

# How to count travelers without tracking them between locations

May 26, 2023

## 1 Introduction

As urbanization continues to grow, increasingly more people use public and individual modes of transportation. To implement policies that increase the sustainability of urban transportation systems, a deeper understanding of travel patterns is essential. Traditional surveys and travel diaries require significant effort and provide only snapshots [1]. Various more recent technologies allow automated counting, e.g., Bluetooth and Wi-Fi detection systems and automated fare collection systems. Information gathered by these systems has proved to be helpful in improving security, physical activity, traffic safety, public transportation, communication infrastructure [6, 4], and the overall quality of life for citizens.

However, counting travelers at specific locations by using smart-card IDs allows tracking their movements between locations. It is, therefore, a sensitive issue, especially if it allows monitoring travelers over an extended period of time: the trade-off for valuable insights into movement patterns is an infringement upon their privacy. It has been shown that a few points are enough to identify individual travelers with simple anonymization and persistent identifiers [3].

To prevent such situations, various regulations have been adopted, including Europe's General Data Protection Regulation (GDPR) [5], which requires parties to obtain explicit consent before collecting and using personal information. Obtaining explicit consent may reduce the completeness of the data collection, which can introduce bias and reduce the representativeness of the results. Even if consent is granted, individuals must trust that their data will be used responsibly and not misused for other purposes.

For these reasons, we challenge the feasibility of robust privacy protection within a system that relies on identifying travelers to count them. Instead, we propose an alternative system that offers statistical counts of travelers as the only accessible information. To implement such a system, we propose to use Bloom filters, which are probabilistic data structures that support set operations, in combination with homomorphic encryption, which is a type of encryption that allows performing operations on encrypted data. We envision a system that provides reliable counts of travelers moving between locations as the only retrievable information [9].

In this work, we explain and briefly evaluate our privacy-preserving method for its accuracy in counting travelers moving between locations, with the aim to show its principal working. As a case study, we consider a subway network where travelers utilize smart-card technology to check in and out of the transportation system.

## 2 System model

Our example proof-of-concept assumes a subway network with an automatic fare collection system. Subway networks usually consist of lines that connect specific origin and destination stations. For each station  $A$ , we assume there is a set of  $n_A$  scanners  $S_A = \{s_1^A, \dots, s_{n_A}^A\}$ , which are used by travelers to check in and out. In our model, we trust the sensors, but not the centralized server. For this reason, we first let a sensor collect detections to then send this collection in encrypted form to the server. The time during which detections are collected and aggregated before sending them to the server is called an **epoch**. Typically, an epoch lasts 5 minutes. As we will discuss in detail below, the server can operate on the encrypted collections of detections, but cannot reconstruct individual detections themselves.

A scanner  $s \in S_A$  reads a card's unique identifier  $cid$ . Each card reading belongs to an epoch  $e \in \mathcal{E}$  corresponding to its timestamp  $t$ , such that  $t_{start}(e) \leq t < t_{end}(e)$ , where  $t_{start}$  and  $t_{end}$  mark the beginning and the end of an epoch and  $\mathcal{E}$  denotes the set of all such epochs. A detection is thus a triplet  $(cid, s, e)$ , representing a card uniquely identified by its identifier  $cid$ , read by scanner  $s$  during epoch  $e$ . By  $\mathcal{D}_{s,e}$ , we denote the set containing all the identifiers detected by a scanner  $s$  during an epoch  $e$ . Let  $\mathcal{D}_e^A$  denote the set of all identifiers detected by any scanner at  $A$  during epoch  $e$ :  $\mathcal{D}_e^A = \cup_{s \in S_A} \mathcal{D}_{s,e}$ .

Using collections of detections provides a powerful mechanism for counting travelers. One simple example is that the size of a set  $\mathcal{D}_{s,e}$  indicates the number of

travelers who passed the scanner  $s$  during the epoch  $e$ . More interestingly, for two different stations, the size of the set  $\mathcal{D}_{e_1}^A \cap \mathcal{D}_{e_2}^B$  represents the number of travelers who were first detected at  $A$  during epoch  $e_1$  and subsequently detected at  $B$  during epoch  $e_2$ , where  $e_2$  occurs after  $e_1$ .

## 3 Method

### 3.1 Bloom filter

The problem with using sets of detections is that they still contain the card identifiers for anyone to see who has access to those sets. This issue can be addressed by using a *representation* for sets, called **Bloom filters** [2]. A Bloom filter has the property that it allows only for membership tests. In other words, the only way to discover which card identifier is stored, is to go over the entire list of possible card identifiers and check for each one of them which identifier the membership tests succeed. Although this already poses a potentially tremendous computational burden for discovering detected identifiers, it is not enough to prevent finding identifiers. To understand how encryption, combined with Bloom filters, can prevent such a discovery, we must first explain what they are.

A Bloom filter is implemented as a binary vector of  $m$  bits, initially all set to zero. Adding an element to the set involves hashing it with  $k$  different hash functions, each returning a position in the vector. Those bits are then set to 1. To determine whether an element is in the set, the same hash functions are applied, and the corresponding bits in the vector are checked. When each bit is also 1, the element is considered to be in the set. An important observation is that there is a chance that two different elements will see exactly the same bits being set to 1. As a consequence, a membership test may return a *false positive*: the element for which the test is computed is factually *not* in the set represented by the Bloom filter. It is for this reason that Bloom filters are said to be probabilistic data structures. Given the maximum acceptable probability  $p$  for false positives, along with the desired number  $n$  of elements to be stored, one can compute the minimal length  $m$  of a Bloom filter, as well as the minimal number  $k$  of hash functions to use:  $m = -\frac{n \cdot \ln p}{(\ln 2)^2}$  and  $k = \frac{m}{n} \cdot \ln 2$ .

