

Clustering and Labeling of a V2V Communication Dataset

based on CAM and DENM Messages with Malicious Data Injection

Analysis and Implications



INTRO



VEHICULAR COMMUNICATION SYSTEMS

Computer networks where vehicles and roadside units (RSU) serve as communicating nodes

VEHICLE-TO-VEHICLE COMMUNICATION (V2V)

Wireless information exchange about the speed and position of nearby vehicles.
Offering great promise in accident prevention.



INTRO



DATASETS

Two different types of messages: Cooperative Awareness Messages (CAM) and Decentralized Environmental Notification Messages (DENM).

OBJECTIVE OF THE WORK

Data preprocessing, clustering, introduce noise into the provided dataset to simulate potential malicious reports, labeling the data based on clustering and identifying potential outliers.



CAM MESSAGES

Cooperative Awareness Message

defined by the European Telecommunications Standards Institute (ETSI) in 2011



- Basic awareness service by sending status data to **nearby nodes**
- Distributing messages about **presence, location, and fundamental status**

**Version, ID,
Generation Time**

ID

Station Type

Reference Position

Optional Parameters

HEADER

BODY

DENM MESSAGES

Decentralized Environmental Notification Message

defined by the European Telecommunications Standards Institute (ETSI) in 2011



- Notification service regarding **road status**
- Support active road **safety** applications

**Version, ID,
Generation Time**

Management

Situation

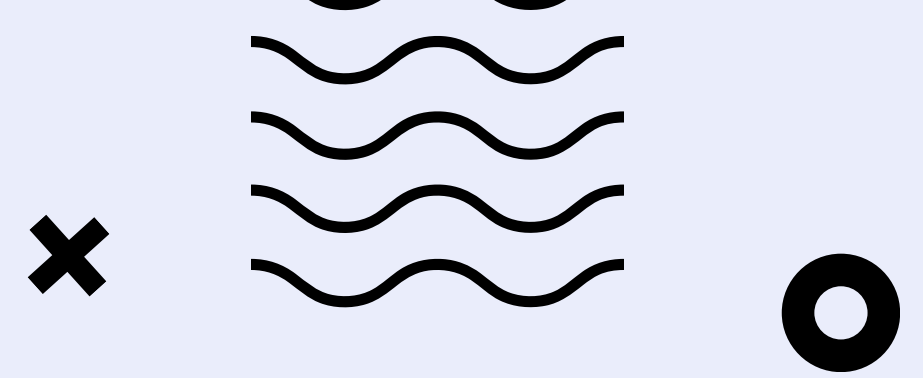
Location

HEADER

BODY



DATA ANALYSIS



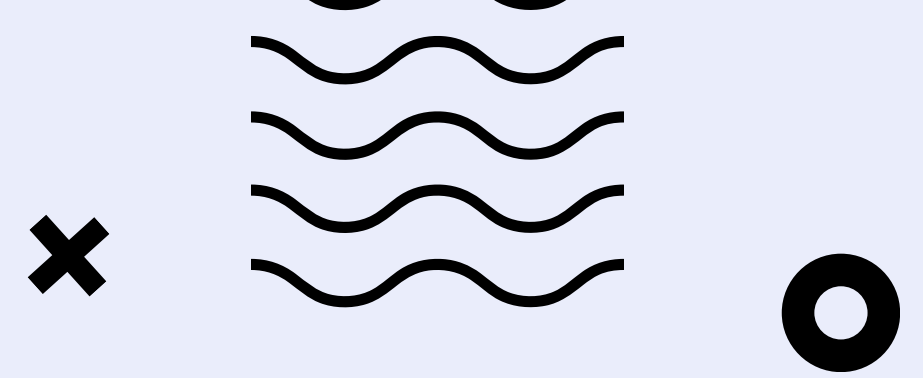
X-Means Algorithm

A **variant of K-Means algorithm**, determines **automatically** the optimal number of clusters in the data without requiring a predefined specification, based on **recursion**. Part of pyclustering open-source library.

- Initially applies **standard K-Means** with an initial cluster count
- Assesses **clustering quality using measures**
- Checks if **splitting clusters** improves overall quality
- **Divides clusters with K-Means** if advantageous



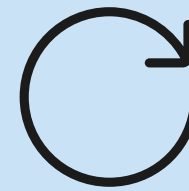
DATA ANALYSIS



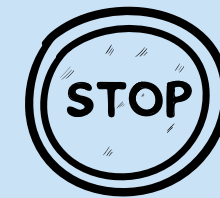
X-Means Algorithm

A **variant of K-Means algorithm**, determines **automatically** the optimal number of clusters in the data without requiring a predefined specification, based on **recursion**. Part of pyclustering open-source library.

- Initially applies **standard K-Means** with an initial cluster count
- Assesses **clustering quality using measures**
- Checks if **splitting clusters** improves overall quality
- **Divides clusters with K-Means** if advantageous



Repeats division and evaluation for existing clusters and **allows potential creation of new sub-clusters for improved clustering**

