

DEPARTMENT OF INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Bachelor's Thesis in Informatics

**Crowdsourcing mobility data with privacy
preservation through decentralized
collection and analysis**

Simon van Endern

DEPARTMENT OF INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Bachelor's Thesis in Informatics

**Crowdsourcing mobility data with privacy
preservation through decentralized
collection and analysis**

**Crowdsourcing von Mobilitätsdaten ohne
Einschränkung der Privatsphäre durch
dezentrales Sammeln und Analysieren**

| | |
|------------------|-------------------------|
| Author: | Simon van Endern |
| Supervisor: | Prof. Dr.-Ing. Jörg Ott |
| Advisor: | Trinh Viet Doan |
| Submission Date: | 30.06.2019 |

I confirm that this bachelor's thesis in informatics is my own work and I have documented all sources and material used.

Munich, 30.06.2019

Simon van Endern

Acknowledgments

Abstract

We propose a method to publish location data without raising privacy concerns.

As still this data could be useful for many stakeholders, we will investigate how on the one hand aggregated data can be published without imposing any privacy risk to the owners of the data and on the other hand develop a prototype of a mobile application through which this location data is aggregated in a decentralized manner so that the raw user data never leaves the users' device.

Contents

| | |
|--|------------|
| Acknowledgments | iii |
| Abstract | iv |
| 1. Introduction | 1 |
| 1.1. Motivation | 1 |
| 1.2. Research Questions | 2 |
| 1.3. Contributions | 3 |
| 2. Related Work | 4 |
| 2.1. Classification of location data usage according to acceptable delay . . . | 4 |
| 2.2. Privacy problems arising from location data | 4 |
| 2.2.1. Risk of privacy intrusion through theft of central databases . . . | 4 |
| 2.2.2. Inference attacks on published anonymized data | 5 |
| 2.3. Countermeasures to prevent inference attacks | 5 |
| 2.4. Limitations of countermeasures | 6 |
| 3. Methodology | 7 |
| 3.1. Aggregation schemes | 9 |
| 3.2. Limiting the spatial area of an aggregation request | 10 |
| 4. Design and Implementation | 12 |
| 4.1. Technology Stack | 12 |
| 4.2. API | 13 |
| 4.2.1. Data Aggregation Design | 15 |
| 4.3. Android Application | 16 |
| 4.3.1. Separation of concerns regarding Data Collection and Aggregation | 18 |
| 4.3.2. Data Collection | 18 |
| 4.3.3. Local Data Aggregation | 19 |
| 4.3.4. Serving Aggregation Requests | 20 |
| 4.4. Server Design and Implementation | 22 |
| 4.4.1. Data Model | 22 |
| 4.4.2. request-chain | 23 |

| | |
|--|-----------|
| 5. Performance and Evaluation | 25 |
| 5.1. Deployment | 25 |
| 5.2. Data Consumption | 25 |
| 5.3. Aggregation Results | 26 |
| 5.3.1. Validity of results | 27 |
| 5.3.2. Implications | 30 |
| 5.3.3. Inference on aggregated data | 30 |
| 5.3.4. Trajectories | 31 |
| 5.4. Further scenarios e.g. Scenario Traffic Jam | 32 |
| 6. Conclusion and Discussion | 35 |
| 6.1. Limitations | 35 |
| 6.2. Future Work | 35 |
| 6.3. Evaluation | 38 |
| 6.4. Limitations | 39 |
| 6.5. Reproducability Considerations | 39 |
| Appendix A. Data Usage Screenshots | 40 |
| List of Figures | 51 |
| List of Tables | 52 |
| Bibliography | 53 |

1. Introduction

1.1. Motivation

"Data is the new oil" [5, 7] is a quote many people agree with. More and more businesses are based not on specific production capacities but on data and especially the ability to process it and the exclusive ownership over it. The success and monopoly of companies like Google, Facebook and Amazon can be attributed to this exclusive ownership at least to a reasonable extend.

While patents that used to power companies' success provide a balance through granting exclusive rights having to make the knowledge public, many companies e.g. Coca-Cola have decided successfully not to file for a patent and thus not having to reveal their knowledge [16]. If that approach is not compromised, it guarantees both, non-disclosure and exclusive rights. Similarly, the non-disclosure of huge data sets collected by Facebook, Google and Amazon circumvent the balance intended by patents. But this unavailability of huge amounts of data to the public is an impediment of innovation and increased growth. For example, cities would benefit from aggregated location data in order to optimize traffic scheduling as also highlighted by Hoh et al. [9]. Nevertheless, the publication of raw data sets is impossible because it severely intrudes the privacy of the owners of the data.

So, even if companies would agree on a publication, a problem arises. There is a conflict between preserving user privacy and publishing user data. But in fact, users' privacy is already compromised even without publication of their data. Already the mere existence of central data sets pose a privacy risk to users, because security issues might allow for theft and unwanted publication of these data. An example is the theft of 14 million users' data from facebook [8].

Some governments and other institutions already publish some of their data sets after anonymizing them and there are crowdsourcing and open source approaches to make data available to everybody [CITE!!!]. Nevertheless, the applied anonymization is often not sufficient or at least critical if the resulting data set should still be useful. Research shows that inferences can be drawn from the published data sets that violate the respective users' privacy [4, 24]. Also, most research is based on datasets that span only over a few days or weeks. When collecting location data over longer periods, the chance of successful inference attacks increases in all cases with more data available.

So, in addition to the privacy risk imposed on the user by the existence of a central data set, publishing anonymized data poses another risk to users' privacy.

1.2. Research Questions

In order to further investigate the trade-off between publication of location data and preservation of user privacy, we pose the following research questions:

RQ1: Is aggregation of location data possible without raw data being accessible to anybody but the owner?

Previous work has shown that in order to overcome the conflict between privacy intrusion and (public) data availability, a solution is needed that works without storing raw data in a central data set. This solution should eliminate the risk of leaking raw user data through theft from a centralized database. It should furthermore eliminate the risk of inference attacks on published believed-to-be-anonymized raw data if only aggregated data is stored on a central server. We investigate whether there is a solution and also what limitations this solution has, respectively which assumptions must hold for the solution to be valid.

RQ2: What types of aggregations can be published?

Clearly, a mean value exposes less details about the underlying raw data than the whole distribution from which e.g. the median value and other percentiles can be computed. We will investigate how many details, respectively which statistical values can be published without user privacy being compromised.

RQ3: What is the risk of inference attacks on aggregated data due to overlappings in the covered timespan?

Previous work has covered inference attacks on anonymized raw data sets. They have shown that inference attacks can be based on overlapping timespans or the same data set being published using different anonymization techniques. We will investigate, whether there is also a risk of inference attacks based on overlapping timespans in aggregated data and examine the limitations of publishing aggregated data.

1.3. Contributions

While most research follows the approach of degrading data in order to provide anonymity, we do not first collect the whole data set and then reduce it to a data set meeting privacy-constraints but we start from the bottom up. First by performing aggregation in a decentralized manner and second by verifying that an aggregation does not violate users' privacy prior to using this aggregation in a production environment. Step by step, the set of possible aggregations to be published can be extended after verifying in a test environment, that the respective aggregation does not intrude the users' privacy. The data stemming from those aggregations will then be available to the public.

We investigate the possibility of storing raw location data only decentralized on the collecting devices. On a central database available to the public, only aggregated data is stored, thus the privacy risk arising from a central database containing the overall raw data set is eliminated. Similar to previous research, we assume the central server to be trusted and also the clients to be trusted. As the aggregation process happens decentrally, the central server will never hold any other than aggregated data and never know about individual raw data.

Our research is organized as follows: First we review related work in the areas of location privacy and anonymization techniques. Chapter 3 outlines the general architecture of our approach using decentralized data analysis and chapter 4 documents the implemented setup and explains design decisions. We present the results of field-testing our setup in 5, chapter 6 discusses limitations as e.g. a trusted server, reproducibility considerations and outlines possible future work.

2. Related Work

2.1. Classification of location data usage according to acceptable delay

In order to review existing approaches and research, we classify location aware services by the acceptable delay of the location information being available. Similar to the classification made by [9], we define three categories:

1. Almost no delay tolerance: e.g. an application showing a pop-up about a nearby venue e.g. a coffee shop when a pedestrian passes by.
2. Some delay tolerance e.g. one minute: An application e.g. Google Maps derives the information of congested traffic from devices reporting their GPS data which show lower than usually on the respective road. As congestions worth reporting last longer than one minute, some delay in the device's information reaching the server is acceptable.
3. Significant delay tolerance of hours, days or even weeks: For historical and statistical use of location data e.g. to find out about popular visiting times of a venue, almost any delay is acceptable.

Some research deals with case one where almost no delay tolerance is acceptable [1, 2]. [15] proposes a solution where not the exact location is sent to a server but the rough region of the user. The server then sends a list of all possible matches like e.g. venues in this area to the client. Locally, this list is then matched with the exact location in order to fulfill the aim of the respective application. Most research though investigates user's privacy in case 3 where the delay of the data being available for processing is not an issue [14, 4, 6, 24]. In the following, we will concentrate on this case as well.

2.2. Privacy problems arising from location data

2.2.1. Risk of privacy intrusion through theft of central databases

Centralized databases containing raw location data expose users to a privacy risk (through theft) [20, 10]. [12] proposes the use of P2P over WIFI and Bluetooth to

decrease the need of central instances. A decentralized analysis approach and its implications for data privacy is also investigated by [20] as an alternative to cloud-based IoT. [13] proposes an approach in which raw data is hidden from the central instance but still aggregated data can be obtained by using encryption methods. Raw data is encrypted using a modified approach of public-private key-pair cryptography in which the sum of two encrypted messages can be decrypted to the sum of the encrypted messages. Furthermore, only a number of messages above a certain threshold can be encrypted this way using the different encryption shares. Another approach by [10] also uses encryption in combination with a middleware. The server storing the data and the participant (e.g. a vehicle) share a symmetric key which is stored securely in the vehicle. The middleware ensures authentication of the participant, and forwards the location to the central server without giving away the vehicle's identity. Nevertheless, several research as e.g. [14, 24, 4] have shown that it is easy to infer identity from such datasets.

2.2.2. Inference attacks on published anonymized data

Research has shown, that from a location data set that is pseudonymous, i.e. the identifiers have been stripped off or the data-set has been anonymized in another way, it is possible to infer the home location of single users through so-called inference attacks [14, 4, 6, 10, 24] and also the work location with a slightly lower probability [4, 6]. The same problem arises when using data collected through crowdsourcing [13]. These home locations or home-work location pairs can then be used to look up the corresponding user's identity e.g. by combining it with publicly available information. One possibility is to reverse code GPS coordinates to addresses and then e.g. search for entries in telephone books to infer the user's identity from its home location [14, 6, 10]. Especially in suburban areas this is quite successful as usually one house can be mapped to only one person or family. This identity can then be linked to other sensitive data, e.g. locations visited by the identified user. The same problem also arises in the area of IoT [20, 10].

2.3. Countermeasures to prevent inference attacks

In general, the problem can be solved by providing k -anonymity for the data-set [22]. A data-set is k -anonymous, if querying for an identifier in this data-set always returns a result set of at least k entries. In order to achieve this, several approaches have been investigated. [14] propose spatial cloaking. K -anonymity is achieved by dropping data points or perturbing them or dropping all data points around a random point around the home location. Also obfuscating locational data close to a home location

and mapping GPS points to the next street crossing is possible to increase anonymity. More sophisticated approaches as e.g. [11] focus on making it less likely to identify GPS points of one trajectory being subsequent and belonging to the same user.

2.4. Limitations of countermeasures

Usually there is a trade-off between the level of anonymity and the usefulness of the data. When k -anonymity is guaranteed, often the resulting data-set becomes useless because the data quality is not sufficient anymore [14, 4, 17, 22, 21]. On the one hand, when the data-set is tried to be kept useful, data suppression algorithms have only limited success and can only reduce, but not eliminate the risk [10]. [11, 1, 10] find that anonymization techniques might score well in densely populated areas or areas with high traffic but poorly in sparsely populated areas especially where a single address can be mapped to a single person or family. Also the applied techniques might not achieve the expected results for individuals who's work and home location are further away than average [6]. Furthermore, taking other sources and databases into account, k -anonymity might be compromised due to quasi-identifiers [21] e.g. a combination of attributes that do not identify an individual but allow linking two different data-sets and by that creating new identifiers.