The size of the set represented by a Bloom filter (i.e., its *cardinality*  $c$ ) can be estimated when knowing only  $k$ ,  $m$ , and the number  $t$  of bits that are set to 1 [7]:

$$c = -\frac{m}{k} \ln \left( 1 - \frac{t}{m} \right) \quad (1)$$

In addition to membership testing, Bloom filters also support union and intersection operations. An intersection of two sets  $\mathcal{D}_A$  and  $\mathcal{D}_B$  can be done by taking their respective Bloom filter representations and conducting a bitwise AND operation. To illustrate, if  $A$  is represented by  $[0, 1, 1, 0, 1]$  and  $B$  by  $[1, 1, 1, 0, 0]$ , then  $A \cap B$  is represented by  $[0, 1, 1, 0, 0]$ . A union is computed through a bitwise OR operation. (Note that for realistic representations of sets, Bloom filters generally have lengths of 1000s of bits.) Whereas unions do not affect the probability of false detections, intersections do. This also means that estimating the size of an intersection when using Bloom filters may easily see deviations. We ran extensive tests and encountered estimates that were 15% off the real size. A more accurate estimation for two intersecting sets is provided by [7], yet no general estimation is known for more than two intersecting sets. For this paper, we will use the simple approximation given by Equation 1.

Ignoring encryption for the moment, detections at a scanner  $s$  are converted into Bloom filters and sent by  $s$  to the server at the end of each epoch. To answer queries, the server may do a series of unions and intersections on various Bloom filters, as we explained in our simple example above. The result is a Bloom filter representing the detections related to the query. At that point, the server could return the estimated cardinality of the set. Unfortunately, the server itself can still, with some computational effort, discover the detected card identifiers. As we mentioned, in our system model, we do not trust the server. This is where encryption comes into play.

### 3.2 Homomorphic Encryption

To prevent the server from discovering identifiers, Bloom filters must be combined with encryption schemes. Homomorphic encryption [8] is a specialized form of encryption that enables mathematical operations to be conducted directly on encrypted data without the need for decryption. The results of these operations are also encrypted, and the output is the same as if the operations had been conducted on unencrypted data.

The following procedure is now followed using homomorphic encryption. Suppose a user  $U$  is interested in knowing how many travelers moved from  $A$  to  $B$ . To that end, she passes an *encryption key* to the server, which is then used to encrypt all Bloom filters from the moment the key is available (note that this means that a user cannot issue queries that relate to the past, i.e., the time before they made the encryption key available). Also note that the user holds the *decryption key*, and is thus the only entity who can decrypt the corresponding encrypted

Bloom filters. Neither the scanners nor the server can decrypt those Bloom filters.

The server now operates on bitwise encrypted Bloom filters and produces a final result, say an encrypted Bloom filter  $BF$  representing a set  $R$ . By simply *adding* the entries of  $BF$ , it can produce an (encrypted) version  $t^*$  of  $t$ , the number of bits that have been set to 1. This value, along with  $k$  and  $m$  can then be handed over to the user  $U$ , who can decrypt  $t^*$  and compute the cardinality  $c$ . The server can also hand out  $BF$  to the user, but not after having shuffled the entries (otherwise, the user could still decrypt  $BF$  and discover detections). Shuffling keeps the same number of (encrypted) bits that have been set to 1, but a shuffled version of  $BF$  has no relationship to  $R$  anymore.

## References

- [1] Kay W Axhausen, Andrea Zimmermann, Stefan Schönfelder, Guido Rindsfuser, and Thomas Haupt. Observing the rhythms of daily life: A six-week travel diary. *Transportation*, 29(2):95–124, 2002.
- [2] Burton H Bloom. Space/time trade-offs in hash coding with allowable errors. *Communications of the ACM*, 13(7):422–426, 1970.
- [3] Yves-Alexandre De Montjoye, César A Hidalgo, Michel Verleysen, and Vincent D Blondel. Unique in the crowd: The privacy bounds of human mobility. *Scientific reports*, 3(1):1–5, 2013.
- [4] Merkebe Getachew Demissie, Santi Phithakkitnukoon, Titipat Sukhvibul, Francisco Antunes, Rui Gomes, and Carlos Bento. Inferring passenger travel demand to improve urban mobility in developing countries using cell phone data: a case study of senegal. *IEEE Transactions on intelligent transportation systems*, 17(9):2466–2478, 2016.
- [5] Yola Georgiadou, Rolf A de By, and Ourania Kounadi. Location privacy in the wake of the gdpr. *ISPRS International Journal of Geo-Information*, 8(3):157, 2019.
- [6] Dmytro Karamshuk, Chiara Boldrini, Marco Conti, and Andrea Passarella. Human mobility models for opportunistic networks. *IEEE Communications Magazine*, 49(12):157–165, 2011.

- [7] Odysseas Papapetrou, Wolf Siberski, and Wolfgang Nejdl. Cardinality estimation and dynamic length adaptation for bloom filters. *Distributed and Parallel Databases*, 28:119–156, 2010.
- [8] Ronald L Rivest, Len Adleman, Michael L Dertouzos, et al. On data banks and privacy homomorphisms. *Foundations of secure computation*, 4(11):169–180, 1978.
- [9] Valeriu-Daniel Stanciu, Maarten van Steen, Ciprian Dobre, and Andreas Peter. Privacy-preserving crowd-monitoring using bloom filters and homomorphic encryption. In *Proceedings of the 4th International Workshop on Edge Systems, Analytics and Networking*, pages 37–42, 2021.