

Privacy-Friendly Mobility Analytics using Aggregate Location Data

Apostolos Pyrgelis
University College London
apostolos.pyrgelis.14@ucl.ac.uk

Emiliano De Cristofaro
University College London
e.decrisofaro@ucl.ac.uk

Gordon J. Ross
University College London
gordon.ross@ucl.ac.uk

ABSTRACT

Location data can be extremely useful to study commuting patterns and disruptions, as well as to predict real-time traffic volumes. At the same time, however, the fine-grained collection of user locations raises serious privacy concerns, as this can reveal sensitive information about the users, such as, life style, political and religious inclinations, or even identities. In this paper, we study the feasibility of crowd-sourced mobility analytics over *aggregate* location information: users periodically report their location, using a privacy-preserving aggregation protocol, so that the server can only recover aggregates – i.e., how many, but not which, users are in a region at a given time. We experiment with real-world mobility datasets obtained from the Transport For London authority and the San Francisco Cabs network, and present a novel methodology based on time series modeling that is geared to forecast traffic volumes in regions of interest and to detect mobility anomalies in them. In the presence of anomalies, we also make enhanced traffic volume predictions by feeding our model with additional information from correlated regions. Finally, we present and evaluate a mobile app prototype, called Mobility Data Donors (MDD), in terms of computation, communication, and energy overhead, demonstrating the real-world deployability of our techniques.

1. INTRODUCTION

The availability of information about people’s locations and movements holds the promise to make urban planning more effective and efficient, and ultimately improve citizens’ quality of life. Prompted by the increased presence of always-on, always-connected devices, the ubiquitous collection of location information enables a number of interesting applications. New research frontiers, e.g., in the field of *anticipatory mobile computing*, make it increasingly possible to use mobile sensing along with machine learning for intelligent reasoning [32]. For instance, contextual location information collected from mobile users can be used to predict future mobility events [26, 39], detect mobility anomalies [30] and enable real-time traffic or event statistics [4].

At the same time, however, large-scale collection of individual users’ fine-grained locations raises serious privacy concerns, as this can reveal sensitive information about the users, such as, life style, political and religious inclinations, or even identities [31, 24]. Although often advocated, anonymization of location traces is moot as these reveal home/work locations, which in turn can be used to re-identify users [18]. In fact, just a few locations are enough to re-identify users [43]. Therefore, in this paper, we set to investigate whether or not mobility analytics can be effectively and efficiently performed over aggregate data. We turn to cryptographic protocols for privacy-friendly data aggregation and use them to privately gather location statistics [9, 23, 28, 34]. Overall, we aim to demon-

strate: (1) the usefulness of mobility analytics over aggregate locations, and (2) the real-world deployability of a scalable system for privacy-friendly location data collection.

Roadmap. We present a crowd-sourced system for privacy-friendly mobility analytics whereby users periodically report locations, but do so using a privacy-preserving aggregation protocol, so that only aggregates can be recovered (i.e., how many but not which users were in a region at a given time). We experiment with real-world mobility datasets obtained from the Transport For London (TFL) authority as well as the San Francisco Cabs (SFC) network, and present a methodology based on time series modeling geared to forecast traffic volumes in regions of interest (ROIs) and to detect mobility anomalies in them. In the presence of anomalies, we also make enhanced traffic volume predictions (achieving up to 50% improvement) by training our model with additional information from correlated regions. Such tasks are particularly useful in modern cities for journey planning [4, 26] and congestion prevention [38]. Finally, we show how to build a privacy-respecting system for data collection. To this end, we present a mobile app prototype, called Mobility Data Donors (MDD), and present an empirical evaluation of its computation, communication, and energy complexities, which attest to the practicality of our vision.

Paper Organization. Next section introduces a few concepts and tools used in our work, then, Section 3 presents our datasets and our methodology for predictive mobility analytics. In Section 4, we discuss the details of our proposed framework and analyze its real-world deployment. Finally, after reviewing related work in Section 5, the paper concludes in Section 6.

2. PRELIMINARIES

2.1 Auto Regressive Moving Average

As we aim to perform analytics on aggregate locations – specifically, predicting traffic volumes as well as detecting mobility anomalies in a Region Of Interest (ROI), such as underground stations (cf. Section 3.2–3.3) – we model ROIs’ time series using Auto-Regressive Moving Average (ARMA). We build on the work by Box et al. [8], who present an iterative method for choosing and estimating ARMA models.

Given a time series Y_t , an ARMA model is a tool for understanding and predicting future values in Y_t . The model is usually denoted as $ARMA(p, q)$, where $AR(p)$ denotes the autoregressive model of order p and $MA(q)$ refers to the moving average model of order q . Specifically, an $ARMA(p, q)$ model is defined as:

$$Y_t = c + \sum_{i=1}^p \phi_i \cdot Y_{t-i} + \epsilon_t + \sum_{i=1}^q \theta_i \cdot \epsilon_{t-i} \quad (1)$$

where c is a constant, ϕ_1, \dots, ϕ_p and $\theta_1, \dots, \theta_q$ are model parameters, and $\epsilon_t, \epsilon_{t-1}, \dots$ are white noise error terms.

2.2 Vector Auto-Regression

We also investigate how to improve traffic volume predictions in the presence of anomalies (cf. Section 3.4), thus, we also attempt to discover correlated ROIs and use their aggregate time series, along with a Vector Auto-Regression (VAR) model, to make enhanced predictions. VARs are statistical models used in econometrics to capture linear interdependencies among multiple time series, and consist a generalization of uni-variate autoregressive models (AR models) that allow more than one evolving variable. All variables in a VAR model are treated symmetrically and each of them has an equation explaining its evolution based on its own lags as well as those of the other model variables. VAR modeling requires the prior knowledge of a list of variables which can be hypothesized to affect each other inter-temporally.

A VAR model describes the evolution of a set of k variables (endogenous variables) over a sample period $t = 1, \dots, T$ as a linear function of their past values. The variables are collected in a vector y_t of size $(k, 1)$, whose i_{th} element y_{it} is the observation of the variable i at time t . A p -th order VAR model, denoted as $VAR(p)$ is given by the equation:

$$y_t = c + A_1 \cdot y_{t-1} + A_2 \cdot y_{t-2} + \dots + A_p \cdot y_{t-p} + e_t \quad (2)$$

where c is a vector of constants with size $(k, 1)$, A_i is a time-invariant matrix of size (k, k) and e_t is a vector of error terms with size $(k, 1)$ where: (a) $E(e_t) = 0$, every error term has mean zero, (b) $E(e_t e_t') = \Omega$, the co-variance matrix of error terms is Ω and (c) $E(e_t e_{t-k}') = 0$, for any non-zero k there is no serial correlation in individual error terms.

2.3 Spearman Correlation

To discover correlated ROIs, we will use Spearman's correlation coefficient, which is a non-parametric measure of the statistical dependence between the ranking of two variables [12]. It provides an estimate of how well the relationship between two variables can be described with a monotonic function and, unlike Pearson, it does not assume that both variables are normally distributed. Given two variables W, Z , the Spearman correlation coefficient is defined as:

$$r_s = 1 - \frac{6 \cdot \sum d_i^2}{n \cdot (n^2 - 1)} \quad (3)$$

where $d_i = rg(W_i) - rg(Z_i)$ is the difference between the two ranks of each observation and n is the number of observations. Similar to other correlation measures, Spearman's obtains values between -1 and $+1$, with 0 implying no correlation, and -1 or $+1$ implying an exact monotonic relationship. Intuitively, positive correlations imply that as W increases, so does Z , while negative correlations mean that as W increases, Z decreases.

