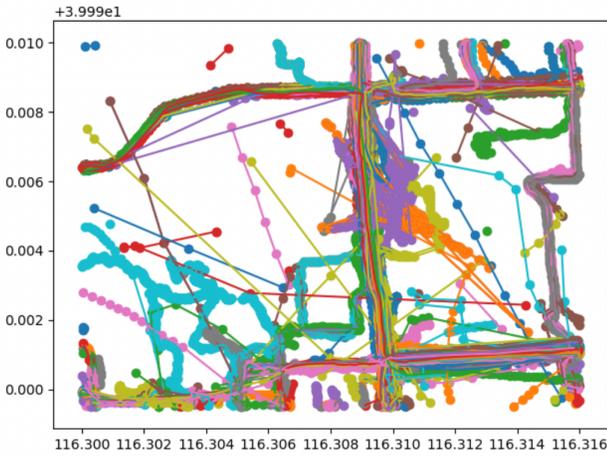


Figure 5: Road network mapping of the original trajectory set

### 6.2 Experimental Process

In the process of realizing k-anonymity privacy protection for trajectory datasets, in order to enhance the data utility of the final anonymous trajectory set, the model in this paper preprocesses the original trajectory set by segmenting the trajectories based on point density. In the preprocessing, firstly, auxiliary points used to
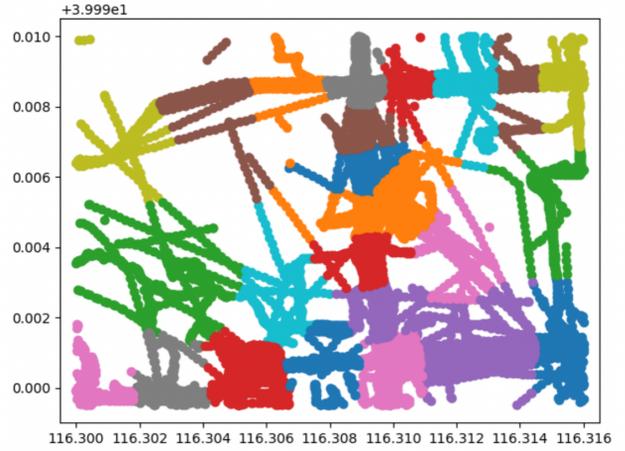


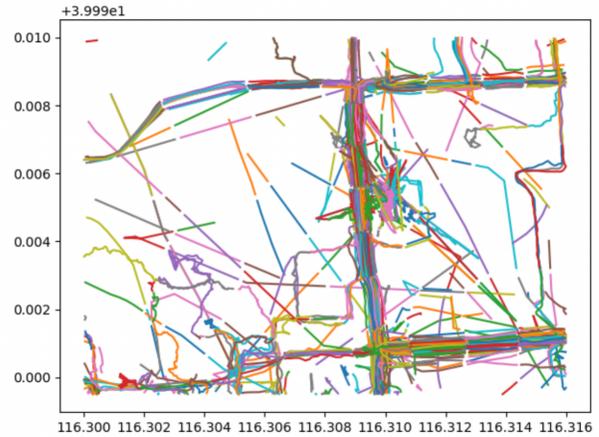Figure 6: Trajectory point set clustering results