3. Methodology

In order to avoid having a central raw dataset and to eliminate the risk of inference attacks on anonymized data, we propose a framework which collects and locally aggregates location data on end user devices (smartphones) through an application designed for this purpose. The raw data will stay on each device and will only be used to serve aggregation requests send from the central server. The aggregation requests have to be defined upfront. Apart from GPS data, we also take other movement data into account - the number of steps and the current activity e.g. walking. These can be used to enrich the pure GPS data. Hence, an example for an aggregation is the determination of the average number of steps per day accross all users participating in the respective aggregation. An exemplary aggregation request sent to a smartphone is depicted in figure 3.1. It specifies the timespan which should be covered by the aggregation, the type of aggregation and the current data combined in the ongoing aggregation.

```
1 {  
2     "start": "2019-05-30",  
3     "end": "2019-06-02",  
4     "type": "steps",  
5     "n" :3,  
6     "value": 2000  
7 }
```

Figure 3.1.: Incoming aggregation request.

```
1 {  
2     "n" :4,  
3     "value": 2500  
4 }
```

Figure 3.2.: Outgoing aggregation response.

The response send back to the server after processing the request only needs to contain the data itself, as the type and timespan can be inferred by the server from the request identifier. Figure 3.2 shows an exemplary response to the server that matches the request in figure 3.1.

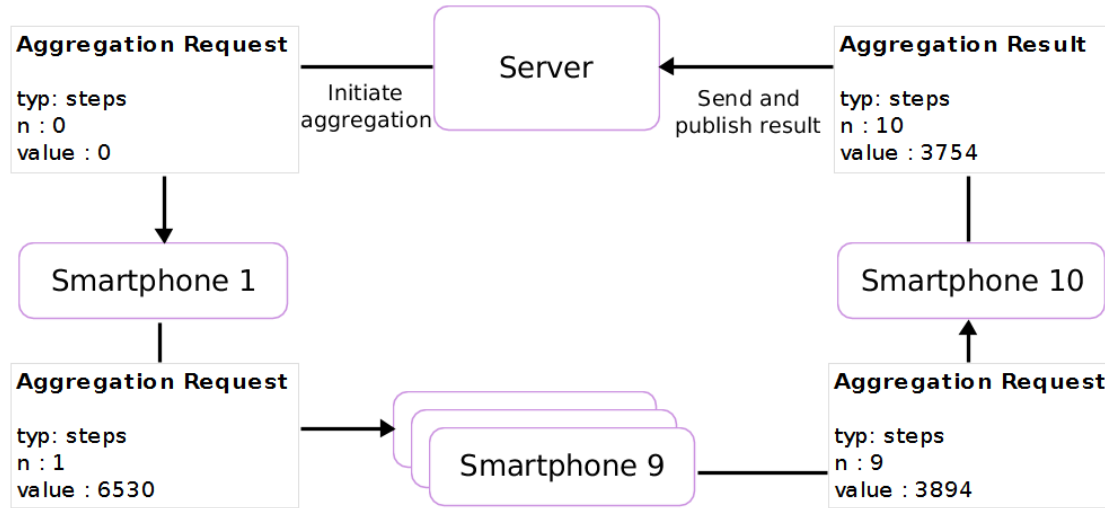


Figure 3.3.: Decentral unencrypted aggregation process via P2P.

Figure 3.3 depicts the overall process of such a decentral aggregation request using an example of 10 participating devices. In order to protect the user's privacy and completely shield the raw data from the server, it would be necessary to pass the request via P2P from one device to another until the last device finally sends the results to the server. On the one hand, P2P on mobile phones though is still hardly possible, if the devices are not locally close to each other. Otherwise, P2P could be established using WIFI and Bluetooth connections [12]. On the other hand, if the server is used to pass an aggregation request from one device to another, it could read the data and compute the respective user's input from the difference. We propose to use encryption in order to hinder the server from reading the data and assume the server to be trusted in the first place. Chapter 6 discusses this limitation. The modified aggregation process using encryption is depicted in figure 3.4. On installation of the application on a smartphone, a public-private key-pair is generated and every installed application registers at the server with this public key. The corresponding private key is stored only locally. On start of an aggregation request, not only the first user but also the subsequent user in the chain who should deal with the aggregation request is determined and the public key of this user is passed along with the aggregation request.

When one end user device needs to send the processed aggregation request to the next phone, it encrypts the data using the provided public key of the next user leveraging the benefits of synchronous keys using the standard hybrid encryption approach¹. This way, the next phone in the aggregation chain will be able to decrypt the request and process the data while the server is unable to read the data until the aggregation request is finally send in plain text for publishing to the server.

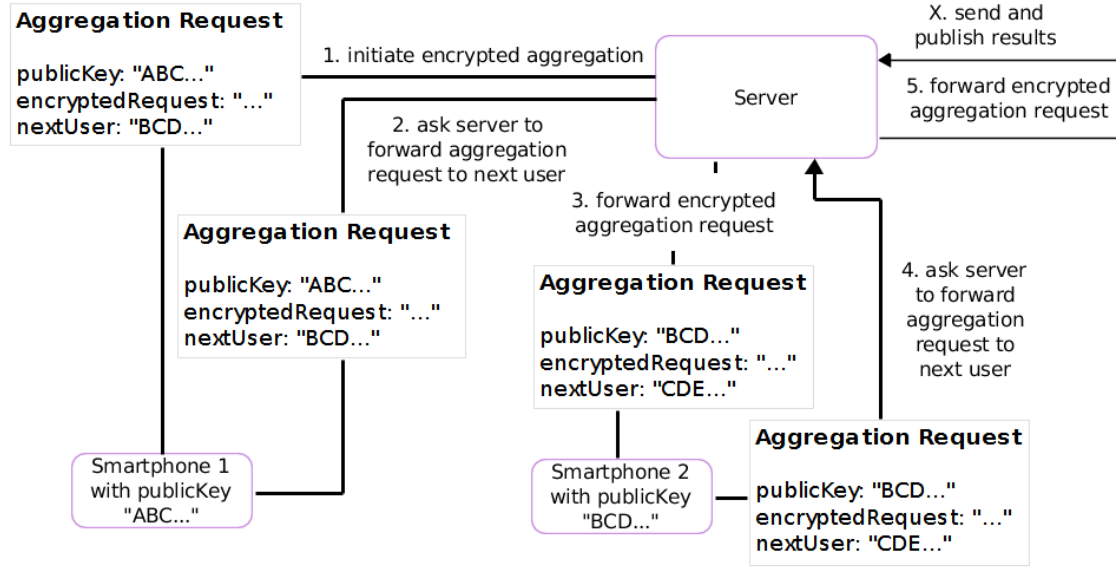


Figure 3.4.: Decentral encrypted aggregation process.

The aggregated results are planned to be available through a public route on the server. Nevertheless, we restrict access in our test run to our research team in order to protect the research participants' privacy in case there is a privacy risk we have not thought of. The aggregation requests themselves are initiated by our research team but can be initiated automatically on a regular basis in the future. The next section will detail the planned aggregations.

3.1. Aggregation schemes

Two types of aggregation requests are of special interest for our research. First, the aggregation of mean values and second, the aggregation of more advanced statistical values such as median or distribution. The latter includes the possibility to calculate

¹In hybrid encryption as used in SSL, the message itself is encrypted with a synchronous key while this key itself is encrypted using the public key

the former and requires the simultaneous availability of each single user's response. The following list is an excerpt of aggregations that would be of interest and can be computed based on GPS location data in combination with the number of steps and the current activity of the respective person:

1. The average number of steps walked across all users participating in the request. (e.g. to calculate how many people reach 10.000 steps per day².)
2. The average time spent walking, running, in a vehicle or on a bicycle.
3. The share of people who combine using a bicycle with using a vehicle such as public transport or a car within one trajectory.
4. The time spent at work and the time people need to travel to work.
5. Locations where many participants spend a significant amount of time on a certain day. (A posteriori event recognition)
6. The percentage of their overall travelling time that people spend on their bike, car, etc.
7. The average speed on roads.

For all aggregations, always both, the mean value and if possible, a complete list of single users' mean values in order to compute other statistical figures, are of interest.

3.2. Limiting the spatial area of an aggregation request

The aggregation requests outlined in the former subsection only provide useful data if the area of the aggregation can be limited e.g. to the scope of a city. Otherwise the resulting data would not allow for comparison and the scope of each aggregation would either be the whole user base which especially in the case of listing values would result in a huge amount of data passed around. Or, the limited number of participants in each aggregation would not be locally close to each other which would not result in useful data and also make it impossible to do aggregations as the average speed on roads, where you need a certain number of participants providing data about this area. In order to limit the area of aggregation and avoid sending the aggregation request to each user and leaving him with determining whether he is inside the area and participates in the request, the initiator of those requests has to know the location of

²It has to be evaluated, which percentage of steps are registered because the phone will not always be on the person.

users, respectively the location where the user is active the most. We do not see this as a violation of the user's privacy due to the following reason: The exact location of the user e.g. the home or work location is not of interest at all. Rather the area for which the user can provide data is of interest. We propose to cluster location areas in a hierarchical structure similar to [15] and determine the granularity of the location published to the server as follows:

1. Each user sends the most coarse locational area possible to the server - e.g. the continent.
2. If more than the required anonymity threshold of e.g. 10 active users are already registered with this area at the server, the server not only links this location to the user but also requests the user to send a less coarse location.
3. The user step by step sends a less coarse location e.g. country, district, etc. until the server denies to link the user to the area because not enough active users have registered with the location on this granularity level. The server nevertheless increases a counter of users who requested access to this area. Once the counter exceeds a certain number, the server sends an aggregation request to participants in the more coarse area encompassing the requested less coarse area to determine the number of active users. If this number is above the required threshold, the granularity level for this location is made available and users can now register with this area and aggregation requests targeting this area can be send.
4. The fulfillment of the threshold has to be checked on a regular basis in order to close areas once the user base sinks below the anonymity threshold.

In a more advanced setting, it should be possible to register with more than one area on the same level to avoid that a user e.g. living close to a city provides only data about the city or the bordering area but not both. Nevertheless, this should be limited to areas bordering each other to avoid user identification similar as in [6] and it has to be investigated whether this indeed does not pose a privacy risk to the user or whether also the combination of areas needs to meet a certain anonymity threshold.

4. Design and Implementation

4.1. Technology Stack

In order to implement the architecture proposed in chapter 3 we choose Android as end user device platform. Android has the highest market share among mobile devices [18] and offers a healthy ecosystem of frameworks and libraries that simplify development. Furthermore we opted for an implementation in Kotlin instead of Java to reduce boilerplate code and improve readability.

As server side technology we choose node.js in conjunction with a mongoDB object store database. A NoSQL database like mongoDB provides flexibility and easy adaption of data schemes without much overhead and thus perfectly fits our prototyping purpose. We choose node.js out of the same reasons. In contrast to a statically typed language, javascript provides more flexibility and ease of change. Furthermore, node.js is often used in combination with mongoDB and provides seamless integration.

The aggregation results are made available to the public via a dedicated route of the server but might also be made available directly through granting access to the respective collection¹ of the database.

The following sections describe the implementation of the Android application and the server and explain design decisions. We start with an outline of the API shared by both, the Android application and the server.

¹A collection in an object store is the equivalent to a table in an SQL database.

4.2. API

The API is designed as a REST API based on JSON data exchange format and consists of the following endpoints:

- POST /user HTTP/1.1
for creating a new user (see figure 4.1 and 4.2).
- GET /requests?pk=... HTTP/1.1
for retrieving aggregation requests for a user (see figure 4.3).
- POST /forward HTTP/1.1
for sending a processed aggregation request back to the server (see figure 4.4).
- GET /aggregations HTTP/1.1
for retrieving all aggregation results (see figure 4.5).
- POST /admin/sampleRequest HTTP/1.1
for initiating a new aggregation (see figure 4.6).

```
1 {  
2   "publicKey": "-----BEGIN PUBLIC KEY-----\nMIIBIjANBgkqhkiG..."  
3 }
```

Figure 4.1.: Body of an HTTP request for creating a new user.

```
1 {  
2   "publicKey": "-----BEGIN PUBLIC KEY-----\nMIIBIjANBgkqhkiG...",  
3   "password": "LSFDfzduSFSozfwf"  
4 }
```

Figure 4.2.: Body of an HTTP response upon successful creation of a new user.