2.4 Privacy-Preserving Data Aggregation

We also use cryptographic protocols for privacy-preserving data aggregation, allowing an untrusted aggregator to gather statistics (e.g., sum or mean) from users in such a way that data of single users is not revealed in the clear, but only the aggregate information can be recovered. These protocols are often used for smart metering [25], participatory sensing [34], or recommender systems [28].

Typically, private aggregation relies on a cryptosystem that is additively homomorphic: users send encrypted data to the aggregator, which does not hold the corresponding decryption key and cannot access *single* users' contributions, however, it can decrypt the sum of *all* users' reports. Specifically, we choose a protocol recently

proposed by Melis et al. [28], as it guarantees: scalability, independence from trusted third parties and/or key distribution centers, and fault tolerance. Scalability is achieved by combining the private aggregation protocol of Kursawe et al. [25] (secure under the Computational Diffie Hellman assumption in the presence of honest-but-curious adversaries) with data structures supporting succinct data representation, i.e., Count-Min Sketches [13]. These introduce a small, upper-bounded error in the aggregation, but reduce the computational/communication complexities of the cryptographic operations from linear to logarithmic in the size of the input. It also features a completely distributed key generation/distribution which, unlike other protocols, e.g. [22, 34], does not require any other authorities. Finally, its fault tolerance protocol addresses one of the main limitations of [25], i.e., if one or more users fails to report their (encrypted) data, the aggregator cannot correctly decrypt the aggregate (since it relies on encryption keys summing up to zero).

Melis et al. [28]'s protocol consists of four phases. **(1) Setup:** Assuming a cyclic group \mathbb{G} of order q for which the Computational Diffie-Hellman problem is hard, and g a generator of this group, each user $U_i \in \mathbb{U} = \{1, \dots, N\}$ generates a private key $x_i \in_r \mathbb{G}$ (i.e., sampled at random from G) and a public key $y_i = g^{x_i} \bmod q$. The public keys are published with the aggregator. **(2) Encryption:** Each user U_i holds an input vector of data points $S = \{S_c \in \mathbb{N}, c = \{1, \dots, T\}\}$. To participate in the privacy-preserving aggregation each user needs to generate blinding factors based on the public keys of the other users in such a way that they all sum up to zero. At round s , for $l = 1, \dots, T$, user U_i computes $k_{il} = \sum_{j=1, j \neq i}^N H(y_j^{x_i} \| l \| s) \cdot (-1)^{i > j} \bmod q$, where H is a cryptographic hash function and $\|$ denotes the concatenation operator. Then, for each entry $\{S_{il}\}_{l=1}^T$, U_i encrypts S_{il} as $b_{il} = S_{il} + k_{il} \bmod 2^{32}$ and sends the resulting ciphertext to the aggregator. **(3) Aggregation:** The aggregator collects the ciphertexts from each user U_i and (obviously) aggregates them. More precisely, for $l = 1, \dots, T$ it computes $C_l = \sum_{i=1}^N b_{il} = \sum_{i=1}^N k_{il} + \sum_{i=1}^N S_{il} = \sum_{i=1}^N S_{il} \bmod 2^{32}$, where C_l denotes the l -th item of the input vector S . **(4) Fault Recovery:** If, during the aggregation phase, only a subset of users \mathbb{U}_{on} successfully submit data, the aggregator sends \mathbb{U}_{on} to each $U_i \in \mathbb{U}_{on}$ and U_i computes, for each $l = 1, \dots, T$, $k'_{il} = \sum_{j=1, j \neq i, j \notin \mathbb{U}_{on}}^N H(y_j^{x_i} \| l \| s) \cdot (-1)^{i > j} \bmod q$. Then each user U_i sends these values back to the aggregator who can now obtain the aggregate counts by computing $C'_l = (\sum_{i \in \mathbb{U}_{on}} b_{il} - \sum_{i \in \mathbb{U}_{on}} k'_{il}) \bmod 2^{32}$.

Groups. Another feature of [28] is the ability to dynamically allocate users in groups, and perform within-group aggregation and then combining statistics from multiple groups, which is crucial to cope with dynamic/mobile settings. It also allows to bound the complexity of the encryption phase, which depends on the number of users in the group.

Input Compression. As mentioned above, [28] uses Count-Min Sketches to guarantee scalability when the input vector (S) is large. Specifically, the encryption phase is modified as follows: each user U_i initializes a Count-Min Sketch vector $X_i \in \mathbb{N}^{d \times w}$ with zero entries, then encodes his original input vector S using the update procedure of Count-Min Sketches [13] while employing the following pairwise hash function: $h(x) = ((a \cdot x + b) \bmod p) \bmod w$ for $a \neq 0$, b random integers modulo a random prime p . Then each user encrypts X_i as in the previously described encryption phase.

If $|S|$ denotes the size of the input vector S , its compact representation with a Count-Min Sketch has size $O(\log(|S|))$. More precisely, given the sketch parameters (ϵ, δ) , the Count-Min Sketch is a vector of size $L = d \times w$ where $d = \lceil \ln(|S|/\delta) \rceil$ and $w = \lceil e/\epsilon \rceil$. For instance, if $\epsilon = \delta = 0.01$, a vector S of size $|S| = 10^4$ can

be encoded as a sketch of size $L = 3,808$, while a vector S of size $|S| = 10^6$ can be represented as a sketch of size $L = 5,168$. Obviously, the Count-Min succinct structure introduces an accuracy error and its parameters (ϵ, δ) give an upper bounded error for the estimated counters \hat{c}_i , amounting to $\hat{c}_i \leq c_i + \epsilon \cdot \sum_j |c_j|$ with probability $1 - \delta$ (with c_i being the true element of the vector).

3. MOBILITY ANALYTICS USING AGGREGATE LOCATIONS

We now present and evaluate our “mobility analytics” algorithms, specifically, predicting traffic volumes at ROIs, discovering and predicting anomalies – all using aggregate location reports. We rely on two real-world datasets obtained, respectively, from Transport for London (TFL) and the San Francisco Cab (SFC) network.

3.1 Datasets

3.1.1 Transport For London (TFL)

London’s transportation system consists of various connected subsystems: London Underground Ltd (LUL), London Transport Buses (LTB), Docklands Light Rail (DLR), Overground (LRC), Tramlink (TRAM), and National Rail (NR), operating in the city under the umbrella of Transport for London (TFL). The most common payment method for TFL fares is the Oyster Card, a pre-paid, RFID-enabled card. We have obtained from TFL data corresponding to all March 2010 trips from all (anonymized) oyster cards, which we pre-process in the following way. First, we discard trips from TRAM due to scarce density and LTB for consistency as travelers only tap-in but do not tap-out for bus trips paid by Oyster. Then, to observe weekly patterns, we focus on the four weeks from Monday March 1 to Sunday 28, 2010. The final dataset consists of approximately 60 million oyster-card trips, performed by almost 4 million unique users, over 582 stations. Each entry in the data describes a unique trip and consists of the following fields: *oyster id*, *start time*, *start station id*, *end time*, and *end station id*. Note that the time resolution of the timestamps is 1 minute.

Next, we aggregate single-trip records by grouping trips start and end times in time epochs of 1 hour, aiming to achieve regularity in transit patterns (similar to [45]), and count the number of passengers that entered (“tap-in”) or exited (“tap-out”) each station during a slot. For each station $i \in \{1, \dots, n\}$ (with n being the total number of stations), we create a time series Y_{it} indicating how many passengers transited through it in a time epoch $t \in \{1, \dots, m\}$ (m denotes the total number of epochs, i.e., 672): $Y_{it} = Y_{it}^{in} + Y_{it}^{out}$, where Y_{it}^{in} indicates the number of tap-in events and Y_{it}^{out} the number of tap-out events, at station i during epoch t .