Figure 7: Road network mapping of the segmented trajectory set

reflect the spatial distribution of trajectories need to be generated on the trajectories, and then all the points in the trajectory set are k-means clustered, and finally, the trajectories are partitioned according to the clustering of the trajectory point set. If the adjacent points on a trajectory are grouped into different clusters, the trajectory is segmented. Figure 5 shows the clustering results of dividing all points in the trajectory set into 27 clusters by k-means after adding auxiliary points, and different colors are used in the figure to identify the different clusters to which the point set belongs. Based on the clustering in Figure 6, the segmented trajectory set generated by segmenting the original trajectory set is shown in Figure 7, and again, the trajectories are distinguished from each other by color. Due to the randomness of the k-means algorithm in selecting the initial clustering centers, the segmentation results usually vary from experiment to experiment, and the number of segmented trajectories is in the range of about 1200 to 1450, as
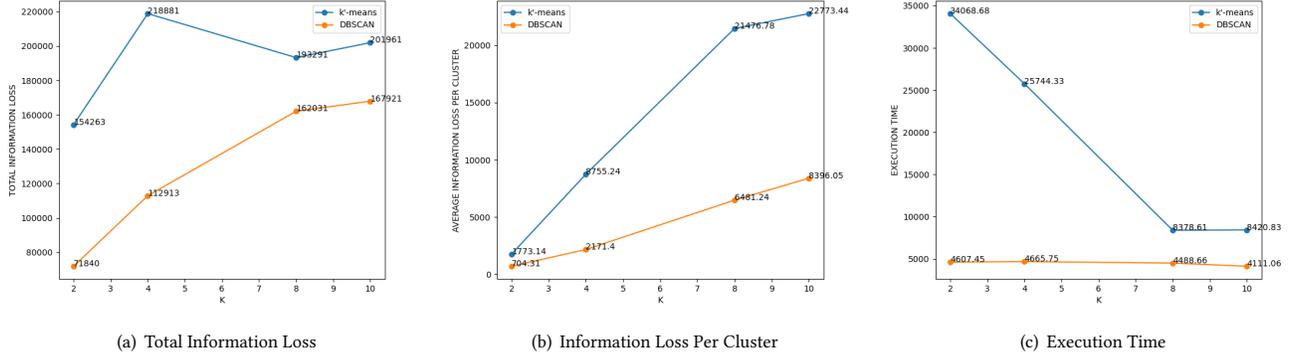
(a) Total Information Loss       (b) Information Loss Per Cluster       (c) Execution Time

**Figure 8: Comparison of the three performances of the two trajectory clustering algorithms without segmentation preprocessing at different $k$ values**
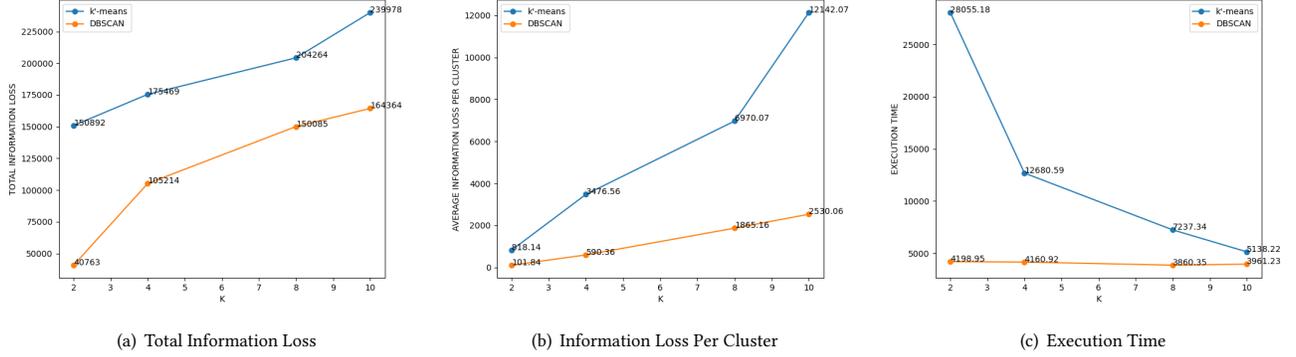


(a) Total Information Loss       (b) Information Loss Per Cluster       (c) Execution Time

**Figure 9: After segmentation preprocessing, the three performance comparisons of the two trajectory clustering algorithms at different $k$ values**

obtained from a large number of reliable repeated experiments. In the experiments in Figures 6 and 7, the number of segmented trajectories increases from 270 to 1372.

In order to make the final published trajectory dataset resistant to re-identification attacks, the trajectory privacy-preserving model needs to anonymize the trajectory set according to the selected k-values in k-anonymity criterion. In this process, the DGH tree generalization model of the trajectory set needs to be established first, i.e., the corresponding coordinate values of the trajectory points are represented by the numbers of the leaf nodes on the latitude and longitude DGH trees. Then the iterative k'-means algorithm and the adaptive DBSCAN algorithm cluster the trajectories, respectively. Finally, the clusters of trajectories formed by the clusters are generalized to obtain the trajectory dataset conforming to k-anonymity and the corresponding loss of generalization information.

## 6.3 Analysis of Experimental Results

In the experiments on trajectory privacy preservation against re-identification attacks, this paper will compare and evaluate two

trajectory clustering algorithms with and without segment preprocessing in three aspects: total information loss, average information loss per cluster, and execution time.

Figure 8 shows the comparison of the values of the three metrics obtained by the iterative k'-means algorithm and the adaptive DBSCAN algorithm for different k-anonymity metrics without trajectory segmentation preprocessing during the experiment, where the k-anonymity criterion takes the $k$ of 2, 4, 8 and 10.