The routes used for interaction with the application (except the one for creating a new user) are secured and can only be accessed if the user can be authenticated successfully. The routes for initiating new aggregations and retrieving all results require administrator authentication (see figure 4.6). The route for the results is designed to be

```

1 [{
2   "publicKey": "-----BEGIN PUBLIC KEY-----\nMIIBIjANBgkqhkiG...",
3   "serverId": "acfbdxxfceola0dkfeecIf",
4   "nextuser" : "-----BEGIN PUBLIC KEY-----\nMIIBIjANBgkqhBms...",
5   "encryptionKey": "HXporXstqrgkfsCpezlqZcpLb...",
6   "iv": "pMxQznHqbzeZuzVLfiFHYQ==",
7   "encryptedRequest": "xtgPteNAeAKEhfPXQeZt..."
8 }]

```

Figure 4.3.: Body of an HTTP response listing all aggregation requests for the specified user.

```

1 {
2   "publicKey": "-----BEGIN PUBLIC KEY-----\nMIIBIjANBgkqhkiG...",
3   "password": "LSFDfzduSFSozfwf",
4   "serverId": "acfbdxxfceola0dkfeecIf",
5   "encryptionKey": "HXporXstqrgkfsCpezlqZcpLb...",
6   "iv": "pMxQznHqbzeZuzVLfiFHYQ==",
7   "encryptedRequest": "xtgPteNAeAKEhfPXQeZt..."
8 }

```

Figure 4.4.: Body of an HTTP request sending the processed aggregation request back to the server.

available to the public without authentication. Nevertheless we restrict access in this test setup to fully protect users' privacy. On creation of a new user, the server response contains a password randomly generated for the user to be used for authentication in all future server communication as e.g. in figure 4.4. All aggregations waiting to be processed by the user and also the respective responses (see figure 4.4 and 4.3) contain the public key of the user to identify which user should handle the request. Furthermore, the request contains a field *encryptionKey* which is a synchronous key, encrypted with the users public key, that has been used in combination with the initialization vector *iv* to encrypt the actual aggregation request described in chapter 3 send along in *encryptedRequest* in it's encrypted version. The *serverId* is the identifier of the request used by the server's database.

```
1 [{
2   "timestamp": 1559216514521,
3   "startedAt": 1559216514685,
4   "start": "2019-05-28",
5   "end": "2019-05-31",
6   "type": "steps",
7   "n": 6,
8   "value": 4271,
9   "valueList": []
10 },
11 {
12   "timestamp": 1559216514521,
13   "startedAt": 1559216514685,
14   "start": "2019-05-28",
15   "end": "2019-05-31",
16   "type": "stepsListing",
17   "n": 5,
18   "value": 0,
19   "valueList": [1661,1246,3195,7714,7842]
20 }]
```

Figure 4.5.: Body of an HTTP response listing all completed aggregations.

```
1 [{
2   "password": "adminPassword",
3   "start": "2019-05-31",
4   "end": "2019-06-05",
5   "type": "activity_2",
6 }]
```

Figure 4.6.: Body of an HTTP request for initiating a new aggregation.

4.2.1. Data Aggregation Design

All aggregations proposed in section 3.1 can be implemented using only 4 fields: An Integer n specifying the current number of participants, a Floating Point Number $value$ e.g. containing the current mean value of the aggregation, a list $valueList$ of Floating Point Numbers that can be reused for different purposes accross requests and a String

type that specifies the request type and by that the used scheme, i.e. how the other fields are populated e.g. the encoding of the *valueList*. In addition, all aggregation requests have a start day and an end day. Days' time is always treated as 00:00 o'clock. Thus, the timespan for an aggregation is always a multiple of 24 hours. Figure 4.5 shows two examples of completed aggregations which highlight the reusability of the 4 fields. They furthermore contain a *timestamp* when the aggregation has been completed and when it had been started (*startedAt*).

Due to limited scope of this thesis, only the following out of the aggregation requests defined in section 3.1 were implemented:

1. Computing the average number of steps walked across all users participating in the request.
2. Computing the average time spent walking, running, in a vehicle or on a bicycle².
3. Collecting a list of the average number of steps walked by each participant during the timespan.
4. Collecting a list of all trajectories registered by the users' phones.

Nevertheless, implementation of other aggregations is straight forward using the provided setup. All aggregations work as a prove of concept of the proposed setup and hence serve answering research question 1. Aggregation 1, 2 and 3 furthermore help to investigate which types of aggregations are possible to be published and allow for an investigation whether overlapping timeperiods allow for inference attacks. Thus they help to answer research question 2 and 3. Especially aggregation 3 in comparison with aggregation 1 and 2 allow for answering research question 2. Aggregation 4 imposes a privacy risk as discussed in [24, 4, 14] on the user and was only implemented in order to get an overview of the data quality of our setup and validate the feasibility of the other aggregation requests proposed in section 3.1.

4.3. Android Application

The Android application targets Android Oreo (API level 27) and requires a minimum API level of 19. Approximately 96.8% of devices run on this or a higher version of Android [3] which allows our application to be installed on the majority of Android devices. Furthermore, the application leverages Google Play Services to obtain GPS and activity data. Without Google Play Services installed, the application will not

²See <https://developers.google.com/android/reference/com/google/android/gms/location/DetectedActivity>

work. In order to ease future adaptability, we chose to use the Dagger2³ framework for dependency injection in order to decouple classes as far as possible. Further use of frameworks and libraries will be explained in the following respective sections.

The application is aimed to collect GPS data, detect the user's current activity and count the user's steps. The data collection process happens in the background without any user interaction required. Aggregation requests to aggregate data accross devices are also served without the need of any user interaction. The android application can be grouped into three main parts of loosely coupled modules:

- A module responsible for collecting and saving raw data.
- A module responsible for locally aggregating raw data.
- A module responsible for communicating with the server and handling aggregation requests.

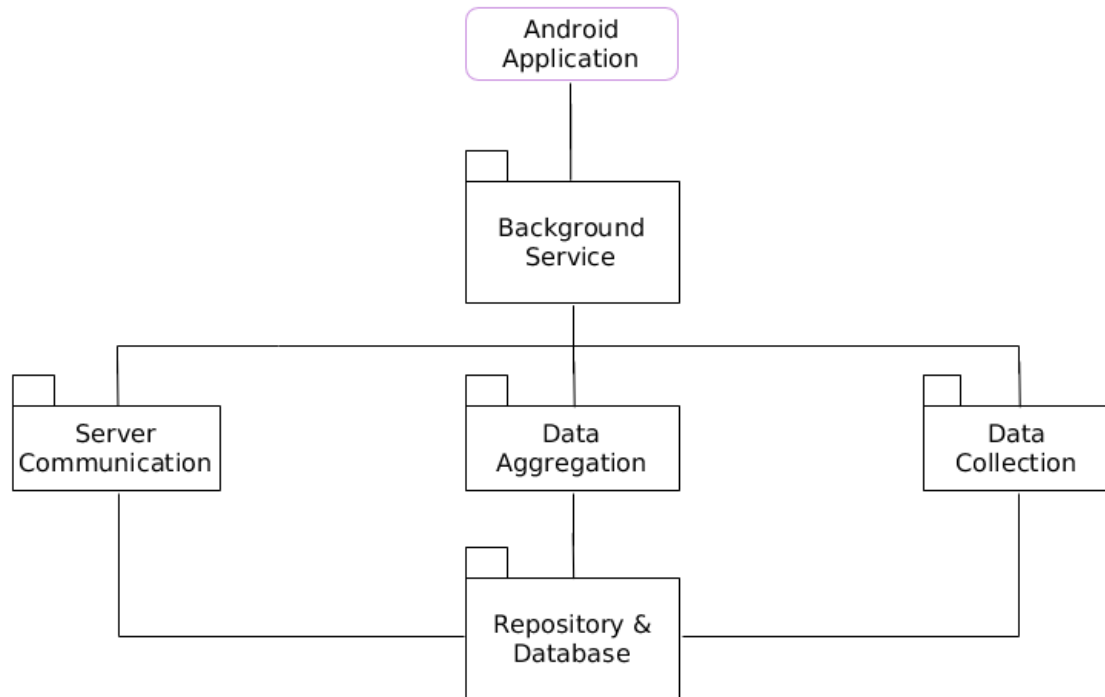


Figure 4.7.: Android architecture

The control flow as depicted in figure 4.7 is as follows: The Android application has only one Main Activity in order to ask the user to allow location access and start the

³See <https://dagger.dev/>

background services. Apart from that, this Activity does not serve any specific purpose. The local aggregation as well as the polling of new requests from the server happens in the background on a 15 minute interval and is initiated by the background service. The Android Workmanager then controls this periodic work without any user interaction being required.

For the App in order to have maximum possibilities collecting especially GPS data and preventing the Android operating system from shutting down when not interacted with by the user (which is usually never the case), a non-dismissible status notification is displayed at all time⁴. Also, the application is registered to be automatically restarted upon boot⁵ and also when the application is closed by the user (e.g. via the task manager) so that once installed, no further user interaction is necessary. Furthermore, the application, respectively each module is heavily unit tested in order to guarantee functionality and facilitate further development by other research teams. Unit and integration tests are based on AndroidJUnit4⁶ and the espresso⁷ framework.

4.3.1. Separation of concerns regarding Data Collection and Aggregation

We choose to separate the aggregation and collection of location data in order to decouple the modules and provide the possibility to extend the model of aggregated data in the future without the need to change the raw data model. Vice versa the data collection process can be modified without impacting the aggregation process.

4.3.2. Data Collection

We use the Android Room Persistence library⁸ to locally store data. The library provides a layer over the standard LiteSQL database commonly used in many Android applications. We collect three types of data:

- Steps: If available, the phone's internal step sensor provides updates on a regular basis. The step sensor always informs about the total number of steps since the

⁴Compare to the non-dismissible status notification displayed by Google Maps when the navigation system is active.

⁵From Android 6.0 on (API level 23), apps' behaviour is restricted by the operating system in order to reduce battery consumption. E.g. all apps are automatically managed by the battery manager which restricts background launches. The user has to switch this option to manual management in order to allow the app to function in the way it was designed. See <https://developer.android.com/guide/background/> and also <https://developer.android.com/training/monitoring-device-state/doze-standby>.

⁶See <https://developer.android.com/reference/android/support/test/runner/AndroidJUnit4>.

⁷See <https://developer.android.com/training/testing/espresso>.

⁸See <https://developer.android.com/topic/libraries/architecture/room>.

last reboot. Upon each time we receive data from the step sensor, this data is stored directly in the *step_counter_table*.

- User's activities: The Google Play Services activity recognition API leverages different data and sensors available on the phone in order to inform about the most probable current activity of the user as one of *still*, *walking*, *running*, *in a vehicle*, *on bicycle*⁹. Whenever there is a change detected, two events are fired - one for exiting the former and one for entering the current activity. The events might not be dispatched instantly but contain the timestamp of the exact occurrence. Upon each received event, this data is stored directly in the *activity_transition_table*.
- GPS positions: GPS data is retrieved through the *FusedLocationProviderClient* which leverages cellphone-tower and WIFI data apart from GPS to determine the position. In order to limit battery consumption, GPS data is only requested every minute if the device's detected activity is *still*¹⁰. If the current detected activity is *walking*, the interval is set to 5 seconds and in any other state, the interval is set to every second. The data is stored in the linked tables *gps_data_table* and *gps_location_table*. We choose to separate the GPS point itself from the timestamp having in mind that future aggregations might need or leverage the separation of spatial data and time and more than one event might be attached to the same GPS point.

4.3.3. Local Data Aggregation

From the received values of the step counter since last reboot saved in *step_counter_table* the daily steps are computed and stored in *steps_table*. The exit and enter events received via the activity recognition framework and stored in the *activity_transition_table* are matched in order to compute activities with start and duration. Those activities are then saved in the *activity_table*. GPS data is used to compute trajectories through the following algorithm:

1. When there are more than 10 minutes between two subsequent GPS points in the sequence of all GPS points to be processed, the sequence is separated into two separate sequences and each is processed separately as a possible trajectory in the next step.
2. First, we identify still moments - periods of no movement - as follows:

⁹See <https://developers.google.com/android/reference/com/google/android/gms/location/DetectedActivity>

¹⁰Nevertheless, if other applications request a GPS position, our application also receives this data, even if it occurs on a faster interval

- a) For each GPS point, we identify a subsequent GPS point that was registered at least two minutes after the first one.
 - b) If the average speed between those two points was below 0.6 m/s, the pair is added to a list to be processed in the next step.
 - c) The list of pairs of GPS points resulting from the last step is fused into sequences of GPS points as long as possible: Whenever two pairs overlap in their timestamp, they are fused to a new pair covering the combined timespan.
3. The resulting GPS pairs of still moments are used to exclude still moments from the original sequence and divide it into subsequences marking trajectories.