In Figure 1, we plot the *hourly* aggregate time series of two stations – Canary Wharf (one of the busiest stations of London) and Clapham Common (a moderately busy station) – showing different patterns during weekdays and weekends, as well as peak commuting hours. In general, we note some weekly/daily seasonality in the stations’ time series as well as stationarity (i.e., no particular trend). We verify the latter by performing the Augmented Dickey-Fuller test [14] which indicates that 93% of tube stations have stationary time series with 95% confidence.

3.1.2 San Francisco Cab (SFC)

We also use the San Francisco Cab (SFC) dataset [33], which contains mobility traces recorded by taxis in San Francisco, between May 17 to June 10, 2008. The dataset contains approximately 11 million GPS coordinates, generated by 536 taxis. To observe weekly patterns in our data we sample the dataset to cover

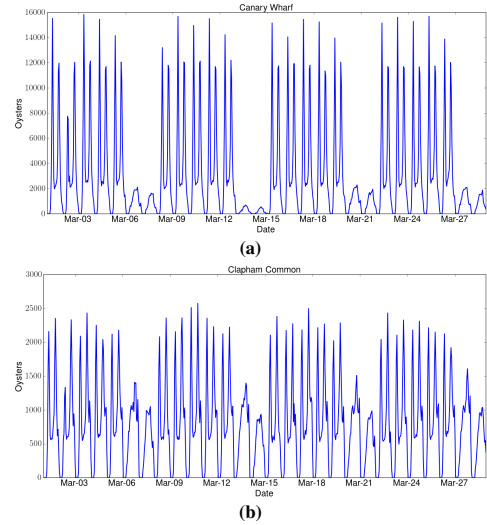


Figure 1: Hourly traffic volume at two TFL stations (March 1–28, 2010).

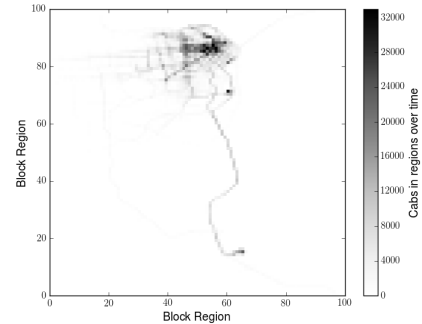


Figure 2: Number of cabs on 100×100 SF grid (May 19 – June 8, 2008).

exactly 3 weeks of data: Monday May 19 to Sunday June 8, 2008. Entries in the dataset include the following fields: *cab identifier*, *latitude*, *longitude* and a *time stamp* in UNIX epoch format.

We follow a similar approach as with the TFL dataset to aggregate the traces, however, since locations are GPS coordinates rather than points of interest, we divide the city of San Francisco into a grid S consisting of 100×100 regions, each covering an area of 0.19×0.14 square miles. We group the GPS traces in one-hour epochs and, for each region $i \in \{1, \dots, n\}$ (with n being the total number of regions, i.e., 10,000), we count the number of taxis that have reported a presence in that block during epoch $t \in \{1, \dots, m\}$ (m being the number of time epochs, i.e., 504), and create a time series Y_{it} as: $Y_{it} = \sum_{j=1}^k p_{jt}$, where k is the total number of taxis (i.e., 536) and $p_{jt} \in \{0, 1\}$ indicates whether taxi $j \in \{1, \dots, k\}$ reported its location at region i during epoch t .

We aggregate the traffic from our 3-week dataset for each region and, in Figure 2, plot the resulting heatmap. Unsurprisingly, the downtown area exhibits the highest traffic volume, with the route to/from SFO airport also clearly visible. In Figure 3, we plot the aggregate time series of two regions, one of the busiest (id = 7160) and a moderately busy one (id = 8554). Once again, weekly and daily patterns can be observed, along with stationarity (96 out of the 100 busiest regions have stationary time series with 99% confidence as indicated by the Augmented Dickey-Fuller test).

3.1.3 Removing Seasonality

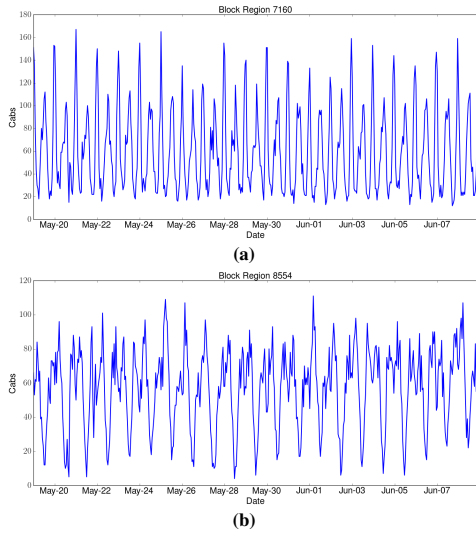


Figure 3: Hourly traffic volume in regions 7160 and 8554 of SFC dataset.

Our preliminary analysis of both datasets shows that aggregate time series of the ROIs (tube stations or regions) exhibit no particular trend but do preserve weekly/daily seasonality. Therefore, as proposed in prior work, e.g., [21], we de-seasonalize each region’s time series via additive decomposition. More specifically:

$$D_{it} = Y_{it} - \bar{Y}_{it} \quad (4)$$

where Y_{it} is the ROI’s original time series i and \bar{Y}_{it} is its seasonality defined as $\bar{Y}_{it} = \frac{1}{w} \sum Y_{idh}$, with w being the number of weeks in the dataset, d the day of the week (i.e., $d \in \{\text{Monday}, \dots, \text{Sunday}\}$), and $h \in \{0, \dots, l\}$ with l denoting the number of epochs in one day (24 since we aggregate hourly). Observe that \bar{Y}_{it} is a time series containing the average value of each specific time slot (e.g., Mondays 3pm – 4pm). As an example, Figure 4 shows Green Park station’s (a station among the busiest TFL stations) aggregate time series in both its original and de-seasonalized form. Note a negative spike on the morning hours of March 8, as the station must have probably had reduced access (e.g. due to partial closure). In general, the de-seasonalized ROIs’ time series show strong auto-regressive structure.

3.2 Predicting Traffic Volumes in ROIs

We now investigate how to make hourly traffic volume predictions on ROIs using the TFL and SFC time series. Such predictions are particularly useful in modern cities for journey planning [26, 39], congestion prevention [38] as well as improving transportation service levels and adjusting staff needs at stations [2].

We focus on the 100 busiest TFL stations and the 100 most popular SFC regions. Since our preliminary analysis shows that ROIs’ time series are stationary and exhibit strong auto-regressive structure, we turn to ARMA modeling (cf. Section 2.1). For each ROI i , we feed the ARMA model with the values of the last 6 days of its aggregate and de-seasonalized time series Y_{it}, D_{it} , respectively. We train the model using the first 5 days of D_{it} and test it against the last (“test day”) of Y_{it} following a recursive approach with a sliding time window to predict its hourly traffic. To do so, for each time slot we combine ARMA model’s predictions on D_{it} with ROI’s seasonality \bar{Y}_{it} , therefore, our predictions are given by $\hat{Y}_{it} = \hat{D}_{it} + \bar{Y}_{it}$, where \hat{D}_{it} is the ARMA prediction on the ROI’s de-seasonalized time series and \bar{Y}_{it} its seasonality.

We then compare our approach against a baseline, i.e., a black

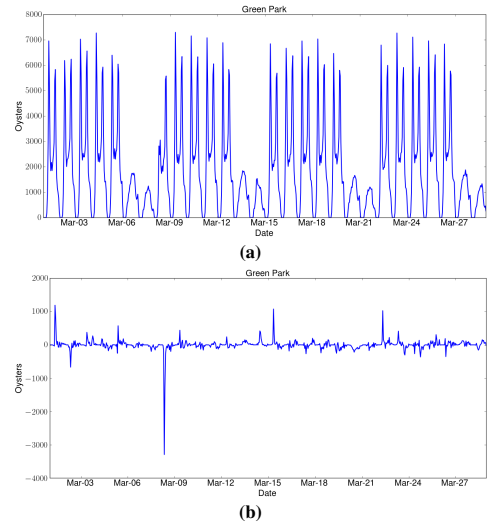


Figure 4: Green Park station’s time series without (a) and with (b) de-seasonalization.