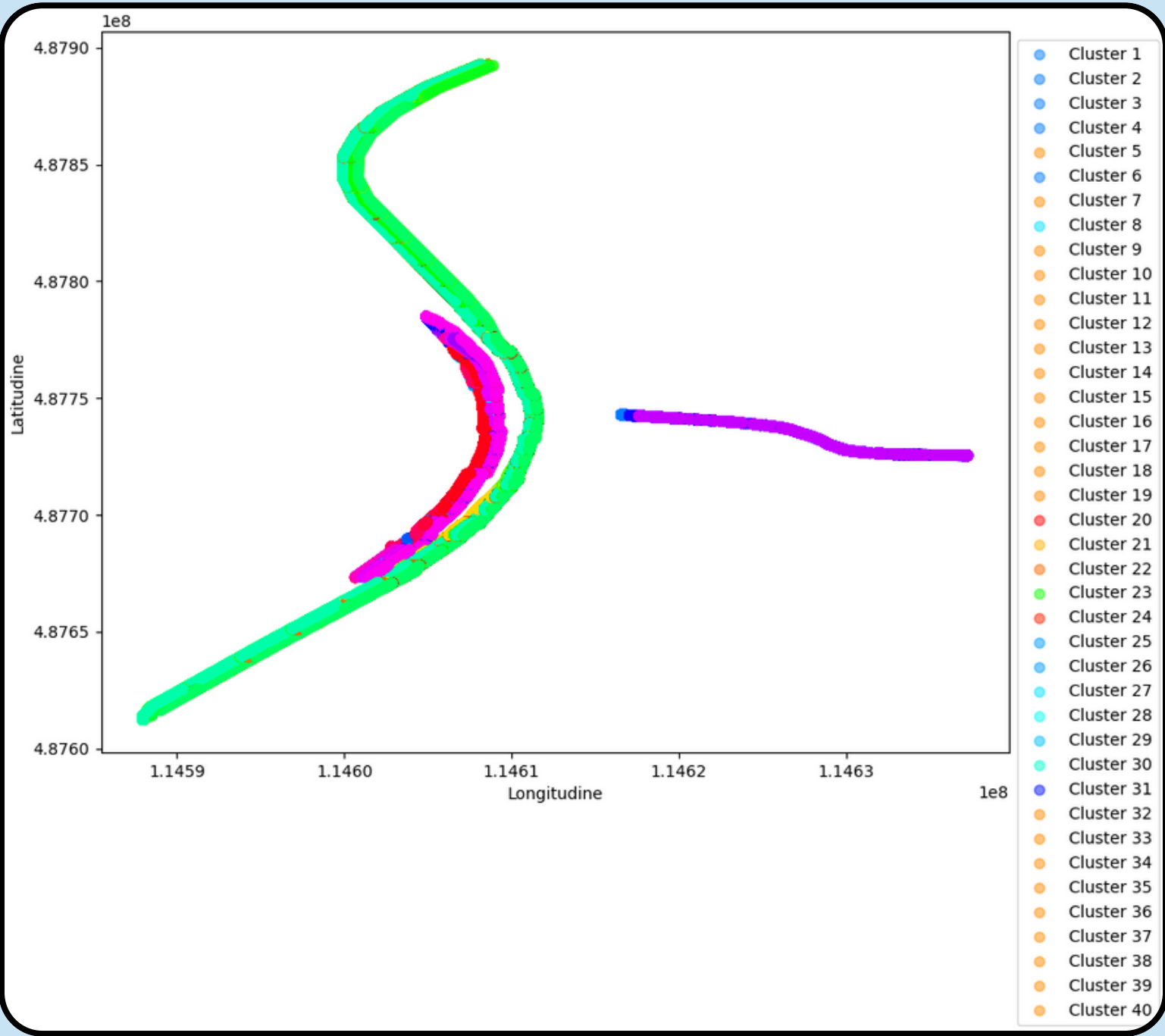


- Stops when **no further cluster division is possible**
- Halts if **the division doesn't significantly enhance results** compared to complexity