Of each trajectory, the start and end location as well as the respective timestamp are then saved in *trajectory_table*. We tested 0.5 m/s, 0.6 m/s and 0.7 m/s as threshold in step 2b) and found 0.6 m/s to be best fit the tested sample. On the one hand the threshold must be low enough to still include slow walking which might be below 1 m/s. On the other hand, the threshold should not be too low because inaccuracy in GPS data might otherwise induce trajectories where the device has actually not moved at all.

Example of algorithm (TODO):

```
{  
Original dataset:  
Latitude Longitude Time  
44 11 10:11:03  
44.5 11.1 10:11:15  
44.4 11.05 10:12:12  
44.3 11.07 10:33:00  
44 11.2 10:34:00  
}
```

4.3.4. Serving Aggregation Requests

We use the retrofit2 framework¹¹ based on OKHTTP¹² to handle communication with our REST server described in 4.4. An HTTP Interceptor is used to modify incoming and outgoing requests. The interceptor decrypts the request body of incoming messages using the private key of the installation before the body is parsed into Java Objects.

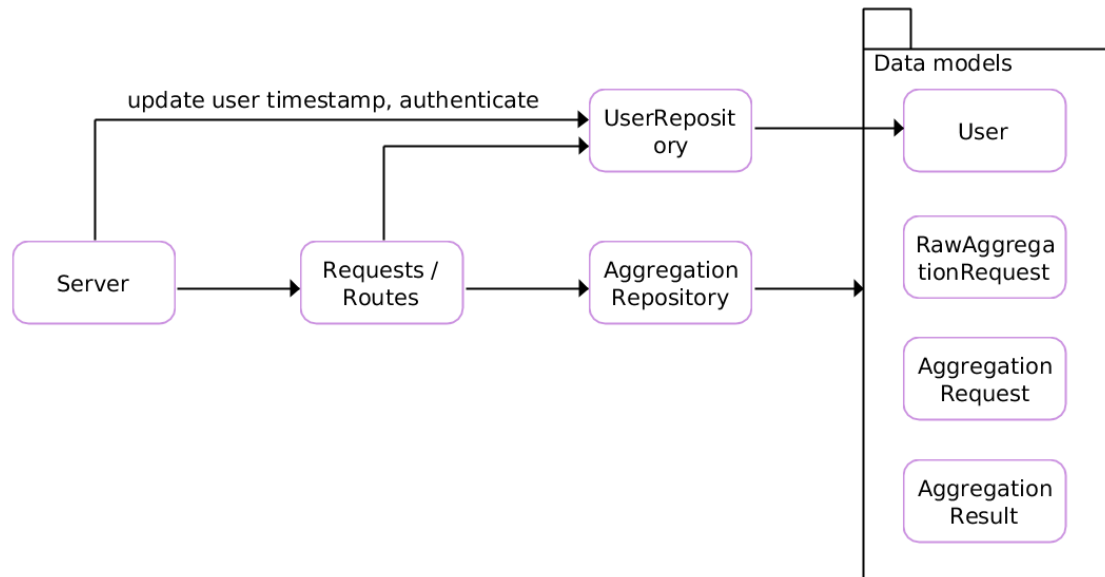
¹¹<https://square.github.io/retrofit/>

¹²<https://square.github.io/okhttp/>

On outgoing messages the interceptor adds authentication before sending them to the server. The app polls for new aggregation requests every 15 minutes. New aggregation requests are first stored locally in the database. Those requests are then processed and the results are again stored locally as pending outgoing requests until they are finally send to the server. Figure XX illustrates this process. This separation of concerns is useful especially in case of an interrupted communication during processing the aggregation request. When the results cannot be send to the server, the app automatically retries the next time that the communication module is invoked. The aggregation itself takes the type parameter of the request to specify which actions to take on the three fields (n:int, value:Float, valueList: List<Float>) shared across all aggregation requests. In case of the tyes "steps" and "activity_X" the field value contains the current mean of the data and the field n is the number of participants so far. In case of "stepsListing" only the field "valueList" is used. Each user's mean value is added to the list. In case of "trajectories", only the field "valueList" is used. Four subsequent elements of the list always represent one trajectory as of latitude of start, longitude of start, latitude of end and longitude of end.

In case that the aggregation should be changed to actually work over P2P e.g. using local WIFI networks this module only has to be adapted to the new routing of requests. No further changes to the application are necessary.

Figure 4.8.: Server architecture



4.4. Server Design and Implementation

The server is build using the event-driven node.js version 10.15.3 leveraging the express¹³ web-server framework and using the mocha¹⁴ testing framework in combination with the chai¹⁵ assertion library for unit and integration testing. The server is designed using a layered architecture as described in figure 4.8. On the lowest level are the data models which define and verify the data schemes defined in subsection 4.4.1. The *commonRepository* and the *userRepository* are build on top of those models and persists data in a mongoDB object store. They also handle e.g. transactions where several objects are modified depending on the result of the previous modification (TODO: implement transactions?). The third level provides the logic to be executed for each endpoint defined in *routes.js* while *server.js* on the top level starts the server on the port specified in *TODO.environment.json-TODO*. It also registers the routes described in section 4.2, adds authentication and interacts directly with the *userRepository* in order to update the respective users *lastSeen* property. Furthermore, it starts a scheduled repeating task to keep the request chain running. This process is described in subsection 4.4.2

4.4.1. Data Model

We organize the data in four collections. The user collection stores the user data which is the public key, the hashed password and "lastSeen" - the timestamp of the last interaction of the user with the server. Aggregation requests are split into two collections. The collection "rawAggregationRequests" stores the initial aggregation request inserted through the admin interface containing the fields start, end, type - the type of the request, the three fields n, value, valueList reused across all aggregations to pass data, the timestamp when the request was filed to the server and a flag indicating whether this request has been started yet¹⁶. Upon start of the aggregation request, a list of the 10 most recently active users is retrieved in order to serve this request. The request body is then encrypted with the first users public key and stored in the collection "aggregationRequests". Each time, a user requests an aggregationRequest, proceeds with it and sends the results back to the server, the result is inserted into the database as a new aggregationRequest. The fields of this collection are

- rawRequestId - The id of the related rawRequest. This field is not available through the API.

¹³See <https://expressjs.com/>

¹⁴See <https://mochajs.org/>

¹⁵See <https://www.chaijs.com/>

¹⁶E.g. when the end data of a newly inserted aggregation request is in the future, the request will be started only when this day has passed

- `started_at` - The timestamp, when the request has been started
- `publicKey` - The public key of the user that should proceed this request
- `nextUser` - The public key of the user that will receive the request afterwards. This is necessary so that the user that should proceed this request can encrypt the processed request with the public key of the next user.
- `previousRequest` - The id of the previous request. This is null, if it is the first request in the chain. This is used for the mechanism taking care if a request is not processed by the user it is pending for.
- `users` - the list of public keys of the following users that will proceed with this request. This field is not available through the API.
- `encryptionKey` - A synchronous key, encrypted with the public key of the user the request is aimed at.
- `iv` - the initialization vector used for synchronous encryption and decryption of the actual aggregation request.
- `encryptedRequest` - The actual aggregation request encrypted with the synchronous key.
- `timestamp` - The timestamp when this object has been created.
- `completed` - A flag indicating whether this aggregation request has already been proceeded by the respective user and the resulting `aggregationRequest` has been received by the server.

The last collection called "aggregationResults" is used to store the results of an aggregation request. Once there are no more users to serve an `aggregationRequest`, the last user sends the final data unencrypted to the server where it is stored as an `aggregationResult`. It contains the same fields as the `rawAggregationRequest` except the `started` flag and additionally a field `started_at` and `timestamp` - indicating when the aggregation request referenced through `rawRequestId` was started and when it was completed.

4.4.2. request-chain

//TODO `server.js` also invokes a scheduled task which re-routes stale requests where the user has not proceeded with the pending request either due to being offline or due to a problem handling the request. When new aggregation requests are started, the `lastSeen` timestamp of users is taken into account to exclude users that have not

4. Design and Implementation

connected for a certain time. Furthermore, the list of users who are selected to serve the new request is ordered by the time the user was last seen.

5. Performance and Evaluation

5.1. Deployment

The proposed Android application¹ and server have been tested on 16 devices for six days from 30.05.2019 until 04.06.2019. However, evaluating the *lastSeen* timestamp of the users showed that only 13 installations were active after 31.05.2019. The server was deployed² in the IBM Cloud as a 128 MB node.js instance. The database was hosted as a free version at [mongodb.com](https://www.mongodb.com).

We ran 7 different aggregation requests defined in section 4.2.1 covering different timespans from single days to the whole time period. Also, we limited the number of participants in each aggregation to 10. The results can be found in the tables in this section³ and will be discussed in chapter 5.3. The results of the collection of trajectories were modified in order to protect the privacy of all research participants. We used this testing period also to improve the performance of the server as well as the Android application and to find and remove bugs. All improvements are incorporated in the most recent commits of the respective GitHub repositories.

5.2. Data Consumption

As expected, the data consumption of the Android application was low. Some research participants provided information about the app's data usage. Table 5.1 shows the 12 collected results of this rather qualitative than quantitative analysis. Screenshots of the provided information are attached in the appendix. On 75% of the devices the data consumption for 6 days was below 21 MB. The highest data consumption was 376 MB and is attributable to an error that occurred on the last day of the testing period.

¹The version deployed during field testing can be found at <https://github.com/SimonVanEndern/location-app/releases/tag/v1.0>

²The version deployed during field testing can be found at <https://github.com/SimonVanEndern/location-server/releases/tag/v1.0> The final version provides the same functionality but includes improvements (e.g. bug-fixes, comments, ...) over the deployed version.

³The results are also available in JSON format at <https://github.com/SimonVanEndern/tum-thesis-latex/tree/master/published/results>

Requests were sent multiple times during the whole period and up to unlimited times on the last day due to this error in combination with another error that was present during the whole timespan. This explains the high data consumption reported also by two other participants as the error occurred at least on three devices. The data suggests that the average data consumption in a production environment would be lower and regarding the size of our setup below 100 MB per month. Furthermore, a mobile connection is not required. In the future, the setup can be changed to use only or mostly WIFI connections. Some few available reports about battery consumption also indicate a rather moderate battery consumption.

| Combined Data Consumption from 30.05.2019 - 04.06.2019 | |
|---|---|
| Data Consumption | Comment |
| 20.5 MB | |
| 8.18 MB | |
| 8.5 MB | From 01.04.2019 - 04.06.2019 only |
| 11.98 MB | |
| 6.29 MB | mobile data only, WIFI unknown |
| 0 MB | mobile data only, WIFI unknown |
| 376 MB | Known and fixed error |
| 6.19 MB | mobile data only, WIFI unknown |
| 2.6 MB | |
| 18.53 MB | 16.5 MB mobile, 2.03 MB WIFI |
| 75.73 MB | From 01.04.2019 - 04.06.2019 only, possible error candidate |
| 49.39 MB | mobile data only, WIFI unknown, possible error candidate |

Table 5.1.: Data consumption of the Android application during the testing period.