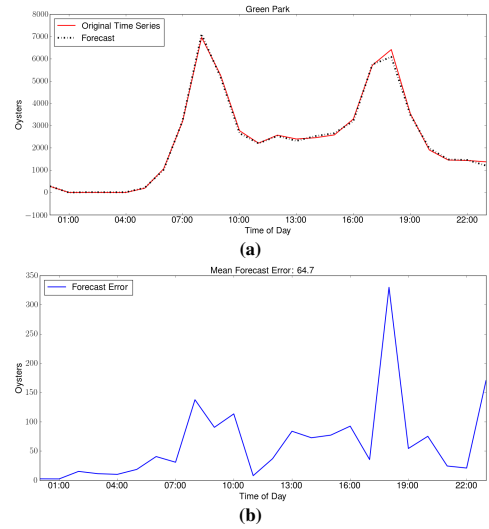


Figure 5: Hourly traffic forecasts for Green Park station on March 25.

box approach where we fit the ARMA model directly on each ROI’s i time series Y_{it} . We evaluate the accuracy of the predictions using the absolute forecast error as $e_{it} = |Y_{it} - \hat{Y}_{it}|$, where Y_{it} is the actual time series value at time slot t (ground truth) and \hat{Y}_{it} is the predicted value for that time slot using our approach. We also convert the error into a percentage error, i.e., $p_{it} = \frac{e_{it}}{Y_{it}} \times 100$. Figure 5(a) plots the traffic volume forecast for Green Park station on March 25, while Figure 5(b) shows the absolute forecast error. Overall, on the TFL dataset, the mean absolute forecast error for March 25, over the 100 busiest stations, is 59.53 ± 42.48 oysters, compared to 545.9 ± 376.8 oysters with the baseline. This corresponds to an error of $19.6\% \pm 59.5\%$ vs $638\% \pm 1619\%$, showing that the seasonality-based method significantly outperforms the baseline.

We follow the same approach for the SFC dataset, predicting the traffic volume of the most popular regions. Figure 6 shows the predictions and the forecast error for the region with identifier 8755, on June 5. The average forecast error over the 100 busiest regions is 5.62 ± 3.12 taxis ($19.7\% \pm 10.3\%$), whereas, the baseline

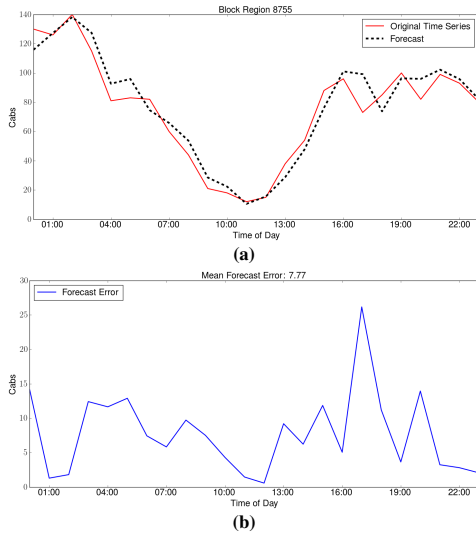


Figure 6: Hourly traffic forecasts for SFC’s region 8755 on June 5.

error is 9.07 ± 2.95 ($35.4\% \pm 13.4\%$), once again showing that predictions can be improved when considering seasonal effects.

3.3 Detecting Traffic Anomalies

Next, we focus on detecting traffic volume anomalies on ROIs’ time series. This is particularly important for traffic provisioning and travel planning as trip recommendations can be made to drivers/commuters during road accidents, incidents or events that cause overcrowding in transportation stations [3, 4, 30].

Once again, we utilize ARMA modeling: our intuition is to train the model for each ROI and rely on the absolute forecast error as a confidence interval for detecting anomalies. More precisely, since the forecast error is normally distributed, we apply the 3σ rule and set upper and lower confidence intervals λ_1, λ_2 . Thus, we detect an anomaly when: $e_t > \lambda_1 \parallel e_t < \lambda_2$, where \parallel is the logical OR operator, $\lambda_1 = \mu + 3\sigma$, $\lambda_2 = \mu - 3\sigma$, and μ, σ are, respectively, average and standard deviation of the forecast error e_t . In a way, we flag as anomalies time slots that our model could not predict with good accuracy.

Subsequently, we experiment with our anomaly detection technique using a similar approach to that described in Section 3.2: we train the ARMA model using data of the first week of each ROI’s time series and test it against the remaining weeks (i.e., for the TFL dataset we have 3 test weeks while for the SFC we have 2), using a sliding window, aiming to identify traffic volume anomalies. We focus on the 100 busiest stations of the TFL dataset and detect 896 anomalies, which roughly corresponds to 2% of all 1-hour slots in the 3 test weeks. In the SFC dataset, over the 100 busiest regions, we find 366 anomalies (i.e., 1% of the 2 test week time slots). We rank each anomaly based on its deviation from the confidence intervals as a measure of its magnitude and keep track of the top 10% of anomalies, i.e., 90 anomalies for the TFL dataset and 30 for SFC: in Section 3.4, we will investigate whether or not we can enhance traffic prediction in the presence of anomalies by combining information from correlated ROIs.

Note that we do not have *ground truth* as to what constitutes an actual “anomaly” in our datasets, so we cannot empirically evaluate how well our approach corresponds to detecting, e.g., events, strikes, disruptions, etc., and anyway this would be out of the scope of our work. In general, we consider an anomaly to be a pattern that does not conform to expected normal behavior [10] and, as such,

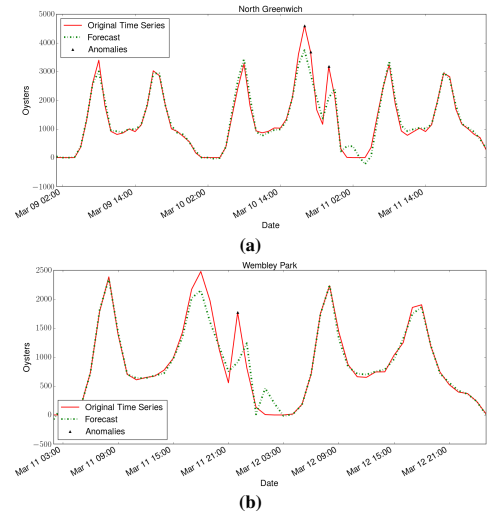


Figure 7: Original time series, forecasts, and detected anomalies in (a) North Greenwich station, March 9–11, and (b) Wembley Park station, March 11–12.

our anomaly detection techniques really consist in automatically *flagging* such patterns using aggregate locations. By focusing on the top events in terms of deviation from the confidence intervals (i.e., unexpectedly increased/decreased traffic patterns in ROIs), we aim to investigate whether collecting information from multiple ROIs can improve traffic volume predictions in the presence of anomalies (see Section 3.4).

We discuss some case studies among the top anomalies that we were able to correlate with external events. As shown in previous work [36], distinct human mobility patterns are observed during events that attract big crowds like football matches or music concerts. In Figures 7(a) and 7(b), we plot “anomalies” we detect in North Greenwich and Wembley Park tube stations, on the evening hours of March 10 and March 11. This seems to correspond to concerts taking place in the O2 and Wembley arenas, which are venues close to those stations. These events likely cause increased traffic spikes at nearby stations. Similarly, Figure 8(a) shows the original aggregate time series of Arsenal station as well as the anomalies detected on it when fitting our model. We can observe that the model detects anomalies on the evening of March 20, when an Arsenal FC soccer game was taking place. Finally, Figure 8(b) does the same for a region (id 8261) in the SFC dataset that is nearby AT&T Park, showing increased taxi traffic on the evening of May 31, when the San Francisco Giants were playing a baseball match.