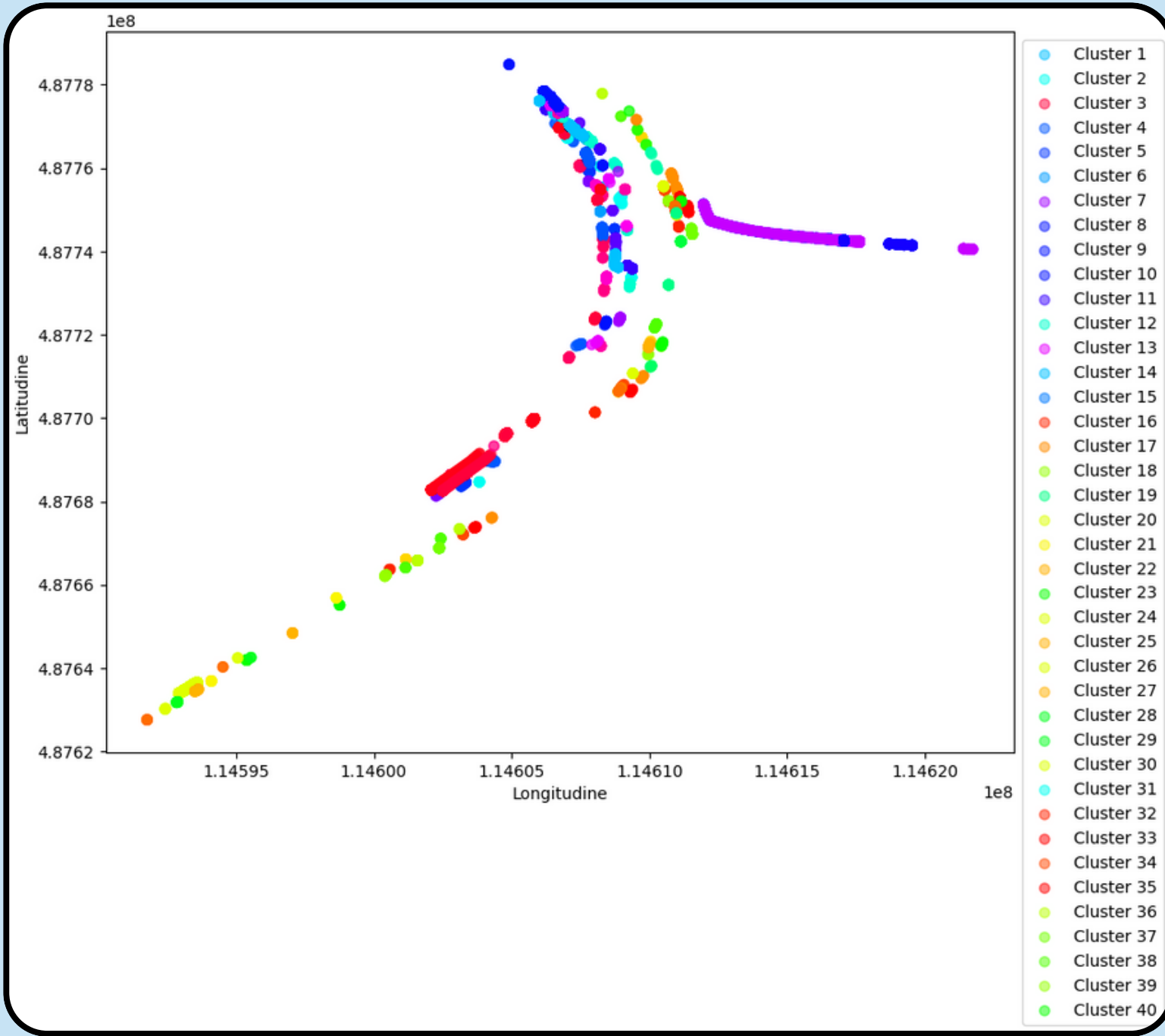


Returns **optimal clusters determined automatically**

DATA ANALYSIS



Clustering of the provided CAM dataset



Clustering of the provided DENM dataset

INTRODUCTION OF MALICIOUS DATA

Variation of coordinates using a **Gaussian distribution** (mean: 0, standard deviation: 1) multiplied by a factor of 100



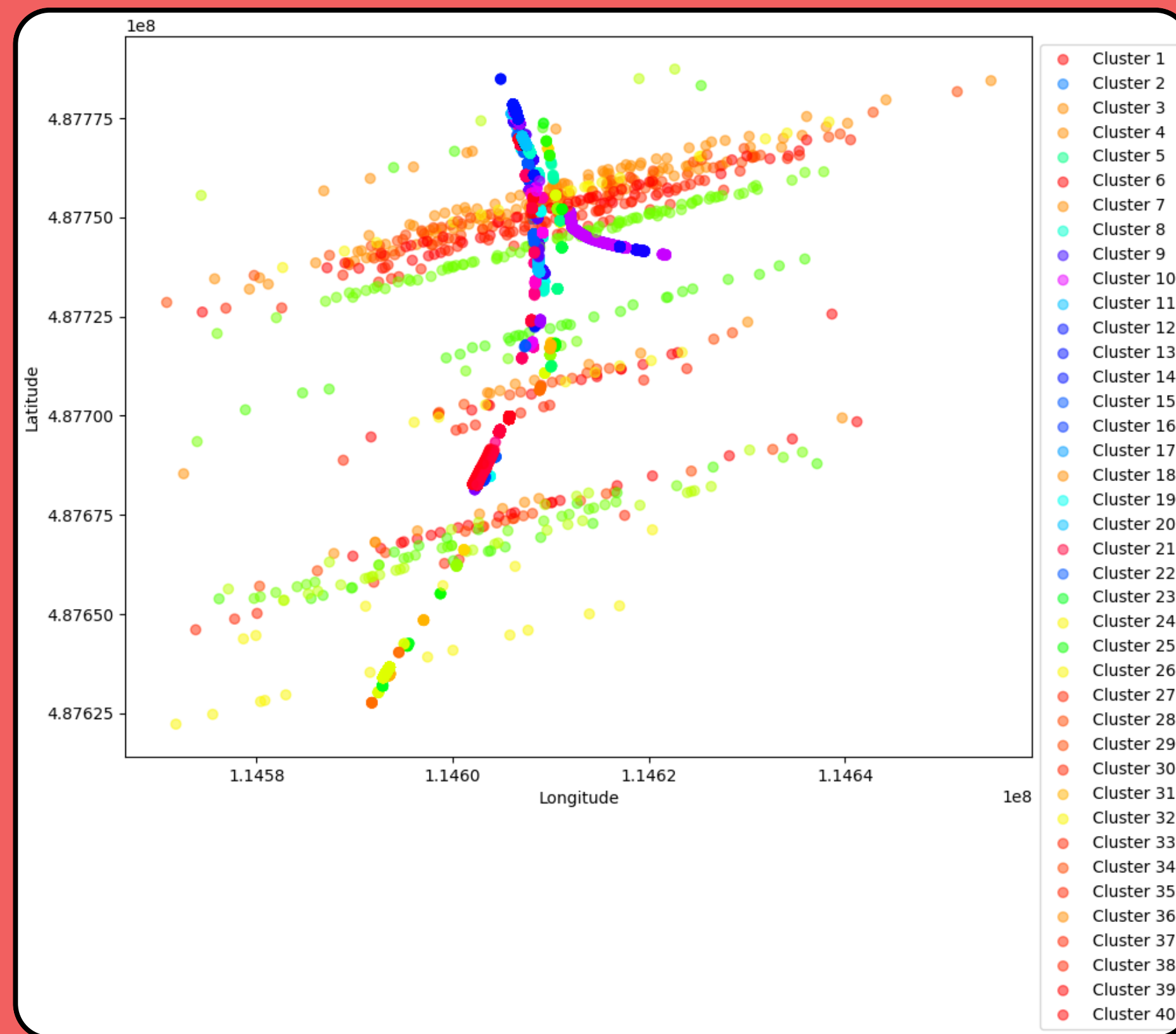
Variations within a range of approximately 100 to 900 meters on the map



- Contamination of data based on the number of sources to simulate malicious vehicles
- Focus contamination on **sources with eventType value 97**, representing the most significant cluster
- Contamination of **20% (8 sources) of this specific eventType**



INTRODUCTION OF



Clustering Following
Data Contamination

EVENT TYPE ANALYSIS



CauseCodeType_dangerousEndOfQueue = 27
CauseCodeType_collisionRisk = 97
CauseCodeType_trafficCondition = 1