5.3. Aggregation Results

We computed the aggregation over several different timespans. We aggregated the average number of steps across all participants, the average time spent walking, running, in a vehicle or on a bicycle and collected a list of each individual's average number of steps. The results can be seen in table 5.5 - 5.4. The results from the experimental aggregation of trajectories will be discussed in chapter 5.3.4.

| Average number of steps per day | | |
|---------------------------------|---|--------|
| Time period / Day | N | Value |
| 30.05.2019 | 5 | 4332 |
| 31.05.2019 | 5 | 3440 |
| 01.06.2019 | 1 | 17 |
| 03.06.2019 | 3 | 455411 |
| 04.06.2019 | 3 | 1123 |
| 30.05.2019 - 31.05.2019 | 6 | 3279 |
| 30.05.2019 - 31.05.2019 | 5 | 4030 |
| 30.05.2019 - 01.06.2019 | 6 | 4271 |
| 30.05.2019 - 03.06.2019 | 5 | 643959 |
| 30.05.2019 - 04.06.2019 | 5 | 3623 |
| 31.05.2019 - 04.06.2019 | 3 | 4681 |
| 01.06.2019 - 04.06.2019 | 2 | 2089 |
| 02.04.2019 - 04.06.2019 | 3 | 326363 |
| 03.04.2019 - 04.06.2019 | 3 | 455904 |

Table 5.2.: Average number of steps per day

| Average time spent walking [min] | | |
|----------------------------------|----|--------|
| Time period / Day | N | Value |
| 30.09.2019 | 10 | 17.63 |
| 31.05.2019 | 10 | 97.25 |
| 01.06.2019 | 4 | 72.85 |
| 03.06.2019 | 9 | 146.49 |
| 04.06.2019 | 8 | 86.77 |
| 30.05.2019 - 31.05.2019 | 10 | 63.20 |
| 30.05.2019 - 31.05.2019 | 10 | 78.41 |
| 30.05.2019 - 01.06.2019 | 10 | 75.14 |
| 30.05.2019 - 03.06.2019 | 9 | 81.94 |
| 30.05.2019 - 04.06.2019 | 8 | 61.86 |
| 31.06.2019 - 04.06.2019 | 7 | 67.63 |
| 01.06.2019 - 04.06.2019 | 7 | 57.44 |
| 02.06.2019 - 04.06.2019 | 7 | 51.75 |
| 03.06.2019 - 04.06.2019 | 7 | 64.04 |

Table 5.3.: Average time spent walking

5.3.1. Validity of results

Excluding the value of 30.05.2019 where the collection of results started and the data of the fraction of the day before installation is not included, the average time spent walking computed in the aggregations ranges from 52 to 146 minutes per day. The average time spent walking computed for the whole timespan where data from 8 devices is included amounts to 62 minutes. At the same time, reported number of steps per day ranges from 1123 to 4682 per day, excluding outliers (see table 5.2). The average number of steps computed for the whole timespan with data from 5 participating devices is 3622. Taking an average speed of roughly 100 steps per minute into account [23], the results from average time spent walking and average number of steps diverge by a factor of 1.7. Considering that apparently only 6 devices had a step sensor on their phone and the sample size of the steps aggregation is only about half of the sample size of the other aggregations, this deviation does not seem out of range. Furthermore there was an error computing the number of steps on at least one phone which caused some of the results to be off. Regarding table 5.4 it is clear, that the error occurred in the local aggregation process and it seems likely that only on one device this error occurred, as in each aggregation only one value is erroneous.

| Listing of average number of steps per day per of each participant | | |
|--|----|-----------------------------------|
| Time period / Day | N | Values |
| 30.05.2019 | 10 | 1661, 1246, 3195, 7714, 7842 |
| 31.05.2019 | 10 | 2747, 775, 3924, 9203 |
| 01.06.2019 | 4 | 17 |
| 03.06.2019 | 9 | 914, 1042, 1364278 |
| 04.06.2019 | 8 | 1103, 13332, 934 |
| 30.05.2019 - 31.05.2019 | 10 | 550, 2905, 3195, 775, 8828, 3419 |
| 30.05.2019 - 31.05.2019 | 10 | 550, 8828, 5145, 3195, 2433 |
| 30.05.2019 - 01.06.2019 | 10 | 8126, 3195, 283, 8828, 1629, 3565 |
| 30.05.2019 - 03.06.2019 | 9 | 598, 3195, 1669, 8828, 3205505 |
| 30.05.2019 - 04.06.2019 | 8 | 757, 1719, 8828, 3614, 3195 |
| 31.05.2019 - 04.06.2019 | 7 | 757, 9203, 4084 |
| 01.06.2019 - 04.06.2019 | 7 | 826, 33511 |
| 02.04.2019 - 04.06.2019 | 7 | 974691, 1231, 3168 |
| 03.04.2019 - 04.06.2019 | 7 | 1364278, 1231, 2202 |

Table 5.4.: Average number of steps per day of each user.

The average time spent running ranges up to only 2 minutes per day with one day having a 0 value across 8 participating devices. Most likely, only a fraction of the research participants goes running on a regular basis and probably even less take their smartphone along. This assumption reflects the low values for the average time spent running. Especially as attributing the whole average value to one device only would mean that one person conducted a run of less than 20 minutes seems unlikely, while it seems likely that quite some persons run for a short amount of time e.g. to catch a train or metro. Nevertheless, we cannot prove this assumption. It furthermore highlights, that for thorough analysis the aggregation of the mean value is not sufficient. Similarly, the average number of steps across all participants of an aggregation is not very robust and does not provide much information. Table 5.4 provides the average number of steps of each individual participating in the request. With the median ranging from 3195 to 3565 in all aggregations with more than 3 valid values and a mean number of 5205 steps in Germany [19] we have no doubt about the correctness of the values except the values resulting from (probably only one) erroneous device.

The average time spent in a vehicle (including means of public transport) ranges from 36 to 120 minutes (excluding the aggregation only covering the first day). The average time spent in a vehicle computed by the aggregation covering the whole timespan is 59 minutes. The average time spent biking ranges from 11 to 31 minutes. The aggregation

| Average time spent running [min] | | |
|----------------------------------|----|-------|
| Time period / Day | N | Value |
| 30.09.2019 | 10 | 1.43 |
| 31.05.2019 | 10 | 1.21 |
| 03.06.2019 | 9 | 1.95 |
| 04.06.2019 | 8 | 0 |
| 30.05.2019 - 31.05.2019 | 10 | 1.96 |
| 30.05.2019 - 31.05.2019 | 10 | 1.89 |
| 30.05.2019 - 01.06.2019 | 10 | 1.73 |
| 30.05.2019 - 03.06.2019 | 9 | 1.79 |
| 30.05.2019 - 04.06.2019 | 8 | 1.55 |
| 31.06.2019 - 04.06.2019 | 7 | 0.83 |
| 01.06.2019 - 04.06.2019 | 7 | 0.99 |
| 02.06.2019 - 04.06.2019 | 7 | 1.05 |
| 03.06.2019 - 04.06.2019 | 7 | 0.16 |

Table 5.5.: Average time spent running.

| Average time spent in a vehicle [min] | | |
|---------------------------------------|----|--------|
| Time period / Day | N | Value |
| 30.05.2019 | 10 | 8.24 |
| 31.05.2019 | 10 | 119.92 |
| 01.06.2019 | 4 | 90.40 |
| 03.06.2019 | 9 | 71.74 |
| 04.06.2019 | 8 | 35.51 |
| 30.05.2019 - 31.05.2019 | 10 | 68.40 |
| 30.05.2019 - 31.05.2019 | 10 | 90.49 |
| 30.05.2019 - 01.06.2019 | 10 | 82.19 |
| 30.05.2019 - 03.06.2019 | 9 | 62.16 |
| 30.05.2019 - 04.06.2019 | 8 | 59.34 |
| 01.06.2019 - 04.06.2019 | 7 | 51.07 |
| 31.05.2019 - 04.06.2019 | 7 | 68.20 |
| 02.04.2019 - 04.06.2019 | 7 | 48.89 |
| 03.04.2019 - 04.06.2019 | 7 | 45.96 |

Table 5.6.: Average time spent in a vehicle.

| Average time spent biking [min] | | |
|---------------------------------|----|-------|
| Time period / Day | N | Value |
| 30.05.2019 | 10 | 2.29 |
| 01.06.2019 | 4 | 14.99 |
| 03.06.2019 | 9 | 31.03 |
| 04.06.2019 | 8 | 23.81 |
| 30.05.2019 - 31.05.2019 | 10 | 11.10 |

Table 5.7.: Average time spent biking.

across the whole timespan was erroneous and had to be discarded. Both aggregations - the time spent biking and the time spent in a vehicle seem reasonable to us, especially as all participants were active in the Munich area which is also a quite bike friendly city. Furthermore, the collected trajectories in section 5.3.4 show two rather long trajectories which probably contributed to the high average time spent in a vehicle in one of the aggregations.

5.3.2. Implications

The same aggregation covering the timespan from 30.05.2019 to 31.05.2019 has been computed twice for 4 of the aggregations. The results show that not only the number of participants varies but also the value varies from 4% to 18%. For each aggregation, the selected participants might be different. With a total number of 16 participants in our research, though there has to be an intersection. Nevertheless, the deviation in the same aggregation is within an acceptable range and will be more robust, when the number of participants is increased. Furthermore, the median value will be even more robust to computing the same aggregation twice. We see the same phenomena e.g. regarding the aggregations for the time spent running on 30.05.2019, 31.05.2019 and the aggregation covering both days. The aggregation covering both days yields a higher value than each single aggregation which is due to the fact of a differing underlying research population.

5.3.3. Inference on aggregated data

Regarding the aggregations computing mean values across all participants, we do not see any possibility to infer information about the users who participated. Especially as not even the identifier of the user used by the server is published alongside with the aggregation. In addition, we showed that when repeating the same aggregation request, the user base changes and the value accordingly. Nevertheless, from the two different values for the same aggregation request, we cannot draw any inference. Also an overlapping in the aggregations e.g. 30.05.2019, 31.05.2019 and an aggregation covering both days as well as the aggregations covering the first two days and the last five days does not allow for inference in the case of mean values aggregated across all participants. The collection of individual mean values in table 5.4 shows that with a high probability the same participant with an average value of 8828 steps per day who participated in the request covering the first two days also participated in the request covering the whole timespan. Based on the information that days with a zero value of the step counter are excluded also locally from the aggregation, we can infer that this person's phone did not register any steps after the first two days, otherwise they value

in the second aggregation would have changed. Similarly, we can infer that the user with 757 steps on average per day who participated in the requests covering the whole 6 day timespan and the aggregation covering only the last 5 days probably did not register steps on the first day. In summary, we can only infer information about the fact of some users who participated in one aggregation, did or did not participate in another aggregation. As this information is endogenous to our data-set, it can in no case be used to be linked to other data-sets in order to infer further information. Assuming that we would also have the collected averages of all aggregations and not only steps. Furthermore, linking different aggregations e.g. the number of steps and the time spent walking is only possible with a limited likelihood due to the following reasons. First, there is no exact match of e.g. time spent walking for a number of steps, so that always a range of values of another aggregation would fit. Next, we cannot even assure that the same user participated in both different aggregations. Assuming nevertheless that there is another data-set where the combination of walked steps, time spent walking, etc. is a quasi-identifier as defined by [21], the likelihood of obtaining a correct link is low due to the afore mentioned reasons and a good level of anonymity regarding k-anonymity would always be supplied naturally. And still, those aggregations do not expose raw location data / GPS data points which are usually the risk.

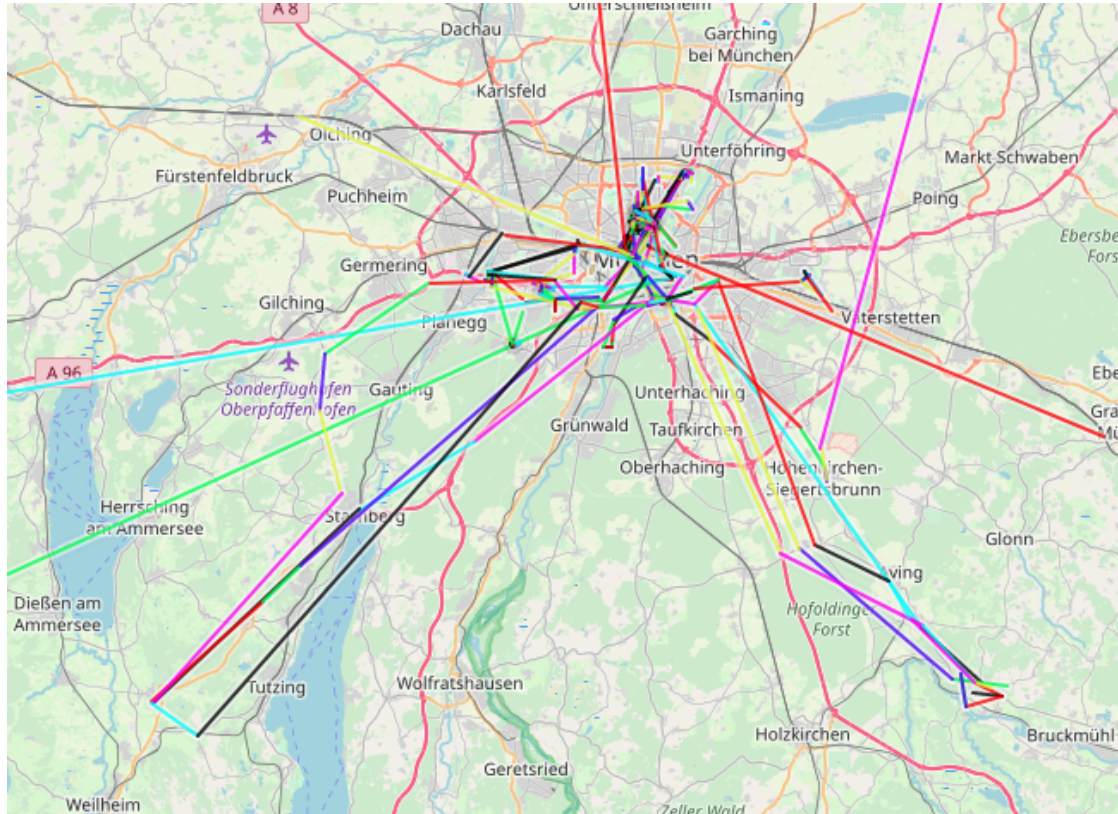
5.3.4. Trajectories

Using the algorithm described in chapter 4.3.3 a total of 406 trajectories were computed from the raw GPS data. The experimental collection of trajectories clearly shows privacy risks as pointed out e.g. by [4], why we only publish a modified results set⁴. Most of the trajectories are shown in figure 5.1 and 5.2. The complete results can also be found in the appendix.