3.4 Predicting Traffic Volumes in Case of Anomalies

We now investigate how to improve ROI traffic volume predictions in the presence of anomalies, using additional information from correlated ROIs. To this end, for each ROI, we use Spearman correlation (see Section 2) in order to discover those ROIs whose traffic can be useful for enhancing our predictions. Subsequently, we train a VAR model – geared to capture linear dependencies among multiple time series – with the time series of a ROI as well as the time series of its correlated ROIs and we compare the prediction results against a *local* model, i.e., an ARMA model trained only with ROI’s past local information (note that the ARMA model described in Sections 3.2 and 3.3 now consists our *baseline*). In the rest of this section, for each of our datasets, we

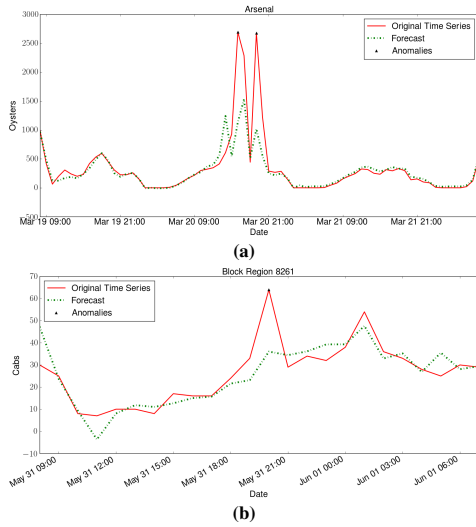


Figure 8: Original time series, forecasts, and detected anomalies in (a) Arsenal station, March 19–22, and (b) SFC region 8261, May 31–June 1.

start by describing our approach on a specific case study and then we generalize our results by focusing on the top events (90 for TFL and 30 for SFC, respectively) that our anomaly detection module has flagged as possible anomalies.

TFL. We first focus on Saturday March 20, when our anomaly detection module spots anomalies – i.e., increased traffic volume – on the Arsenal station, likely caused by an Arsenal FC soccer game. We zoom in on the two hours before and after the game (15:00–17:00, resp., 19:00–21:00, respectively) when the majority of Arsenal fans, exit from, resp., enter Arsenal station. We follow a similar aggregation approach as that described in Section 3.1 although now for each station i , we keep two separate time series: one counting passengers entering the station (Y_{it}^{in}) and one counting those exiting it (Y_{it}^{out}). Once again, we de-seasonalize each station’s entering/exiting time series as discussed in Section 3.1.3.

To discover stations correlated with Arsenal, we compute the Spearman correlation (see Section 2) between the de-seasonalized time series of passengers entering/exiting Arsenal as well as the de-seasonalized time series of all the remaining stations in the TFL network, by sliding them up to 1 hour earlier/later. Our results show that the traffic exiting at Arsenal is highly correlated with the traffic entering at various other TFL stations including Arncliffe Grove, King’s Cross, Leicester Square, Blackhorse Road and Cockfosters (i.e., stations on the same line with Arsenal or on a line connected with the line of Arsenal). We then set to improve the traffic volume predictions of passengers exiting at Arsenal station before the match, by feeding a model with the de-seasonalized entering time series of the correlated stations. To do so, we use a vector autoregression model (VAR) (see Section 2) which describes the evolution of a set of variables over the same sample period as a linear function of their past values. Figure 9 shows traffic volume predictions for passengers exiting Arsenal station using the ARMA model with local information (i.e., the station’s past exiting time series, which now consists our *baseline*) and using the VAR model enhanced with additional information from the 10 most correlated stations. We observe that the *enhanced* model makes significantly better predictions between 15:00–17:00 on March 20, where we observe increased traffic due to the game. We measure the average forecast error of the exit traffic predictions of Arsenal on that day as 133.9 ± 270.6 oysters (i.e. $93\% \pm 185\%$) when making predic-

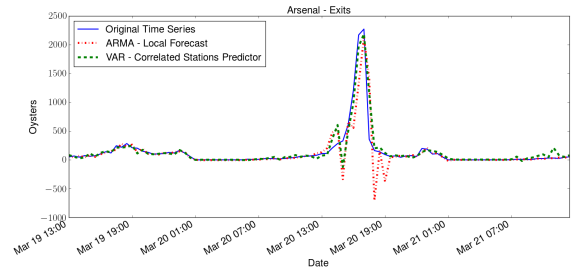


Figure 9: Forecast of passengers exiting Arsenal station March 19–21 based on a local ARMA model only vs. the VAR model with information from correlated stations.

tions with local information only and 65.26 ± 135.04 oysters (or $59.1\% \pm 57.4\%$) when enhancing it with the correlated stations.

In relation to the same event, we discover high correlation between the time series of passengers entering at Arsenal and those exiting at Arncliffe Grove a time slot later. Thus, we enhance the local ARMA prediction of traffic exiting volume at Arncliffe Grove by feeding the model with the de-seasonalized traffic entering Arsenal. Again, we use VAR to calculate the appropriate values for the parameters. The enhanced model makes better predictions between 19:00 and 21:00 where there was increased exiting traffic due to the Arsenal game. Indeed, the average forecast error for March 20 decreases from 23.38 ± 48.79 oysters to 11.22 ± 18.61 , a 52% average improvement in traffic volume predictions.

Next, we apply our approach for *all the top 90 anomalies* found in the TFL dataset (cf. Section 3.3), i.e., for each station under the presence of an anomaly, we predict its exiting or entering traffic (depending on which direction the anomaly has been detected) using a VAR model trained with the exiting or entering time series (again, depending on the anomaly direction) from 10 correlated stations. We measure the average forecast error over the day of the anomaly and we compare it against a local approach (ARMA model), where predictions are made using only past station’s information. Overall, for the 90 anomalies of the TFL dataset, we observe a $29\% \pm 13\%$ improvement in our predictions when employing the VAR model, indicating that sharing information between correlated ROIs improves prediction quality.

SFC. We follow a similar approach for anomalies detected on the SFC dataset. Our correlation analysis shows that the de-seasonalized time series of neighboring regions have high correlation, as it is likely that they are connected by the same roads. For instance, if we focus on the anomaly detected in region 8556 between May 27–28, we observe a 41% improvement in predictions when training a VAR model including additional information from 5 correlated regions (i.e., block regions with ids 8557, 8657, 8558, 8655 and 8555) compared to the baseline, i.e., the ARMA model that predicts using only local information. Similar to our TFL experiments, we generalize this approach by trying to improve the predictions for the *top 30 anomalies* of the SFC dataset, enhancing our model with information from 5 correlated regions. In this case, we obtain a $18\% \pm 14\%$ average improvement on the prediction.

Discussion. We observe that our techniques discussed above yield better improvements over the TFL dataset compared to SFC. A possible explanation arises from the nature of the TFL dataset and the way passenger trips are aggregated, i.e., at station level while preserving the notion of *direction* (number of passengers exiting or entering a station). This allows us to perform a more fine-grained correlation analysis in comparison to the SFC dataset, where as we aggregate the GPS locations, we lose the granularity of each taxi’s

trajectories (i.e., taxi moving from one region to another). While this is a good feature vis-a-vis privacy guarantees, it motivates the need to gather (privacy-friendly) aggregate location statistics while preserving directions.