As shown in Figure 8(a), the total generalized information loss of the trajectory set increases with the increases of $k$. In contrast, the protection model for trajectory, which is clustering by the adaptive DBSCAN algorithm, produces lower information loss than the iterative k'-means algorithm at all four values.

As shown in Figure 8(b), the average generalized information loss per cluster of the trajectory set also increases with increasing $k$. The information loss generated by the iterative k'-means algorithm is two to four times higher than that of the adaptive DBSCAN algorithm. The difference becomes more pronounced as the values increase.

**Table 1: KPDP performance with Partition model and without Partition model**

| Total Information Loss | | | | | Average Information Loss Per Cluster | | | | |
|---|---|---|---|---|---|---|---|---|---|
| k value | clustering algorithm | without Parition model | with partition model | reduction(%) | k value | clustering algorithm | without Parition model | with partition model | reduction(%) |
| k=2 | k'-means | 154263 | 150892 | 2.19 | k=2 | k'-means | 1773.14 | 818.14 | 53.86 |
| | DBSCAN | 71840 | 40763 | 43.26 | | DBSCAN | 704.31 | 101.84 | 85.54 |
| k=4 | k'-means | 218881 | 175469 | 19.83 | k=4 | k'-means | 8755.24 | 3476.56 | 60.29 |
| | DBSCAN | 112913 | 105214 | 6.82 | | DBSCAN | 2171.4 | 590.36 | 72.81 |
| k=8 | k'-means | 193291 | 204264 | -5.68 | k=8 | k'-means | 21476.78 | 6970.07 | 67.55 |
| | DBSCAN | 162031 | 150085 | 7.37 | | DBSCAN | 6481.24 | 1865.16 | 71.22 |
| k=10 | k'-means | 201961 | 239978 | -18.82 | k=10 | k'-means | 22773.44 | 12142.007 | 46.68 |
| | DBSCAN | 167921 | 164364 | 2.12 | | DBSCAN | 8396.05 | 2530.06 | 69.87 |

The execution time of the two trajectory clustering algorithms in the model is shown in Figure 8(c), with a decreasing trend of the algorithm execution time when increasing. In the experiments for each value, the execution time of the adaptive DBSCAN algorithm is stable within 5000 seconds, while the execution time of the iterative k'-means algorithm is much higher than the other algorithm for values 2 and 4, and relatively lower and smoother for values 8 and 10, but still higher than the other algorithm.

In Figure 9, a comparison of the values of the three metrics obtained by the iterative k'-means algorithm and the adaptive DBSCAN algorithm with different k-anonymity criteria for the dataset preprocessed by trajectory segmentation at the time of the experiment is shown, where the k-anonymity criteria take the values of 2, 4, 8 and 10.

Similar to the overall trend and comparison in Figure 8, the adaptive DBSCAN algorithm outperforms the iterative k'-means algorithm in three aspects: total generalized information loss (Figure 9(a)), average information loss per cluster (Figure 9(b)), and execution time (Figure 9(c)). Overall, the total generalized information loss and the average information loss per cluster of both trajectory clustering algorithms subsequently increase with the increase of $k$, and the execution time decreases as the value of $k$ increases gradually.

In contrast, compared with Fig 8, in both the adaptive DBSCAN algorithm and the iterative k'-means algorithm, the total information loss and average information loss per cluster of the final generated anonymized dataset after segmentation preprocessing by the partition model are relatively small, and the consumed execution time is also reduced to different degrees. Especially it is evident in the per-cluster average information loss metric of the adaptive DBSCAN algorithm, which can be obtained from Figure 9(b), that the information loss per cluster obtained by the adaptive DBSCAN algorithm decreases by about 86%, 73%, 71% and 70%, respectively, when the values of 2, 4, 8 and 10 are taken. Compared with the trajectory set without segmentation in Figure 8(b).

Such a decrease is because the preprocessing of the partition model can make the clustered trajectories closer within the clusters. Besides, the smaller the value $k$, the larger the number of clusters after segmentation, and the closer the trajectories that make up the clusters will be. The more obvious is the effect of the partition

model in reducing the information loss of generalization within the clusters.