In more than 90% of all cases, local trajectories obtained from one of the researchers participating devices could match an activity for more than 90% of the time. This data suggests that these trajectories can easily be linked locally and a change of transport system e.g. metro to bike can be identified. For example at the locations Giselastraße, Odeonsplatz and central station, many trajectories end or start. Mapping the current activity to those trajectories enables aggregations as mentioned in chapter 3.1. Noting that locally all data points of the trajectories are available, it is also easy to compute the distance travelled or infer a more granular transport definition e.g. train, bus, car, etc. Also, it is possible to compute how many people combine e.g. bike and car or public transport and furthermore identify, which station is most likely (in case of

⁴Whenever the trajectory started or ended clearly in a precise private location e.g. housing or work-place, we slightly modified this trajectory. So the results do not represent actual trajectories. Nevertheless, the meaning of the results should not have changed.

Figure 5.1.: Trajectories



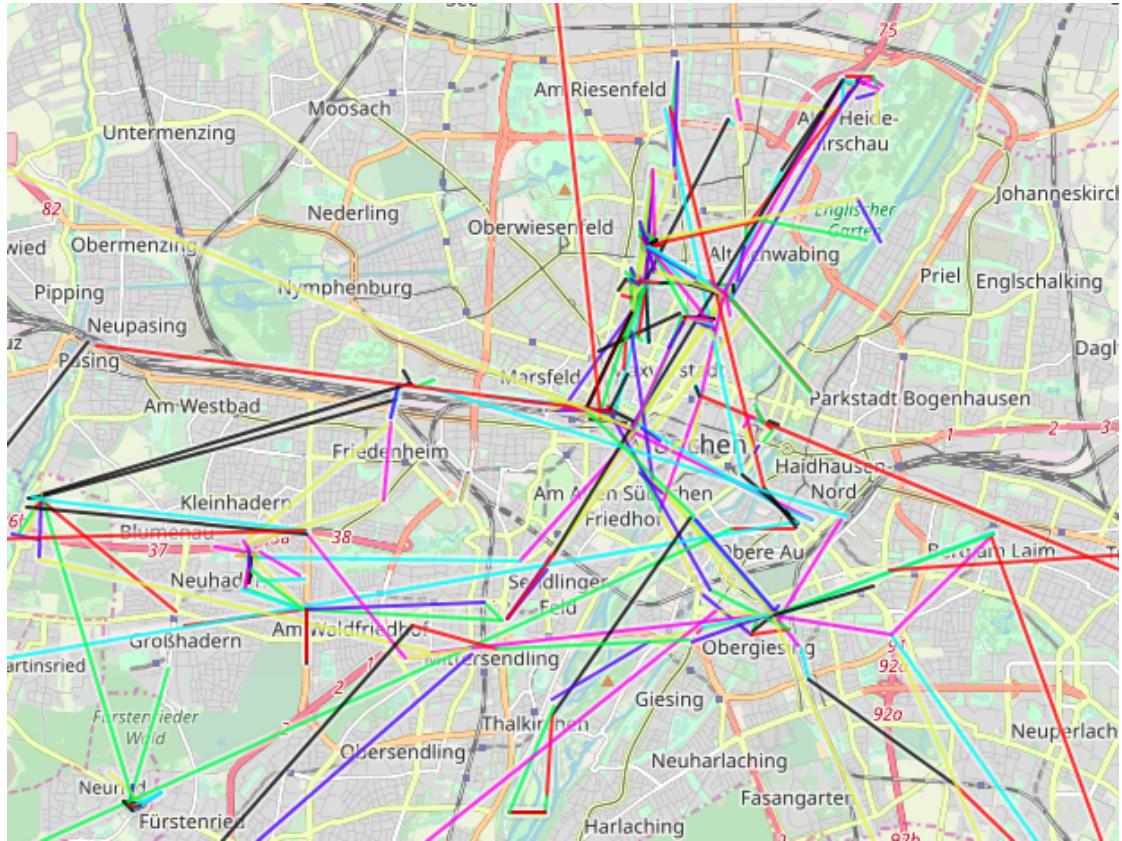
public transport) to be combined with bike.

The trajectories also show, that in the not anonymized data set we had at hands, often trajectories of which the endpoints did not clearly indicate an ongoing, the numeration of the trajectory indicated that they were connected. This 1. highlights the importance of list order randomization and 2. shows that our algorithm of finding trajectories still has possibilities for improvement. (TODO: Randomize order in published data set)

5.4. Further scenarios e.g. Scenario Traffic Jam

Also it would be possible to create a "travelled road" or "travelled public transport" map by aggregating the trajectories data so that there would not be multiple trajectories but each road either marked used or not used. This way it would be possible to identify that a user lives in a certain area but could not link to the work location due to obfuscation with other routes joining. This map would on the other hand allow to have a feeling of

Figure 5.2.: Trajectories



which areas are covered by the data / app. This map could also be extended with the average speed on the respective road depending on vehicle or bike and also compare whether bike is faster.

A typical scenario of google maps is to notify users about traffic jams and suggest alternate routes. The calculation of alternate routes taking traffic jams into account can clearly happen locally with the maps data. Google maps works when offline. The data about all current traffic jams can also be made publicly available through a server. Generating the data can also happen without exposing raw data: The user downloads a map containing data of the usual speeds at each street. While the user is driving, the app registers the speed and compares it in the background to the normal speed. If the speed is significantly lower, the user chooses a random list of known users and sends the signal as a request for those users to the server. They randomly according to a fixed percentage choose to inform the server about the traffic jam or forward the

signal another time. The signal contains a unique id thus that the server even when receiving it multiply times knows it is from one user. If more than a threshold of signals is received, a traffic jam is "created". Also the request is not forwarded anymore after a certain time to stop it from spreading unlimited.

6. Conclusion and Discussion

We have shown that our hypothesis holds but only for aggregated data. This is fine because except in one experimental setting as with trajectories, there is no need for (anonymized) raw data. Also the experimental setting could be replaced by directly implementing the aggregations and testing with a greater user base. In theory, the anonymity and preserved privacy that hold for the tested aggregations holds also for the other proposed aggregations and aggregations we have not evaluated here. In order to follow with this research, we provide the setup to easily implement and test those and further aggregations in future research in order to further support our hypothesis.

So, publishing pure mean values of analysis does not impose any privacy risk to the user.

Furthermore we showed exemplary with the steps aggregation, that also the collection of the users' average values is possible without posing a privacy risk to the user. This implies, that not only the aggregation of mean values but also the collection of locally aggregated data is possible without privacy concerns. We do not see any possibility to infer that the same user took part in different aggregations. There might be a chance in case of a very high steps value and a very high mean time spent walking or running, but this would only allow to assume the same user being present in both aggregations which allows for no further inference.

6.1. Limitations

- As mentioned in XX, our system is based on trust. In case of user data being compromised, this significantly impacts some of the results. Nevertheless, collecting the list of mean values is far less error prone as the outlier could also be identified.

6.2. Future Work

- While XX has found that inference attacks can be based on the same dataset being published two times with different anonymization techniques applied and

XX shows that anonymized datasets that overlap poses a risk, it still has to be investigated whether overlapping aggregated data as in our case can pose a risk.

- Some of the techniques identified as useful, such as spatial cloaking, ... should be applied in our setting.
- Our framework would also allow to pre-populate (simulated) smartphones with artificially generated or otherhow collected data in order to test and verify the functionality.
- Another use of our framework is the area of decentralized computation. Problems might be solved locally and collected by the server afterwards and the pieces put together still with anonymity for the users.
- Ask the user for his / her home and work location or infer it from the data in order to process aggregation requests as mentioned in section XX.
- Store the users location on the server (granularity level approach) in order to allow for aggregation requests targeted at specific areas and not the overall user base. (levels and level database and unlocked levels collection necessary)
- Android: Only send a final aggregation back to the server when criteria like minimum n, ... are met.
- Server: Apply that the request has a counter how many times it was actually retrieved, so that after e.g. 10 times it was retrieved but never answered, the request is rolled back because apparently the user somehow cannot process the request.
- Have a nice activity informing the user / displaying some information to the user.
- (Put somewhere else!! TODO) The app can send traffic alerts to the server if it is on a route where usually traffic is far faster. (This also via other nodes in order to not letting the server know who is on this route.)
- a scheduler which automatically creates aggregation requests on a regular basis so that not as now Postman has to be used to start requests. The Postman collections used during field testing can be found here : XXX
- Generate userList of aggregation request dynamically.
- Pre-populate (raw)aggregation requests so that the first users cannot infer data from the users before them with high probabilities. Upon receiving the final

result, the server can look up the used initialization values from the rawAggregationRequest and calculate the actual correct result and insert it into the database.

- Generate some use for the user of the application e.g. through showing locally aggregated data and comparing it to aggregated data publicly available e.g. in order to show how many percent walked more this day than the respective user. This is also an approach to motivate app installations in case of a broader user base necessary.
- delete local data after some time.
- Evaluate how many people have to participate in an aggregation request to be representative / how to choose users participating in it in order to not be biased (e.g. when taking always the most recent active users, the users only online a few times a day are discriminated against)
- Investigate overlapping
- When returning list data to the server, the list order should be randomised by the user before sending.
- Verify public key e.g. through telephone number passing and verifying this way and then building a network of verifications.
- Adhere to standards of schema.org
- In the future it should be implemented that devices change the public / private key pair from time to time
- It could also be implemented to verify another public key but requesting an SMS, (harder with changing public-private key-pairs)
- Possible future aggregation: List of all activity times e.g. for walking to find out mean, medium, ...
- When dynamically adding users to the aggregation request, we need another field for setting the limit so that when the limit is reached, the user sends back the result.
- Pre-fill the aggregation request so that the second user cannot read the first user's share.

- The application can be modified to be a framework that can easily be incorporated into other apps.
- It also allows for computing how many people go by car, bike, or mixed to work.

We suggest to do especially two things in any further research project: Implement error logging and sending those errors to the server in order to find and remove bugs. Secondly bring the application to the playstore so that updates are possible (even without user interaction in case of automatic updates) and be able to activate a broader user base while still working on the final version to be tested.

Also, the approach could easily be integrated into existing projects like open maps or be build modular in order to allow using it as a library with other applications.