4. A SYSTEM FOR PRIVACY-FRIENDLY MOBILITY ANALYTICS

After having assessed the usefulness of collecting and using aggregate locations for mobility analytics, we now set to investigate how to *enable* such collection in a privacy-friendly way. To this end, we design a distributed, collaborative framework whereby users install an application – called *Mobility Data Donors* (MDD) – that regularly monitors their locations, stores it locally, and periodically reports it to our server in a privacy-friendly way.¹ Privacy is guaranteed through aggregation, by means of the scalable private aggregation protocol presented in Section 2.4, thus, the server only learns aggregate information, i.e., how many (but not which) users were in a particular region or entered/exited a particular underground station in an interval of time. Once the server has received the aggregate location data (i.e., counts of users’ presence in ROIs), it can use it for mobility analytics applications.

As discussed earlier, protecting privacy of user locations is critical, as sensitive data about individuals, such as their religion [31] or their identity [18] can be inferred, and even a few locations are enough to re-identify users from anonymized traces [43]. The ability to privately collect location reports enables applications that would otherwise be impossible due to privacy concerns. For instance, obtaining data from TFL typically requires several rounds of NDAs and the promise not to re-distribute the data: although TFL could publish aggregate statistics, it is unlikely they would do so in real-time (a crucial aspect for mobility analytics) and anyway this would only capture one aspect of urban mobility—i.e., underground/overground trips but not, e.g., taxis or buses. In general, collecting locations directly from the users, without requiring them to forego their privacy, paves the way for a number of novel and interesting analytics, which we are confident our work will support.

4.1 Data Collection

To support private data collection, we need a secure aggregation protocol that allows a server to only learn aggregate locations. As discussed in Section 2.4, we choose the one by Melis et al. [28] as it supports scalability and fault-tolerance without the need for a trusted third party, which are fundamental factors for the success of a distributed, crowd-sourcing system.

Our system model mirrors to that of Melis et al. [28], i.e., it consists of a server, or aggregator, that facilitates networking and collects aggregate location counts from a set of mobile users running the MDD app. There is no other trusted authority. As in [28], we assume the aggregator and the users to be honest-but-curious, i.e., they follow protocol specifications and do not misrepresent their inputs, but try to extract information from other parties.

When installed, the MDD app generates a private/public key pair (see setup phase in Section 2.4) and communicates its public part to the aggregator. After setup, the app runs on the background, regularly collecting GPS coordinates. At predefined time slots (by default, every hour), the privacy-preserving aggregation is triggered by the server, provided that there are at least τ users connected, which are randomly assigned to groups of u users. In the default setting, the app maps GPS coordinates to a grid of $p \times p$ cells of ρ square miles, and builds a $p \times p$ matrix corresponding to the grid, setting to 1 items corresponding to ROIs the user has visited in the

specified time slot (and 0 otherwise).² The values of u , τ , and ρ are passed onto the user to inform them of the granularity of the data collection, and give them the option to withdraw (minimum acceptable values can be adjusted from the MDD’s settings). Next, as per [28], the app generates blinding factors (summing up to zero) based on the keys of the users in the same group, and encrypts each entry in the matrix. Finally, it sends the encrypted matrix to the aggregator who obliviously aggregates all (encrypted) matrices and decrypts the aggregate location counts used for the analytic tasks.

Besides recording coordinates, and mapping them onto a grid, the app can also recognize points of interest, such as train/underground stations, which is particularly useful for mobility analytics on transport datasets such as the TFL data. In this case, the aggregation takes place on a vector where each item corresponds to a point of interest and is set to 1 if the user has visited it in the specified time slot.

4.2 Experimental Evaluation

Next, aiming to assess the real-world deployability of our techniques, we empirically evaluate the performances of the MDD app, in terms of computation, communication, and energy overhead. Specifically, we evaluate the overhead imposed by the cryptographic operations needed for the privacy-preserving data collection. We use the mobility datasets from Section 3.1 as guidelines for simulating the system. For our experiments, we use the prototype implementation, in Javascript/Node.js, of the protocol by Melis et al. [28] and have adapted its client-side to run on Android using Apache Cordova.³ The cryptographic operations are implemented using elliptic curve cryptography, specifically, the Ed25519 elliptic curve [6] (supporting 256-bit points and offering 128-bit security) from the Elliptic.js library.⁴

For the sake of our evaluation, we run the experiments on a mid-range (rather than a high-end) Android device, as we do not want to limit deployment only to (possibly higher-income) users that have newer phones. We use a Samsung Galaxy A3 device, equipped with a 1.2 GHz quad-core Snapdragon 410 processor and 1.5GB RAM, running Lollipop v5.0.2. For our energy consumption analysis, we utilize PowerTutor [1], an Android app for power monitoring. Note that, although a Javascript implementation of the cryptographic operations might not be optimal in terms of efficiency (e.g., compared to Java), it offers portability among different mobile OSes. Anyway, we have actually benchmarked a Java implementation of the same operations and obtained similar results.

TFL. We start our experiments with the TFL use-case. Recall that the TFL data involves 582 ROIs (stations), so each user device, for each time slot t , encrypts a matrix of size $2 \cdot 582$, with the first row indicating entering the station and the second exiting it. We here remind that the complexity of the aggregation protocol depends on how many users are assigned to the same group, since the blinding factors are derived from public keys of other users in the group. In Figure 10(a), we plot the execution time, measured on our Android device, of the encryption phase vis-a-vis the number of users in the group. As expected, running times grow linearly in the size of the group. For instance, the encryption performed by each mobile device takes 4.2s with 100 mobile users and 42s with 1,000. Therefore, one should probably keep groups at around 200 users, which offers a reasonable trade-off between granularity (in terms of privacy) and efficiency. Obviously, even if, say, 1 million

²Note that the app is easily tunable so that, instead of binary values, the matrix encodes, e.g., duration of user’s presence in each ROI, whether the user has entered or exited a cell, etc.

³<https://cordova.apache.org/>

⁴<https://github.com/indutny/elliptic>

¹The app Android prototype is available upon request.

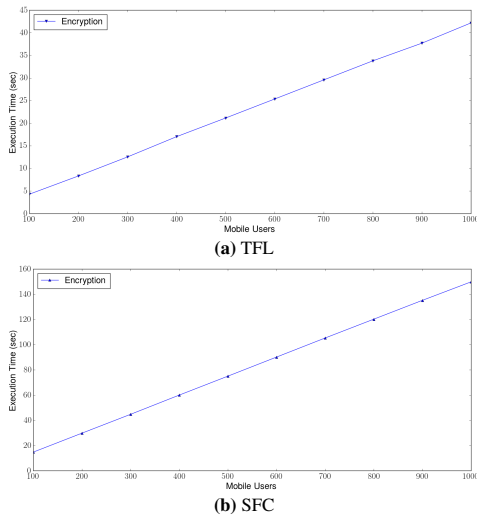


Figure 10: Execution times of the encryption phase for increasing number of users in (a) TFL and (b) SFC settings.

London commuters were to participate, the system simply scales by running multiple parallel instances with each group, and combining multiple aggregates from 5,000 200-user groups.

When assigning users to groups of size 200, in each round of the protocol, each device has to download $10.7KB$ worth of public keys, and transmitting encrypted values for $2 \cdot 582$ entries results in an overhead of $4.54KB$ (independently of how many users are in the group). Finally, we measure the energy consumption to be $862mJ$ for the encryption part, $609mJ$ to download public keys via Wi-Fi, and $322mJ$ to transmit the encrypted matrix. We observe that the energy overhead is quite small for a modern mobile phone (for reference we note that downloading a $20KB$ web page with Google Chrome via WiFi consumes approximately $800mJ$).