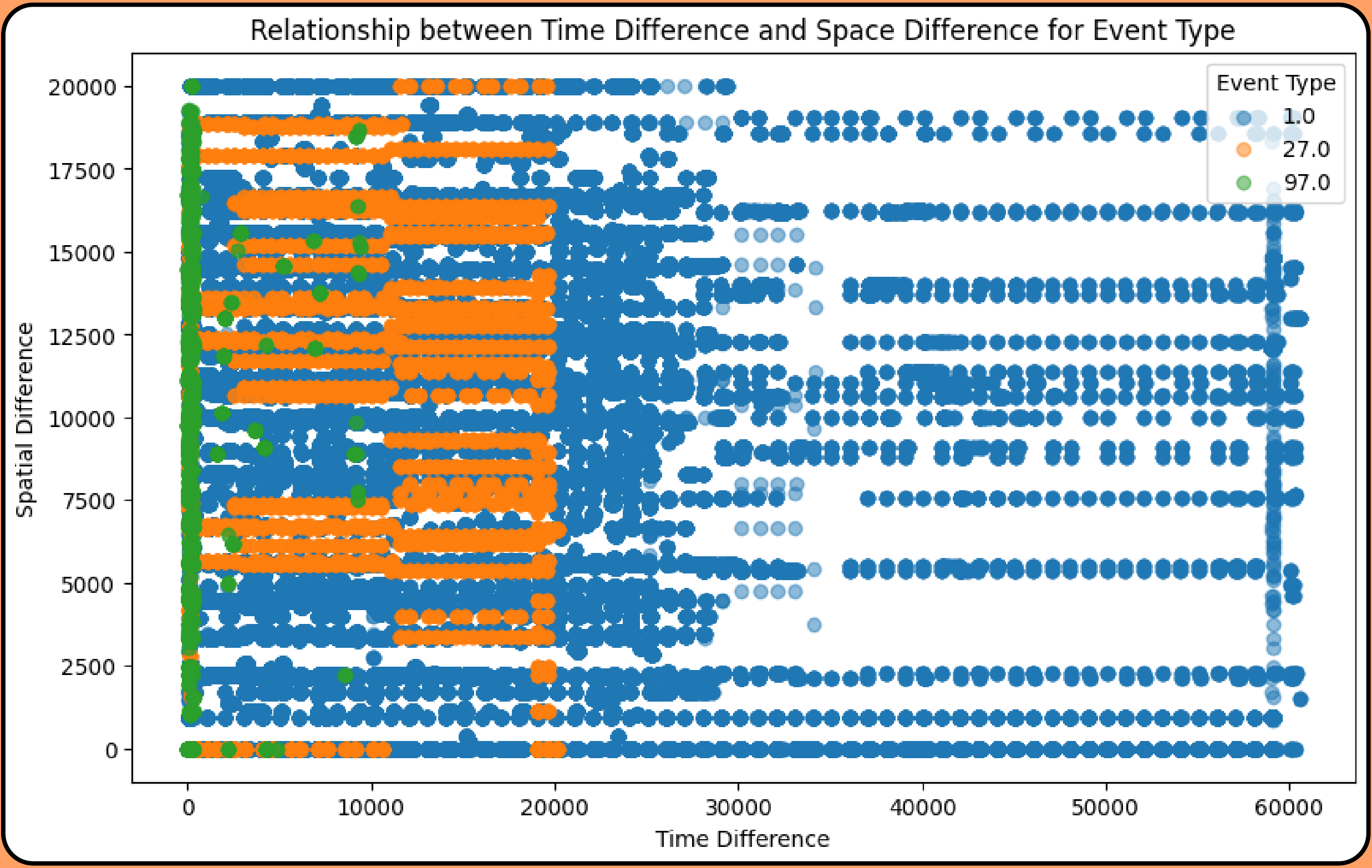
The time differences between **when a vehicle perceives an event and when it reports it.**