6.3. Evaluation

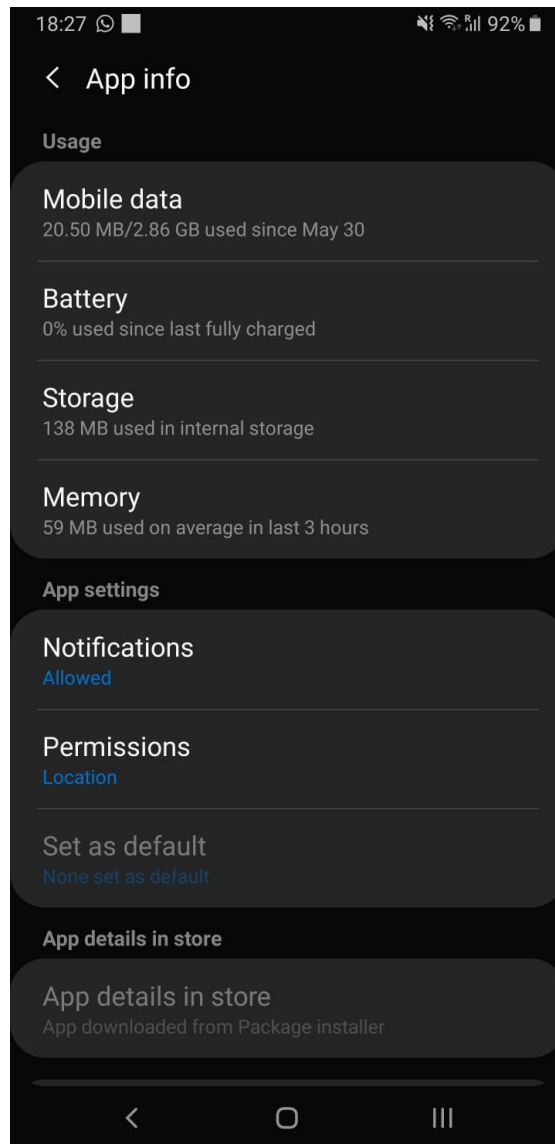
The results support our hypothesis that data can be analyzed decentrally and that aggregated data can be published without any privacy concerns. The aggregations of mean values clearly leave no doubt about full privacy protection as there even is no personal data involved anymore. The listing of mean values of average number of steps per participant allows for more advanced statistical analysis while at the same time the values cannot be mapped to persons. Even when conducting the same request twice, due to users being chosen dynamically, one could most probably see if the same user participated in the second request but nothing else. When aggregating over another time period (that might have an intersection with the other one), there is not even the chance to identify whether the same user participated in both aggregations. As requests are started not at the same time (TODO!!!), and users are in a future setting allocated dynamically to the request, there is also no chance to link the data from different aggregations. From the number of steps one could infer the time somebody spent walking, but as it is not given whether the steps were conducting walking, running or both, this linking would result in a very poor performance and also only reveal that to a very low probability, the user with X steps in one aggregation is the same user spending X minutes walking the same day. The listing of trajectories created a dataset that clearly shows the vulnerability highlighted in XX and XX. Nevertheless, this is no aggregation but just a collection of raw data with stripped of identifiers and timestamps anonymized to a daily basis. The results show clearly that our setup is sufficiently accurate to field test the other aggregations proposed in XX and further prove our thesis. The data shows, that e.g. A change of transportation system can clearly be identified.

6.4. Limitations

As stated in the introduction, our approach is based on trust among the clients and the server. Nevertheless, while e.g. [13] face the same problems in case of a compromised server - namely that the server can create artificial participants and thus obtain the raw values from each user, our setup is more general, allows for more complex aggregations and can be adapted to work via P2P or to establish trust through something like XXX