SFC. Next, we evaluate complexities considering the SFC use-case, for which we divide the city of San Francisco into a grid of 100×100 cells. In this setting, our mobile app, for each time slot t , encrypts a matrix of size 10,000. Figure 10(b) displays the resulting computational overhead for an increasing number of users. Once again, we observe execution times linear in the number of users in the system, i.e., the encryption phase requires about 14s with 100 and 149s with 1,000 users. We also measure communication and energy overhead, assuming groups of 100 users. Obtaining the public keys of users in the group requires downloading $5.37KB$, while transferring the encrypted 100×100 matrix requires $39KB$, which, once again, is acceptable for a mobile app. The energy consumption of the cryptographic operations in each protocol round is $485mJ$. Finally, the communication operations (via WiFi) require $306mJ$ and $2769mJ$, respectively.

Succinct Data Representation. We observe that the number of users in each group mainly affects the computation overhead, while the communication and energy overheads are primarily influenced by the size of the input. In particular, from our experiments above, we notice that the computation/communication/energy is appreciably low when groups are in the order of hundreds of users and matrices are in the order of thousands. This means that if aggregation is performed over larger inputs, the protocol would quickly incur high energy, communication, and computation overheads, and this would remarkably limit the deployability of our techniques.

In fact, even if in the use-cases considered in this paper the overhead is appreciably low, we would not be able to extend to, e.g.,

building origin-destination matrices [38, 45, 29], rather than only keeping counts. Origin-destination information is particularly useful to obtain finer grained statistics, e.g., discovering *similar* locations for personalized recommendations [11], and/or modeling the effects of disruptions in a transportation network as in [38]. More specifically, an origin-destination matrix for the TFL data would be of size 582×582 . At each time slot, the app sets element (i, j) in the matrix to 1 if the user commuted from station i to station j . For groups of 100 users, this would result in a $7.6s$ computation overhead (with $3,387mJ$ energy overhead) and a relatively high $1.5MB$ communication overhead at each protocol round (with $72,704mJ$).

However, as discussed in Section 2.4, we can use Count-Min Sketches to reduce complexities from linear to logarithmic in the size of the input. That is, we can compress the 582×582 origin-destination matrix into one of size 272×11 , yielding $4.6s$ computation (with $461mJ$), and $11.6KB$ (with $823mJ$) overhead when setting ϵ and δ parameters to 0.01. Similarly, for the SFC dataset, an origin-destination matrix would become very large, i.e., $10^4 \times 10^4$. Encrypting such a matrix would yield order of 1,000s computation overhead (and $1,000J$) and transferring it would introduce a $39MB$ overhead (and $2,835J$). Obviously, these numbers are prohibitive large for a mobile application. Whereas, with succinct data representation the matrix size could reduce to $1.9MB$ (with $138,137mJ$ transmission overhead), and a computation overhead of $18s$ ($5,422mJ$) when using a Count-Min Sketch with ϵ, δ parameters both set to 0.01.

Accuracy. Finally, we assess the accuracy loss introduced by the succinct data representation. We start by measuring the error in the aggregate counts on the TFL data, even though, as discussed above, the size of the input is sufficiently small that we do not actually need input compression. Regardless, we do so to determine the viability of using Count-Min Sketches, i.e., whether or not the error they introduce would reduce the effectiveness of the analytics, and the results are very encouraging. When setting the ϵ, δ parameters to the standard 0.01, we notice that the overall accuracy error in the resulting aggregate counts, over all time slots and all stations, is in the order of 1%. We also measure how the error “propagates” on the predictive analytics and, again, do not observe any statistically significant difference: for instance, when we make traffic forecasts for the 100 busiest stations on March 25 (as done in Section 3.2), we measure the average forecast error as 60.81 ± 39.62 using Count-Min Sketches and 59.53 ± 42.48 without it.

5. RELATED WORK

Mobility Analytics. Zhong et al. [45] propose a metric for capturing variability in commuting trips and analyze urban mobility patterns in London, Singapore, and Beijing, using one-week data of underground journeys. They show that regularity is exhibited when considering time intervals longer than 15 minutes, and demonstrate that peak hours are those with the least variability during a day. Although their analysis results provide useful insights for our work, the authors do not present any methodology for predicting mobility. Silva et al. [38] introduce a general framework for predicting traffic volumes in the London underground: they build a predictive model for each pair of stations under normal conditions (called the natural regime) and then extend it to model disruptions like station or line closures. Their approach is substantially different from ours since disruptions are actually part of a ground truth dataset they obtain from London’s transport authority.

Horvitz et al. [20] propose JamBayes, a probabilistic traffic forecasting system deployed in the Seattle Greater Area. They collect highway traffic data and contextual data (e.g., city events), and use

Bayesian structure search to model bottlenecks. Additionally, using data of historic traffic *surprises*, in combination with recent data before an event, their system learns Bayesian networks that infer the likelihood of a future surprise event. Garzó et al. [17] use distributed streaming algorithms to process large scale mobility data and make user mobility predictions on large metropolitan areas. They evaluate their location prediction methods on a 2-week mobility dataset obtained from the Orange D4D challenge [7]. Yavaş et al. [42] use data mining to predict user movements in a mobile computing system and evaluate their algorithms on a simulated mobility dataset. Overall, prior work on mobility analytics differ from ours as they do not consider collecting data directly from users (nor the privacy implications thereof).

Detecting Traffic Anomalies. In [30], Pan et al. combine mobility data along with social media to uncover the road network sub-graph associated with an anomaly, based on the routing behavior of drivers. Their system is evaluated on the Beijing taxi traces dataset. Similarly, Thom et al. [41] present a system geared to detect spatio-temporal anomalies by performing clustering on geolocated Twitter messages and visualize them using tag clouds. They experiment with three case-studies: an earthquake on the US East Coast, London riots, and hurricane Irene. Barria et al. [5] present an anomaly detection algorithm for road traffic using microscopic traffic variables like relative speed of vehicles, inter-vehicle time gap, and lane changing, and evaluate their approach using real-world video images from a highway segment in Bangkok. Zheng et al. [44] investigate whether collective detection of anomalies from multiple spatio-temporal datasets is possible. They propose a probabilistic anomaly detection method based on a spatio-temporal likelihood ratio test and evaluate it on five datasets from New York City. Sun et al. [40] build Markov models on user mobility patterns in a cellular network, aiming to detect intrusions in the network. Note that, although we also focus on identifying event mobility anomalies, unlike these works, we do so using aggregate crowd-sourced location data – specifically, collected directly from users in a privacy-preserving way.

Private Statistics. Prior work has also proposed a number of tools to privately gather location statistics, however, they do not demonstrate how these statistics can actually be used for performing mobility analytics, which is one of our main goals. Ho et al. [19] apply differential privacy to discover interesting geographic locations on aggregate location data, whereas [27] relies on synthetic data generation for publishing statistical information about commuting patterns. Brown et al. [9] propose Haze, a system for privacy-preserving real time traffic statistics based on jury voting protocols and differential privacy. Their system hides individual data while allowing aggregate information to be collected at the service provider. Similarly, Popa et al. [34] present PrivStats, a system for computing aggregate statistics over location data achieving privacy and accountability. Kopp et al. [23] also propose a framework enabling the collection of quantitative visits to sets of locations following a distributed approach. Shi et al. [37] show how an untrusted data aggregator can learn statistics over multiple participants’ private data using cryptographic techniques along with a data randomization procedure for achieving distributed differential privacy, while Melis et al. [28] demonstrate how to combine privacy-preserving aggregation with succinct data structures (Count-Min Sketches [13]) to efficiently compute statistics whilst provably protecting privacy of single data points. They also consider aggregating location information as a possible application of their protocols but do not perform any analytics. PASTE [35] introduces a solution in a similar setting whereby distributed differential privacy is used on time series data using a Fourier perturbation algorithm.