EVENT TYPE

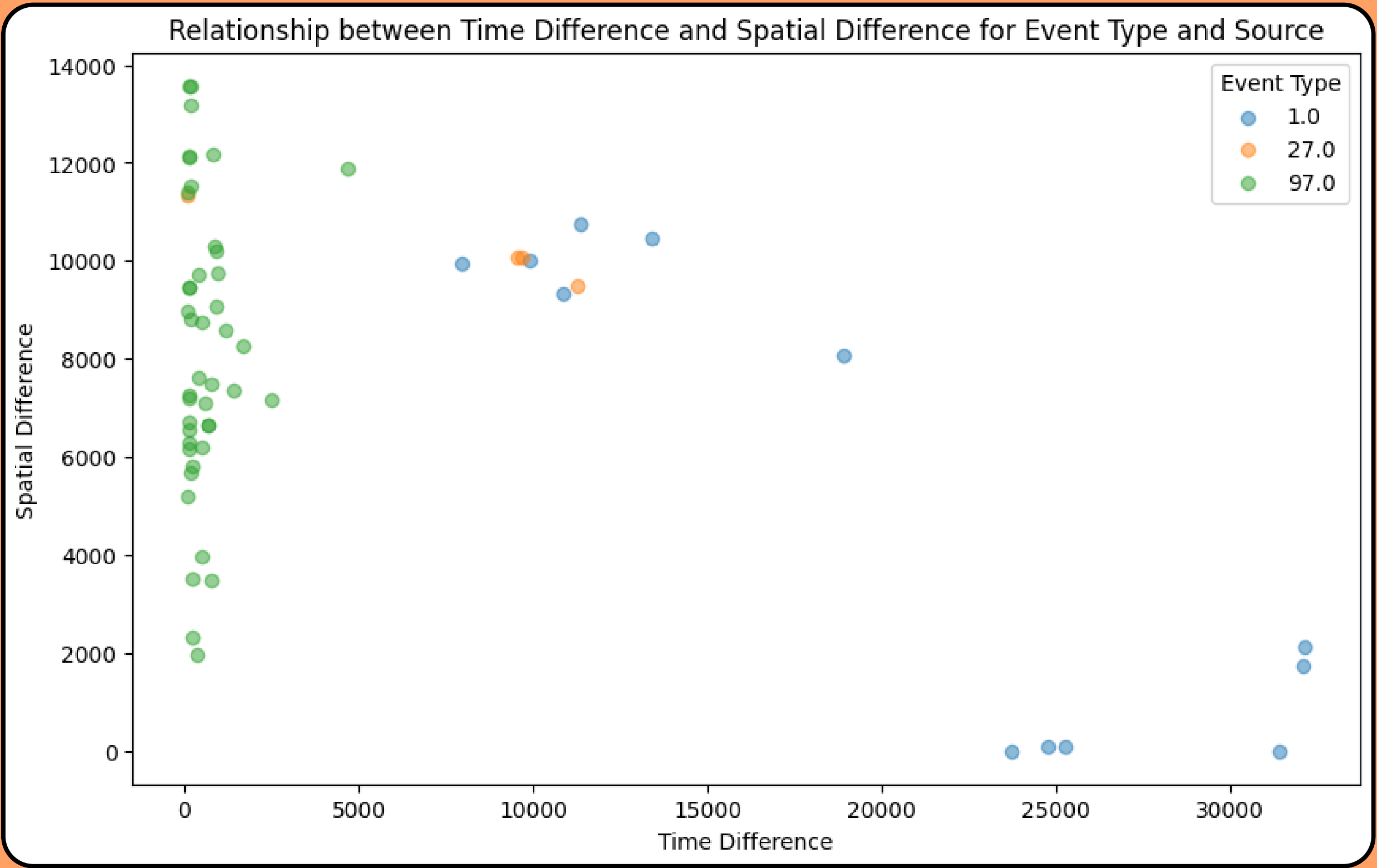
A

Ca
Ca
Ca



EVENT TYPE

A



Ca
Ca
Ca

OUTLIER DETECTION



- **DBSCAN** (Density-Based Spatial Clustering of Applications with Noise) **algorithm**

- Groups points based on the data density in space

- Two parameters:

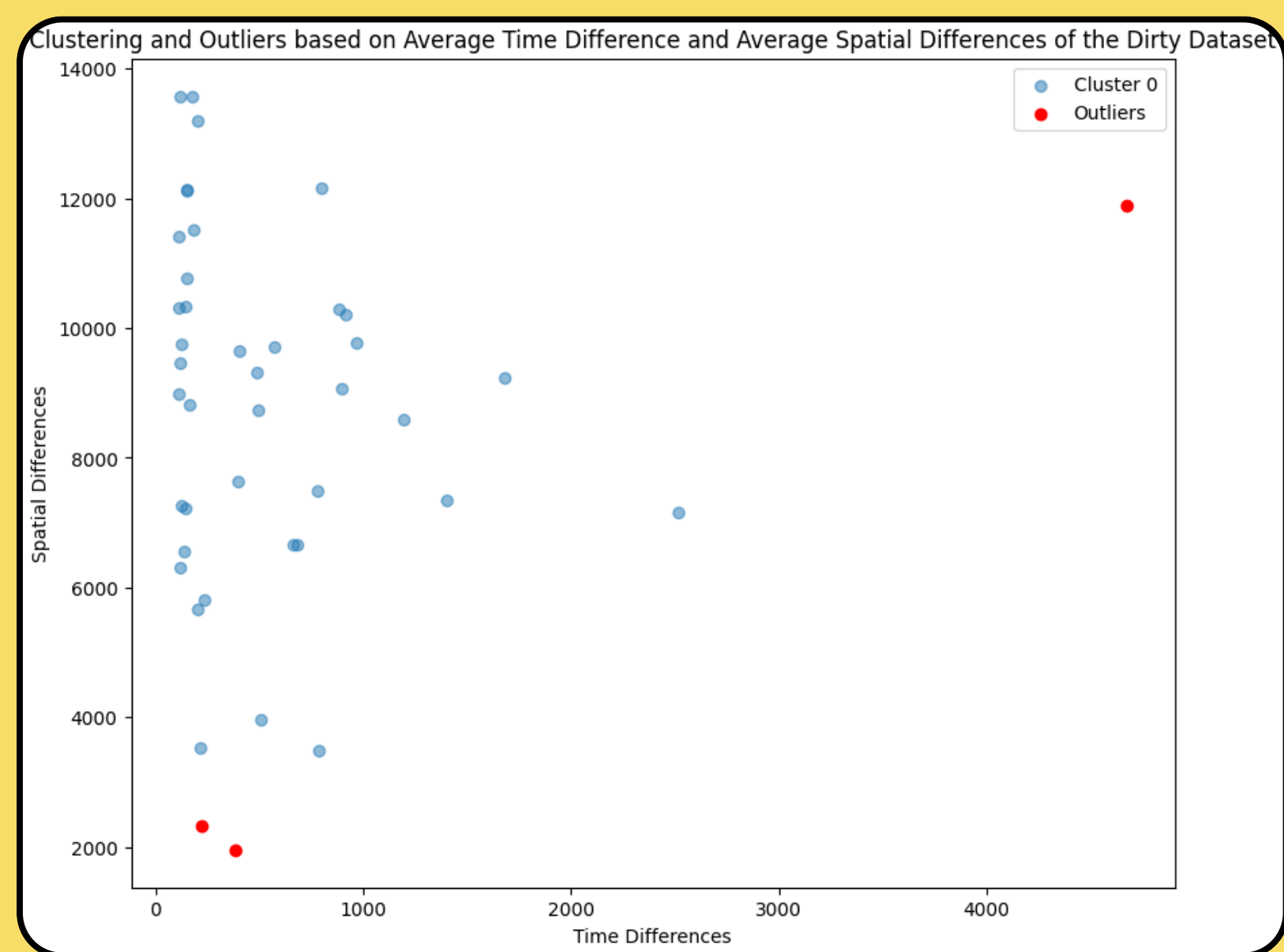
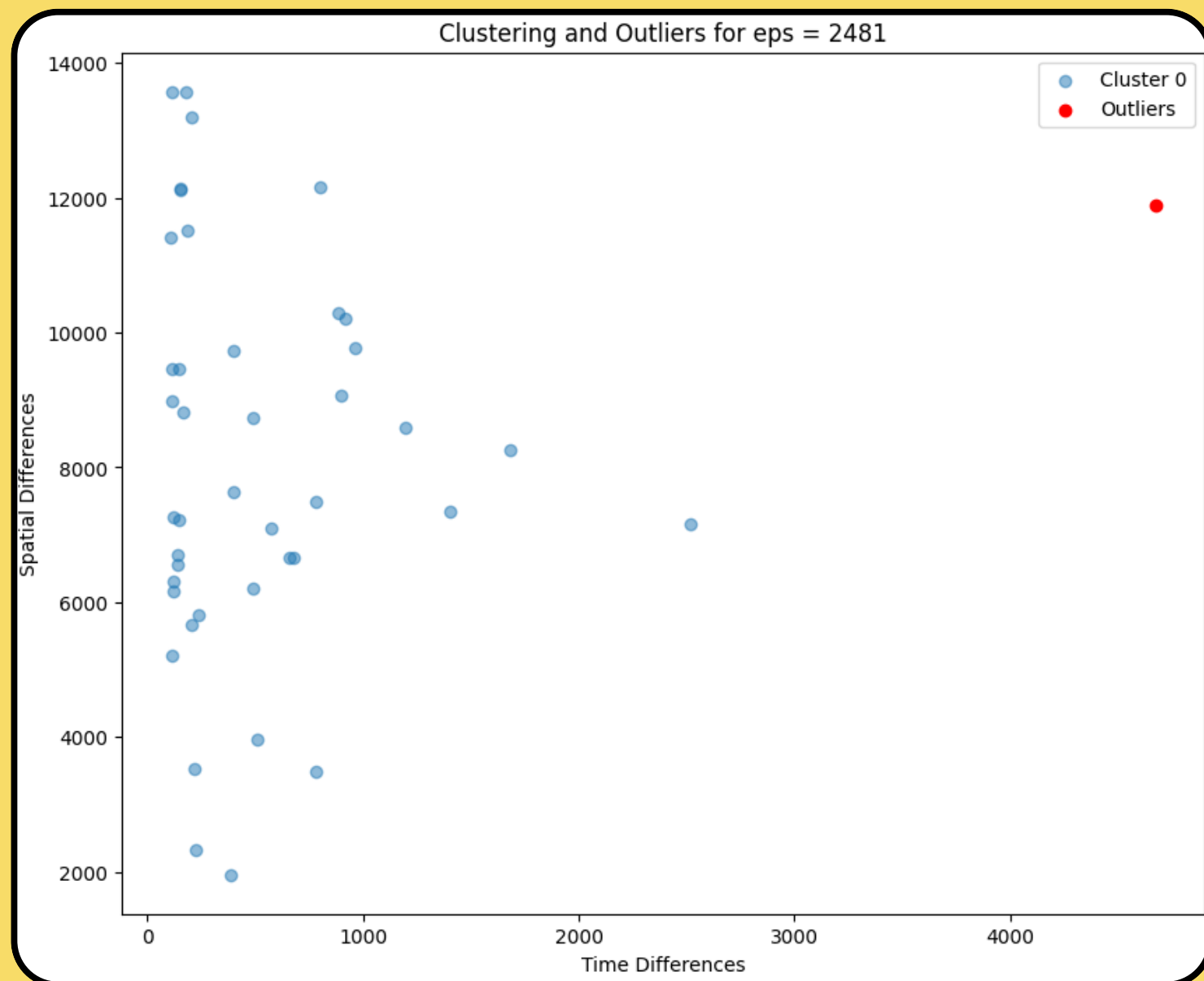
- **epsilon (ϵ)**, the maximum distance between two points to consider them part of the same cluster;
- **minPoints**, the minimum number of points required to form a cluster.