The experimental results shown in Figure 8 and Figure 9 are consistent with the expectations of KPDP design in this paper. For the case where the generalized information loss increases with the value of $k$, this is because an increase in the value of $k$ directly leads to an increase in the number of trajectories in each cluster, resulting in a more extensive total information loss and average information loss per cluster. The superior performance of the adaptive DBSCAN algorithm in the three metrics is attributed to the ability of the algorithm to cluster trajectories close to each other more scientifically and efficiently than the iterative k'-means algorithm, with lower time complexity.

According to a series of experiments, it is proved that the way of adjusting the cluster centers in the k'-means algorithm is not fully applicable to the trajectory clustering process, while the adaptive DBSCAN algorithm forms each cluster by the expansion of the density connection between trajectories, which not only reduces the information loss but also can effectively speed up the processing of the model. The effectiveness of adding a segment preprocessing step in both trajectory clustering algorithms is due to the fact that after preprocessing, the relatively long trajectories in the dataset are avoided to be aligned and combined with shorter trajectories in the trajectory clustering and generalization process, so the information loss from the final generalization is reduced. In addition, because the long trajectories in the dataset are split into relatively short trajectories, the situation that two long trajectories are aligned with each other will be significantly avoided in the alignment, so the execution time of the trajectory clustering algorithm is also shortened.

In summary, the adaptive DBSCAN algorithm and the trajectory set segmentation preprocessing step proposed in this paper to have superior performance in controlled experiments under different scenarios, validating the theoretical expectation of reducing generalization information loss and speeding up model processing when designing the model. In the privacy-preserving phase of trajectory resistance to re-identification attacks, the trajectory preprocessing and adaptive DBSCAN algorithm for trajectory clustering to form anonymous trajectory datasets in Figure 9 has significant advantages in terms of data utility and running time for each value.

# 7 CONCLUSION

In this paper, we proposed a trajectory privacy protection framework against re-identification attacks, which can effectively anonymize the spatiotemporal trajectory dataset. We innovated a point density-based trajectory segmentation preprocessing mechanism to enable accurate clustering and generalization of trajectories. Furthermore, we applied DBSCAN in machine learning to trajectory clustering and presented the adaptive DBSCAN algorithm, which minimizes the generalization information loss to acquire higher data utility while ensuring the k-anonymity of the generated trajectory dataset. Extensive experiments on a realistic dataset also showed that there is the superiority of the short execution time of our approach compared with previous works.

## REFERENCES

[1] Adeel Anjum and Guillaume Raschia. 2017. BangA: an efficient and flexible generalization-based algorithm for privacy preserving data publication. *Computers* 6, 1 (2017), 1.

[2] Alastair R Beresford and Frank Stajano. 2004. Mix zones: User privacy in location-aware services. In *IEEE Annual conference on pervasive computing and communications workshops, 2004. Proceedings of the Second*. IEEE, 127–131.

[3] Claudio Bettini, Sergio Mascetti, X Sean Wang, Dario Freni, and Sushil Jajodia. 2009. Anonymity and historical-anonymity in location-based services. *Privacy in location-based applications: research issues and emerging trends* (2009), 1–30.

[4] Alina Campan, Nicholas Cooper, and Traian Marius Truta. 2011. On-the-fly generalization hierarchies for numerical attributes revisited. In *Secure Data Management: 8th VLDB Workshop, SDM 2011, Seattle, WA, USA, September 2, 2011, Proceedings 8*. Springer, 18–32.

[5] Xi Chen, Chen Wang, Shanjiang Tang, Ce Yu, and Quan Zou. 2017. CMSA: a heterogeneous CPU/GPU computing system for multiple similar RNA/DNA sequence alignment. *BMC bioinformatics* 18 (2017), 1–10.

[6] Richard Chow and Philippe Golle. 2009. Faking contextual data for fun, profit, and privacy. In *Proceedings of the 8th ACM workshop on Privacy in the electronic society*. 105–108.

[7] Biswanath Chowdhury and Gautam Garai. 2017. A review on multiple sequence alignment from the perspective of genetic algorithm. *Genomics* 109, 5-6 (2017), 419–431.

[8] A Ercument Cicek, Mehmet Ercan Nergiz, and Yucel Saygin. 2014. Ensuring location diversity in privacy-preserving spatio-temporal data publishing. *The VLDB Journal* 23, 4 (2014), 609–625.

[9] Jiaxin Ding. 2015. Trajectory mining, representation and privacy protection. In *Proceedings of the 2nd ACM SIGSPATIAL PhD Workshop*. 1–4.