6.5. Reproducability Considerations

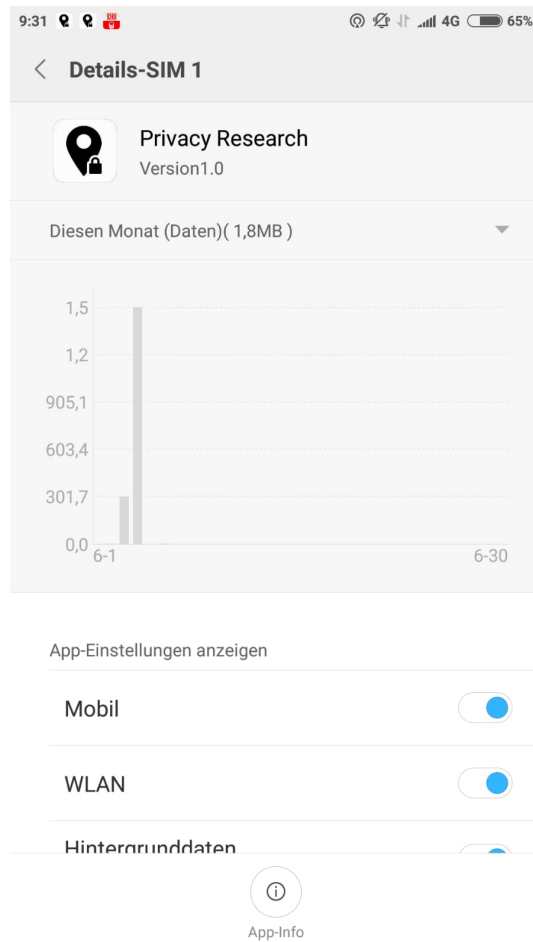
A. Data Usage Screenshots



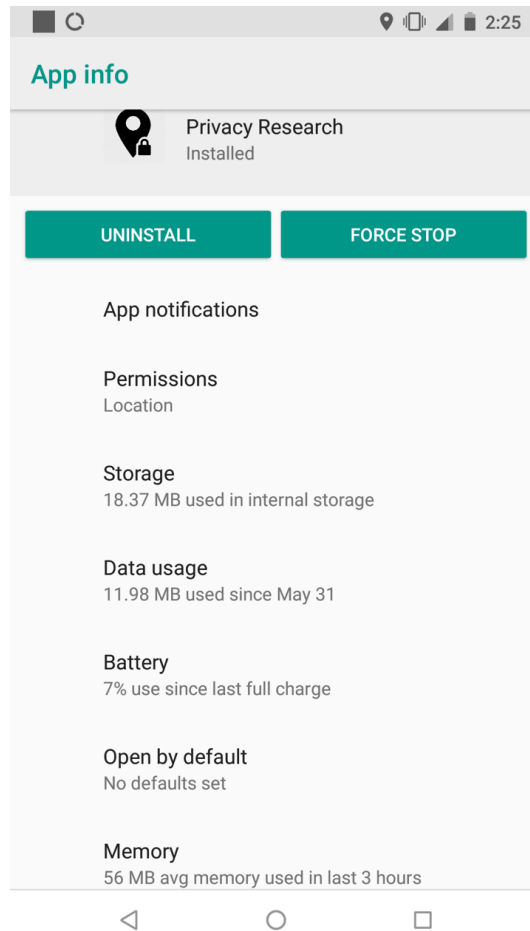
A. Data Usage Screenshots

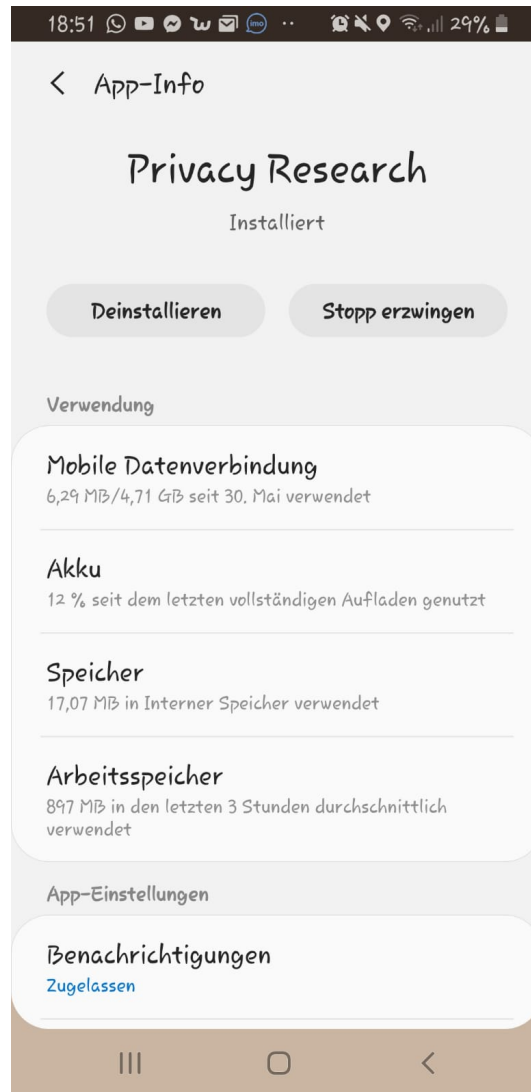


A. Data Usage Screenshots

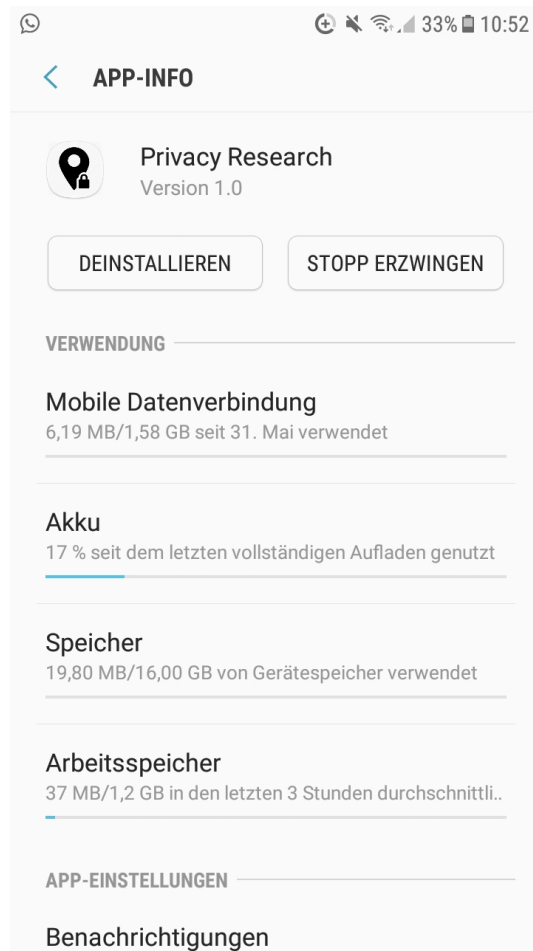


A. Data Usage Screenshots

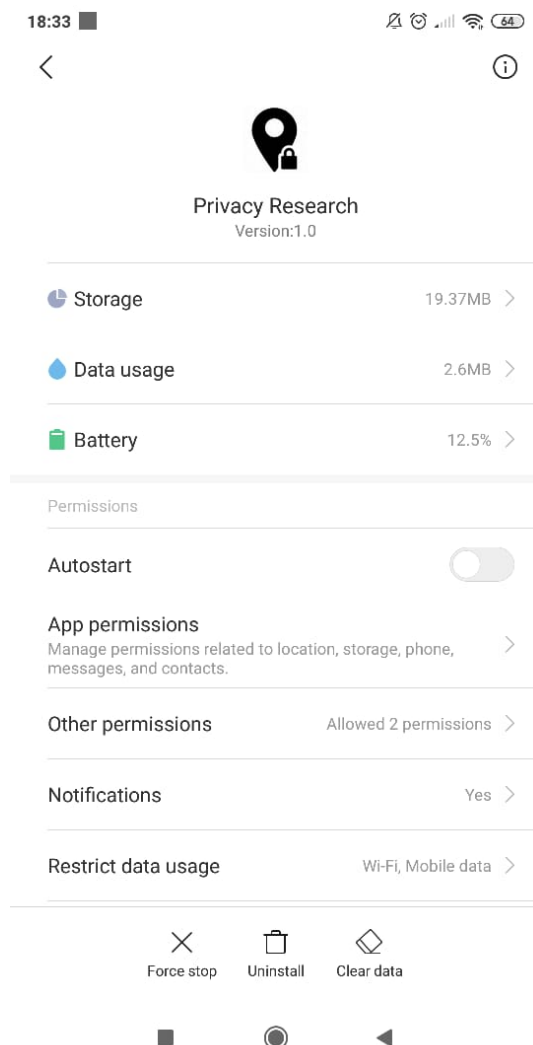




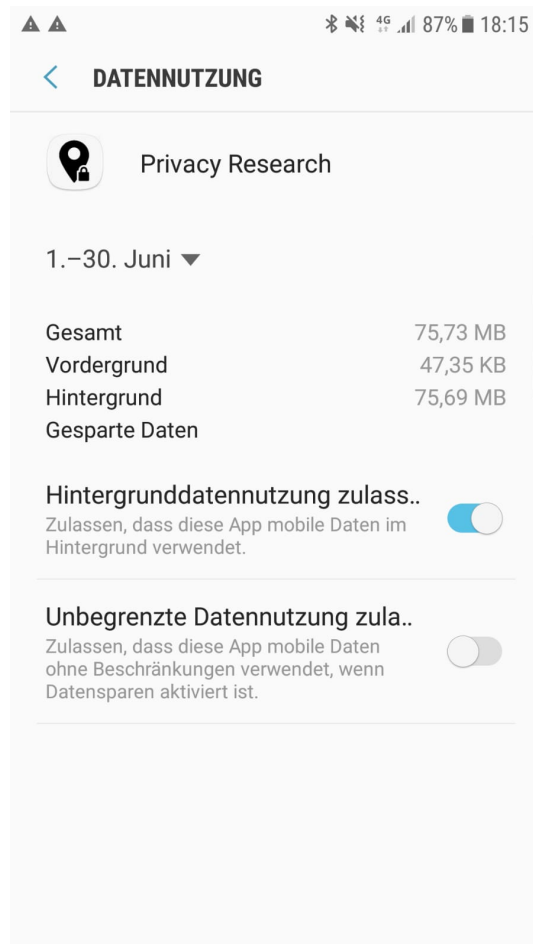
A. Data Usage Screenshots

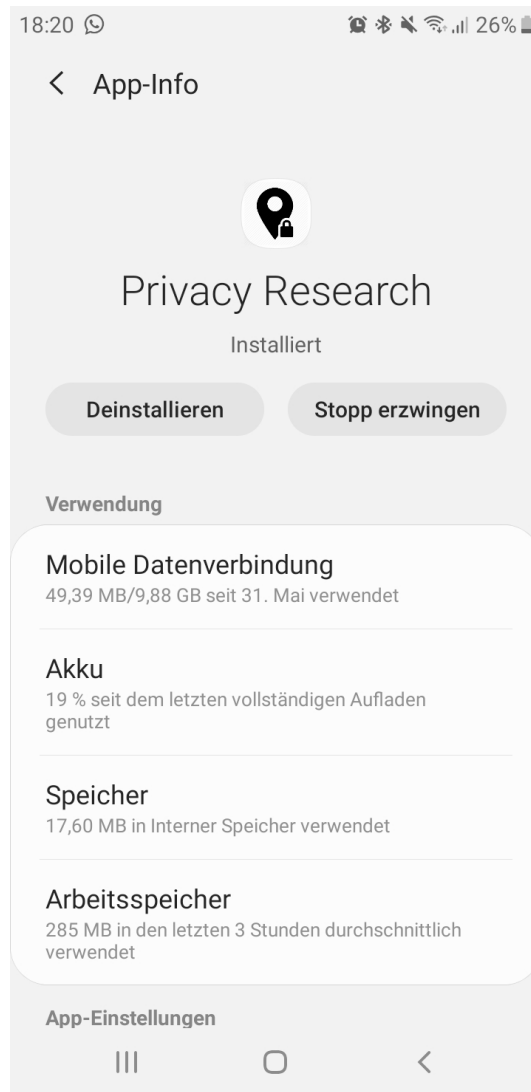


A. Data Usage Screenshots

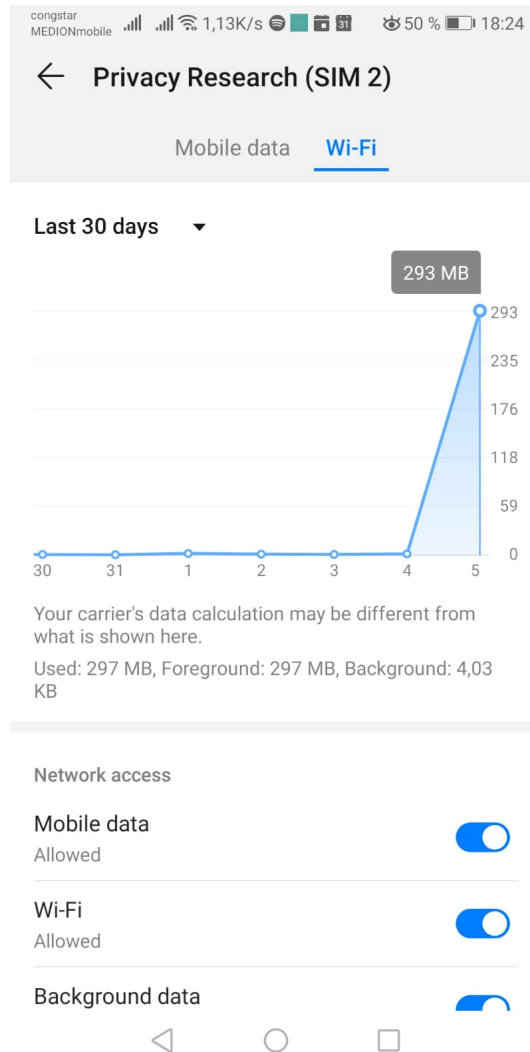


A. Data Usage Screenshots

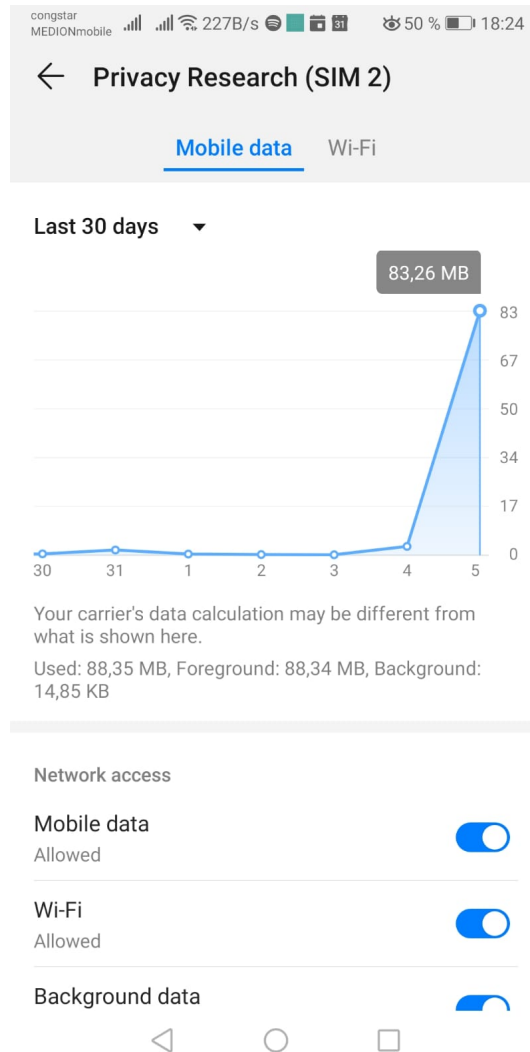




A. Data Usage Screenshots



A. Data Usage Screenshots



List of Figures

| | |
|---|----|
| 3.1. Incoming aggregation request. | 7 |
| 3.2. Outgoing aggregation response. | 7 |
| 3.3. Decentral unencrypted aggregation process via P2P. | 8 |
| 3.4. Decentral encrypted aggregation process. | 9 |
| 4.1. Body of an HTTP request for creating a new user. | 13 |
| 4.2. Body of an HTTP response upon successful creation of a new user. . . . | 13 |
| 4.3. Body of an HTTP response listing all aggregation requests for the speci- fied user. | 14 |
| 4.4. Body of an HTTP request sending the processed aggregation request back to the server. | 14 |
| 4.5. Body of an HTTP response listing all completed aggregations. | 15 |
| 4.6. Body of an HTTP request for initiating a new aggregation. | 15 |
| 4.7. Android architecture | 17 |
| 4.8. Server architecture | 21 |
| 5.1. Trajectories | 32 |
| 5.2. Trajectories | 33 |

List of Tables

| | |
|---|----|
| 5.1. Data consumption of the Android application during the testing period. | 26 |
| 5.2. Average number of steps per day | 27 |
| 5.3. Average time spent walking | 27 |
| 5.4. Average number of steps per day of each user. | 28 |
| 5.5. Average time spent running. | 29 |
| 5.6. Average time spent in a vehicle. | 29 |
| 5.7. Average time spent biking. | 29 |

Bibliography

- [1] A. R. Beresford and F. Stajano. "Location privacy in pervasive computing." In: *IEEE Pervasive computing* 1 (2003), pp. 46–55.
- [2] A. R. Beresford and F. Stajano. "Mix zones: User privacy in location-aware services." In: *IEEE Annual Conference on Pervasive Computing and Communications Workshops, 2004. Proceedings of the Second*. IEEE. 2004, pp. 127–131.
- [3] G. Developers. Accessed: 2019-06-07.
- [4] K. Drakonakis, P. Ilia, S. Ioannidis, and J. Polakis. "Please Forget Where I Was Last Summer: The Privacy Risks of Public Location (Meta)Data." In: *CoRR abs/1901.00897* (2019). arXiv: 1901.00897.
- [5] T. Economist. Accessed: 2019-06-20.
- [6] P. Golle and K. Partridge. "On the Anonymity of Home/Work Location Pairs." In: *Proceedings of the 7th International Conference on Pervasive Computing*. Pervasive '09. Nara, Japan: Springer-Verlag, 2009, pp. 390–397. ISBN: 978-3-642-01515-1. DOI: 10.1007/978-3-642-01516-8_26.
- [7] T. Guardian. Accessed: 2019-06-20.
- [8] T. Guardian. Accessed: 2019-04-10.
- [9] B. Hoh and M. Gruteser. "Protecting location privacy through path confusion." In: *First International Conference on Security and Privacy for Emerging Areas in Communications Networks (SECURECOMM'05)*. IEEE. 2005, pp. 194–205.
- [10] B. Hoh, M. Gruteser, H. Xiong, and A. Alrabady. "Enhancing security and privacy in traffic-monitoring systems." In: *IEEE Pervasive Computing* 5.4 (2006), pp. 38–46.
- [11] B. Hoh, M. Gruteser, H. Xiong, and A. Alrabady. "Preserving Privacy in Gps Traces via Uncertainty-aware Path Cloaking." In: *Proceedings of the 14th ACM Conference on Computer and Communications Security*. CCS '07. Alexandria, Virginia, USA: ACM, 2007, pp. 161–171. ISBN: 978-1-59593-703-2. DOI: 10.1145/1315245.1315266.

- [12] W. A. Jabbar, M. Ismail, and R. Nordin. "Peer-to-peer communication on android-based mobile devices: Middleware and protocols." In: *2013 5th International Conference on Modeling, Simulation and Applied Optimization (ICMSAO)*. IEEE. 2013, pp. 1–6.
- [13] H. Kajino, H. Arai, and H. Kashima. "Preserving worker privacy in crowdsourcing." In: *Data Mining and Knowledge Discovery* 28.5-6 (2014), pp. 1314–1335.
- [14] J. Krumm. "Inference Attacks on Location Tracks." In: *Pervasive Computing*. Ed. by A. LaMarca, M. Langheinrich, and K. N. Truong. Berlin, Heidelberg: Springer Berlin Heidelberg, 2007, pp. 127–143. ISBN: 978-3-540-72037-9.
- [15] M. F. Mokbel, C.-Y. Chow, and W. G. Aref. "The new casper: Query processing for location services without compromising privacy." In: *Proceedings of the 32nd international conference on Very large data bases*. VLDB Endowment. 2006, pp. 763–774.
- [16] H. B. Review. Accessed: 2019-06-20.
- [17] P. Samarati and L. Sweeney. "Protecting privacy when disclosing information: k-anonymity and its enforcement through generalization and specialization." In: *Proceedings of the IEEE Symposium on*. 1998.
- [18] Statista. Accessed: 2019-06-07.
- [19] Statista. Accessed: 2019-06-23.
- [20] M. Stolpe. "The internet of things: Opportunities and challenges for distributed data analysis." In: *ACM SIGKDD Explorations Newsletter* 18.1 (2016), pp. 15–34.
- [21] L. Sweeney. "Achieving k-anonymity privacy protection using generalization and suppression." In: *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 10.05 (2002), pp. 571–588.
- [22] L. Sweeney. "k-anonymity: A model for protecting privacy." In: *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 10.05 (2002), pp. 557–570.
- [23] T. N. Y. Times. Accessed: 2019-06-23.
- [24] H. Zang and J. Bolot. "Anonymization of location data does not work: A large-scale measurement study." In: *Proceedings of the 17th annual international conference on Mobile computing and networking*. ACM. 2011, pp. 145–156.