- Experiment varying **epsilon values** to find the optimal configuration

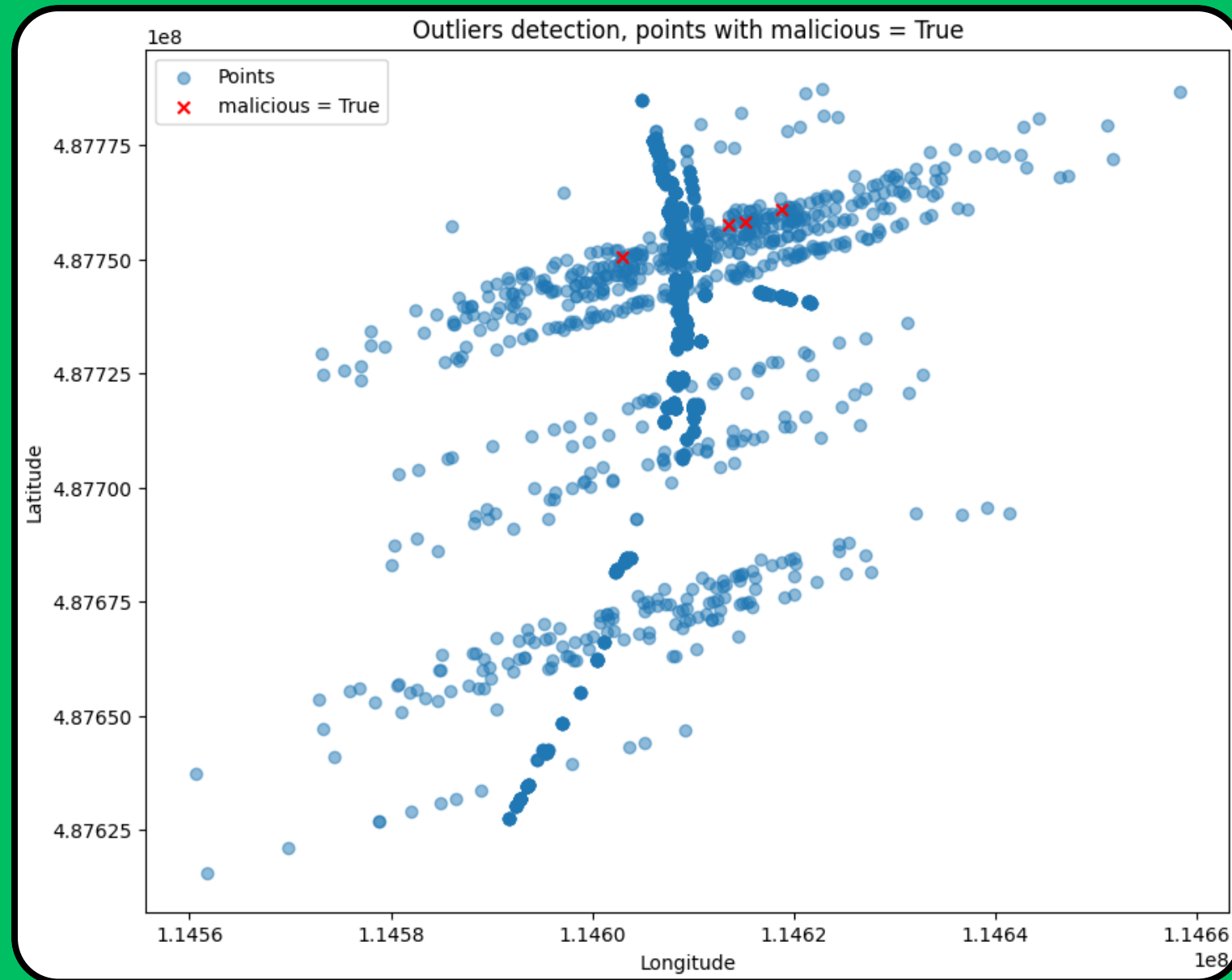
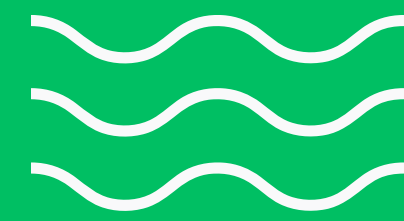
- The best epsilon is the fewest possible outliers, the data we are working with is clean!