[10] Hyo Jin Do, Young-Seob Jeong, Ho-Jin Choi, and Kwangjo Kim. 2016. Another dummy generation technique in location-based services. In *2016 International Conference on Big Data and Smart Computing (BigComp)*. IEEE, 532–538.

[11] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. 1996. A density-based algorithm for discovering clusters in large spatial databases with noise.. In *kdd*, Vol. 96. 226–231.

[12] Benjamin CM Fung, Ke Wang, and Philip S Yu. 2005. Top-down specialization for information and privacy preservation. In *21st international conference on data engineering (ICDE'05)*. IEEE, 205–216.

[13] Gabriel Ghinita, Maria Luisa Damiani, Claudio Silvestri, and Elisa Bertino. 2009. Preventing velocity-based linkage attacks in location-aware applications. In *Proceedings of the 17th ACM SIGSPATIAL international conference on advances in geographic information systems*. 246–255.

[14] Marco Gramaglia, Marco Fiore, Alberto Tarable, and Albert Banchs. 2017. k$^{\tau, \epsilon}$-anonymity: Towards Privacy-Preserving Publishing of Spatiotemporal Trajectory Data. *CoRR* abs/1701.02243 (2017). arXiv:1701.02243 http://arxiv.org/abs/1701.02243

[15] Marco Gruteser and Dirk Grunwald. 2003. Anonymous usage of location-based services through spatial and temporal cloaking. In *Proceedings of the 1st international conference on Mobile systems, applications and services*. 31–42.

[16] Vijay S Iyengar. 2002. Transforming data to satisfy privacy constraints. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*. 279–288.

[17] Kaifeng Jiang, Dongxu Shao, Stéphane Bressan, Thomas Kister, and Kian-Lee Tan. 2013. Publishing trajectories with differential privacy guarantees. In *Proceedings of the 25th International Conference on scientific and statistical database management*. 1–12.

[18] John Krumm. 2007. Inference attacks on location tracks. In *Pervasive Computing: 5th International Conference, PERVASIVE 2007, Toronto, Canada, May 13-16, 2007. Proceedings 5*. Springer, 127–143.

[19] Quan Le, Fabian Sievers, and Desmond G Higgins. 2017. Protein multiple sequence alignment benchmarking through secondary structure prediction. *Bioinformatics* 33, 9 (2017), 1331–1337.

[20] Jae-Gil Lee, Jiawei Han, and Kyu-Young Whang. 2007. Trajectory clustering: a partition-and-group framework. In *Proceedings of the 2007 ACM SIGMOD international conference on Management of data*. 593–604.

[21] Kristen LeFevre, David J DeWitt, and Raghu Ramakrishnan. 2005. Incognito: Efficient full-domain k-anonymity. In *Proceedings of the 2005 ACM SIGMOD international conference on Management of data*. 49–60.

[22] Kristen LeFevre, David J DeWitt, and Raghu Ramakrishnan. 2006. Mondrian multidimensional k-anonymity. In *22nd International conference on data engineering (ICDE'06)*. IEEE, 25–25.

[23] Huaxin Li, Haojin Zhu, Suguo Du, Xiaohui Liang, and Xuemin Shen. 2016. Privacy leakage of location sharing in mobile social networks: Attacks and defense. *IEEE Transactions on Dependable and Secure Computing* 15, 4 (2016), 646–660.

[24] Bo Liu, Wanlei Zhou, Tianqing Zhu, Longxiang Gao, and Yong Xiang. 2018. Location privacy and its applications: A systematic study. *IEEE access* 6 (2018), 17606–17624.

[25] Ashwin Machanavajjhala, Daniel Kifer, Johannes Gehrke, and Muthuramakrishnan Venkitasubramaniam. 2007. l-diversity: Privacy beyond k-anonymity. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 1, 1 (2007), 3–es.

[26] J MacQueen. 1967. Classification and analysis of multivariate observations. In *5th Berkeley Symp. Math. Statist. Probability*. University of California Los Angeles LA USA, 281–297.

[27] Mohamed F Mokbel. 2007. Privacy in location-based services: State-of-the-art and research directions. In *2007 International Conference on Mobile Data Management*. IEEE Computer Society, 228–228.

[28] Elham Naghizade, Lars Kulik, and Egemen Tanin. 2014. Protection of sensitive trajectory datasets through spatial and temporal exchange. In *Proceedings of the 26th International Conference on Scientific and Statistical Database Management*. 1–4.

[29] Pierangela Samarati. 2001. Protecting respondents identities in microdata release. *IEEE transactions on Knowledge and Data Engineering* 13, 6 (2001), 1010–1027.