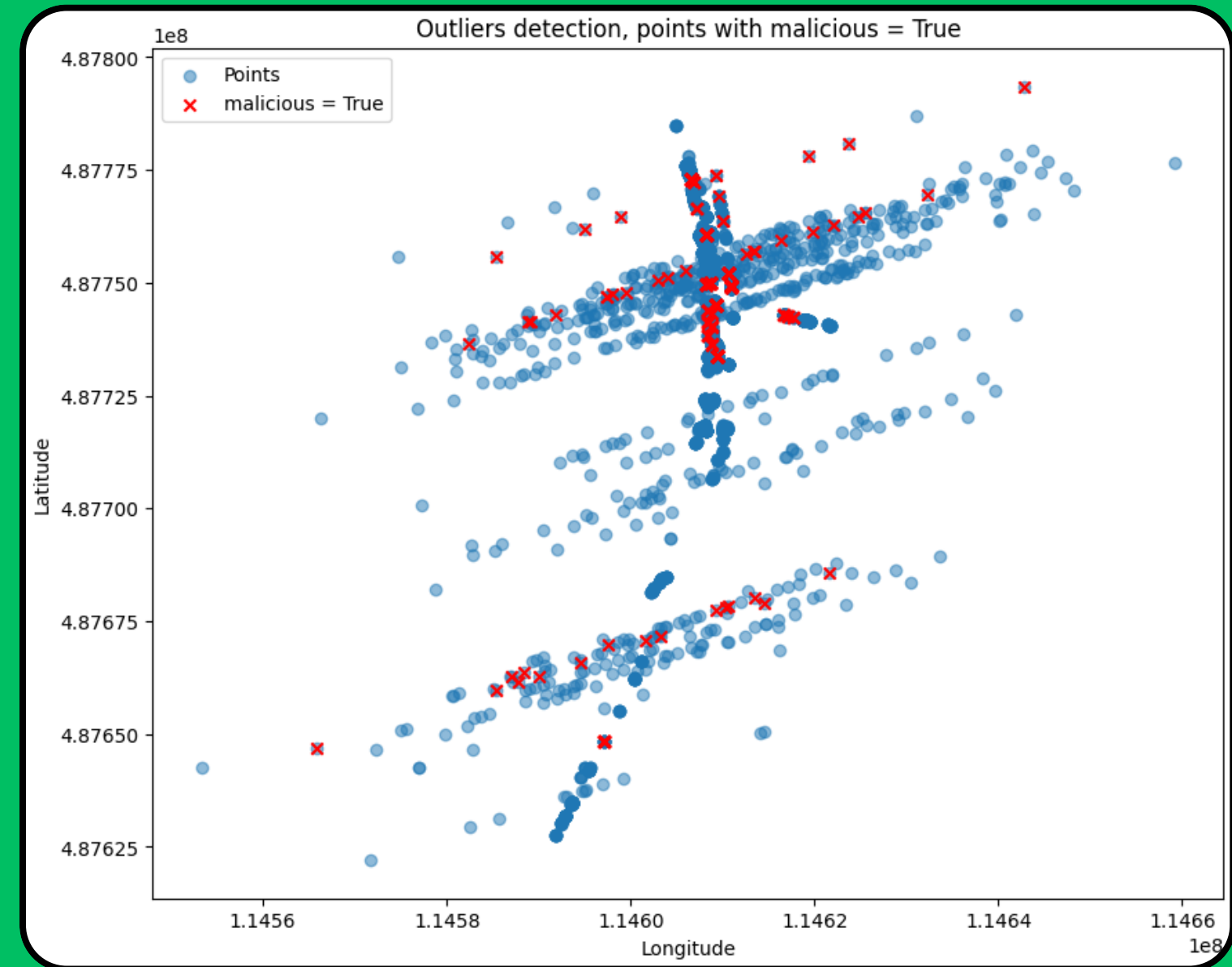
OUTLIER



LABELING - EXPERIMENTS_x

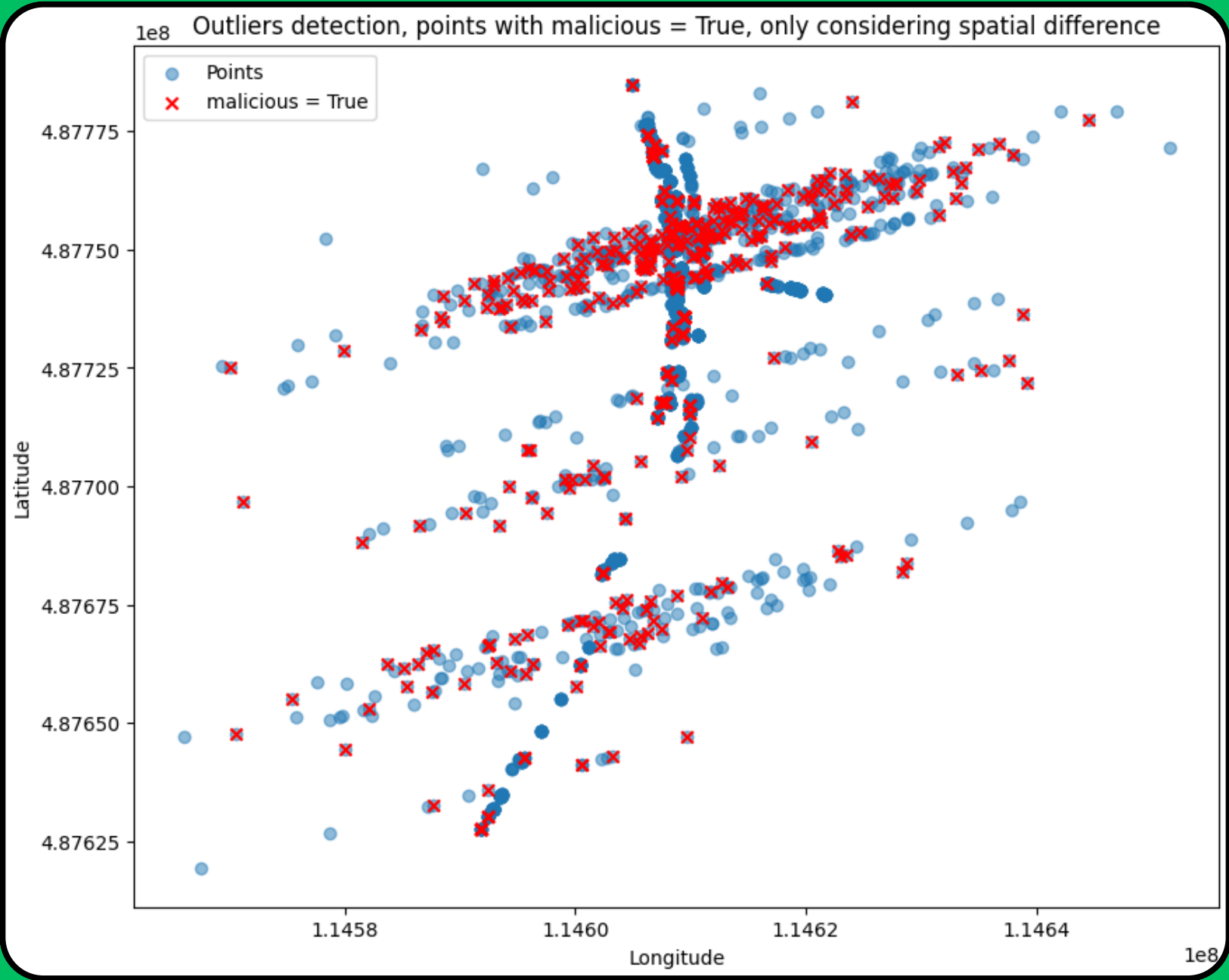


Outliers detection for epsilon = 2481 with the dirty dataset



Outliers detection for epsilon = 1281 with the dirty dataset

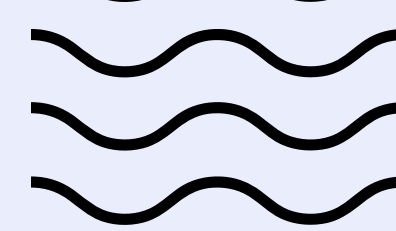
LABELING - EXPERIMENTS_x



TP	FN	FP	TN
38.66%	61.34%	14.44%	85.56%

Outliers detection, points with malicious = True, only considering spatial difference

CONCLUSIONS



Challenges in outlier identification, especially for **points near trajectories**.

Substantial improvement observed by removing the **temporal component**.

Best-performing experiment captured only 38.66% outliers.

Consider utilizing a **more extensive dataset** with diverse event types to enhance the robustness and generalizability of the model.

Explore more advanced **machine learning and deep learning** techniques.



Thanks for the attention!