[30] Pierangela Samarati and Latanya Sweeney. 1998. Protecting privacy when disclosing information: k-anonymity and its enforcement through generalization and suppression. (1998).

[31] Sina Shaham, Ming Ding, Bo Liu, Shuping Dang, Zihuai Lin, and Jun Li. 2020. Privacy preserving location data publishing: A machine learning approach. *IEEE Transactions on Knowledge and Data Engineering* 33, 9 (2020), 3270–3283.

[32] Reza Shokri, Julien Freudiger, and Jean-Pierre Hubaux. 2010. *A unified framework for location privacy*. Technical Report.

[33] Reza Shokri, George Theodorakopoulos, Jean-Yves Le Boudec, and Jean-Pierre Hubaux. 2011. Quantifying location privacy. In *2011 IEEE symposium on security and privacy*. IEEE, 247–262.

[34] Latanya Sweeney. 2002. Achieving k-anonymity privacy protection using generalization and suppression. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 10, 05 (2002), 571–588.

[35] Latanya Sweeney. 2002. k-anonymity: A model for protecting privacy. *International journal of uncertainty, fuzziness and knowledge-based systems* 10, 05 (2002), 557–570.

[36] Nilothpal Talukder and Sheikh Iqbal Ahamed. 2010. Preventing multi-query attack in location-based services. In *Proceedings of the third ACM conference on Wireless network security*. 25–36.

[37] Acar Tamersoy, Grigorios Loukides, Mehmet Ercan Nergiz, Yucel Saygin, and Bradley Malin. 2012. Anonymization of longitudinal electronic medical records. *IEEE Transactions on Information Technology in Biomedicine* 16, 3 (2012), 413–423.

[38] Zhen Tu, Kai Zhao, Fengli Xu, Yong Li, Li Su, and Depeng Jin. 2017. Beyond k-anonymity: protect your trajectory from semantic attack. In *2017 14th Annual IEEE International Conference on Sensing, Communication, and Networking (SECON)"*. IEEE, 1–9.

[39] Zhen Tu, Kai Zhao, Fengli Xu, Yong Li, Li Su, and Depeng Jin. 2019. Protecting Trajectory From Semantic Attack Considering $k$ -Anonymity, $l$ -Diversity, and $t$ -Closeness. *IEEE Transactions on Network and Service Management* 16, 1 (2019), 264–278. https://doi.org/10.1109/TNSM.2018.2877790

[40] Marius Wernke, Pavel Skvortsov, Frank Dürr, and Kurt Rothermel. 2014. A classification of location privacy attacks and approaches. *Personal and ubiquitous computing* 18 (2014), 163–175.

[41] Toby Xu and Ying Cai. 2008. Exploring historical location data for anonymity preservation in location-based services. In *IEEE INFOCOM 2008-The 27th Conference on Computer Communications*. IEEE, 547–555.

[42] Saba Yaseen, Syed M Ali Abbas, Adeel Anjum, Tanzila Saba, Abid Khan, Saif Ur Rehman Malik, Naveed Ahmad, Basit Shahzad, and Ali Kashif Bashir. 2018. Improved generalization for secure data publishing. *IEEE Access* 6 (2018), 27156–27165.

[43] Jing Yuan, Yu Zheng, Xing Xie, and Guangzhong Sun. 2011. Driving with knowledge from the physical world. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*. 316–324.

[44] Jing Yuan, Yu Zheng, Chengyang Zhang, Wenlei Xie, Xing Xie, Guangzhong Sun, and Yan Huang. 2010. T-drive: driving directions based on taxi trajectories. In *Proceedings of the 18th SIGSPATIAL International conference on advances in geographic information systems*. 99–108.
[45] Yu Zheng, Quannan Li, Yukun Chen, Xing Xie, and Wei-Ying Ma. 2008. Understanding mobility based on GPS data. In *Proceedings of the 10th international conference on Ubiquitous computing*. 312–321.
[46] Yu Zheng, Xing Xie, Wei-Ying Ma, et al. 2010. GeoLife: A collaborative social networking service among user, location and trajectory. *IEEE Data Eng. Bull.* 33, 2 (2010), 32–39.
[47] Yu Zheng, Lizhu Zhang, Xing Xie, and Wei-Ying Ma. 2009. Mining interesting locations and travel sequences from GPS trajectories. In *Proceedings of the 18th international conference on World wide web*. 791–800.