Finally, Fan et al. [15] propose FAST, an adaptive system for releasing real-time aggregate statistics with differential privacy. Their approach is based on a *trusted* central authority that adaptively samples the time series according to detected data dynamics to minimize the overall privacy budget. They employ Kalman filtering to predict data at non-sampling points and estimate the true values from perturbed ones at sampling in order to improve the accuracy of data release. In follow-up work [16], they present a generic differentially-private framework for anomaly detection on aggregate statistics, focusing on detecting epidemic outbreak: real-time aggregate data is perturbed using FAST [15] and released to an untrusted entity that performs the anomaly detection task. Whereas, we do not use differential privacy to protect users’ privacy, as this would require the presence of a trusted aggregator and introduce a trade-off between privacy and utility that is challenging to tune.

6. CONCLUSION

This work investigated the feasibility of performing crowd-sourced mobility analytics over aggregate location data, in a setting where users periodically report locations to a server, in such a way that the server can only recover aggregates, thanks to the use of a privacy-preserving aggregation protocol. We experimented with real-world mobility datasets obtained from the Transport For London authority as well as the San Francisco Cabs network, and demonstrated that aggregate location data can be useful for predictive analytic tasks like forecasting traffic volumes in *regions of interest* (ROIs) and detecting anomalies in them, using a methodology based on time series modeling with seasonality. In the presence of traffic anomalies, we also showed how to enhance their traffic volume predictions using additional information from correlated ROIs. Finally, we proposed a privacy-respecting system for data collection, and prototyped a mobile application – Mobility Data Donors (MDD) – which we empirically evaluated in terms of computation, communication, and energy overhead.

As part of future work, we plan to evaluate our methodology on different location datasets as well as perform a thorough (differential) privacy analysis of releasing datasets composed of aggregate locations, focusing on group sizes and semantic characteristics of ROIs such as size and density, and their evolution over time.

Acknowledgments. The authors wish to thank Luca Melis, George Danezis, Mirco Musolesi, and Licia Capra for their useful feedback and/or help with the experiments and the datasets. This research is supported by a Xerox University Affairs Committee grant on “Secure Collaborative Analytics.”

7. REFERENCES

- [1] Power Tutor – A power monitor for Android-Based Mobile Platforms. <http://ziyang.eecs.umich.edu/projects/powermentor/>.
- [2] TFL Budget Report 2015-2016. <http://content.tfl.gov.uk/board-20150326-part-1-item08-tfl-budget-2015-16.pdf>.
- [3] Transport For London - Travel Alerts. <https://twitter.com/tfltravelalerts?lang=en>.
- [4] Waze. <https://www.waze.com>.
- [5] J. A. Barria and S. Thajchayapong. Detection and classification of traffic anomalies using microscopic traffic variables. *IEEE Transactions on Intelligent Transportation Systems*, 2011.
- [6] D. J. Bernstein, N. Duif, T. Lange, P. Schwabe, and B.-Y. Yang. High-speed high-security signatures. *Journal of Cryptographic Engineering*, 2012.
- [7] V. D. Blondel, M. Esch, C. Chan, F. Cl  rot, P. Deville, E. Huens, F. Morlot, Z. Smoreda, and C. Ziemlicki. Data for

Development: The D4D Challenge on Mobile Phone Data. *arXiv 1210.0137*, 2012.

- [8] G. E. Box, G. M. Jenkins, G. C. Reinsel, and G. M. Ljung. *Time Series Analysis: Forecasting and Control*. John Wiley & Sons, 2015.
- [9] J. W. Brown, O. Ohrimenko, and R. Tamassia. Haze: Privacy-preserving real-time traffic statistics. In *ACM SIGSPATIAL*, 2013.
- [10] V. Chandola, A. Banerjee, and V. Kumar. Anomaly detection: A survey. *ACM computing surveys (CSUR)*, 2009.
- [11] M. Clements, P. Serdyukov, A. P. de Vries, and M. J. Reinders. Personalised travel recommendation based on location co-occurrence. *arXiv 1106.5213*, 2011.
- [12] G. W. Corder and D. I. Foreman. *Nonparametric Statistics: A Step-by-Step Approach*. John Wiley & Sons, 2014.
- [13] G. Cormode and S. Muthukrishnan. An Improved Data Stream Summary: The Count-Min Sketch and Its Applications. *Journal of Algorithms*, 2005.
- [14] D. A. Dickey and W. A. Fuller. Distribution of the estimators for autoregressive time series with a unit root. *Journal of the American Statistical Association*, 1979.
- [15] L. Fan and L. Xiong. Real-time aggregate monitoring with differential privacy. In *ACM CIKM*, 2012.
- [16] L. Fan and L. Xiong. Differentially private anomaly detection with a case study on epidemic outbreak detection. In *IEEE ICDMW*, 2013.
- [17] A. Garzó, A. A. Benczúr, C. I. Sidló, D. Tahara, and E. F. Wyatt. Real-time streaming mobility analytics. In *IEEE International Conference on Big Data*, 2013.
- [18] P. Golle and K. Partridge. On the anonymity of home/work location pairs. In *Pervasive Computing*. Springer, 2009.
- [19] S.-S. Ho and S. Ruan. Differential privacy for location pattern mining. In *SIGSPATIAL Workshop on Security and Privacy in GIS and LBS*, 2011.
- [20] E. J. Horvitz, J. Apacible, R. Sarin, and L. Liao. Prediction, expectation, and surprise: Methods, designs, and study of a deployed traffic forecasting service. *arXiv 1207.1352*, 2012.
- [21] S. Hylleberg. *Seasonality in regression*. Academic Press, 2014.
- [22] M. Jawurek and F. Kerschbaum. Fault-tolerant privacy-preserving statistics. In *PETS*, 2012.
- [23] C. Kopp, M. Mock, and M. May. Privacy-preserving distributed monitoring of visit quantities. In *ACM SIGSPATIAL*, 2012.
- [24] J. Krumm. Inference attacks on location tracks. In *Pervasive*, 2007.
- [25] K. Kursawe, G. Danezis, and M. Kohlweiss. Privacy-friendly aggregation for the smart-grid. In *PETS*, 2011.
- [26] W. H. Lam, K. Chan, M. L. Tam, and J. W. Shi. Short-term travel time forecasts for transport information system in hong kong. *Journal of Advanced Transportation*, 2005.
- [27] A. Machanavajjhala, D. Kifer, J. Abowd, J. Gehrke, and L. Vilhuber. Privacy: Theory meets practice on the map. In *IEEE International Conference on Data Engineering*, 2008.
- [28] L. Melis, G. Danezis, and E. De Cristofaro. Efficient Private Statistics with Succinct Sketches. In *NDSS*, 2016.
- [29] M. A. Munizaga and C. Palma. Estimation of a disaggregate multimodal public transport origin–destination matrix from passive smartcard data from santiago, chile. *Transportation Research Part C: Emerging Technologies*, 2012.
- [30] B. Pan, Y. Zheng, D. Wilkie, and C. Shahabi. Crowd sensing of traffic anomalies based on human mobility and social media. In *ACM SIGSPATIAL*, 2013.
- [31] V. Pandurangan. On Taxis and Rainbows. <https://tech.vijayp.ca/of-taxis-and-rainbows-f6bc289679a1>.
- [32] V. Pejovic and M. Musolesi. Anticipatory mobile computing: A survey of the state of the art and research challenges. *ACM Computing Surveys*, 2015.
- [33] M. Piorkowski, N. Sarafijanovic-Djukic, and M. Grossglauser. CRAWDAD Dataset. <http://crawdad.org/epfl/mobility/20090224>, 2009.
- [34] R. A. Popa, A. J. Blumberg, H. Balakrishnan, and F. H. Li. Privacy and accountability for location-based aggregate statistics. In *ACM CCS*, 2011.
- [35] V. Rastogi and S. Nath. Differentially private aggregation of distributed time-series with transformation and encryption. In *ACM SIGMOD*, 2010.
- [36] G. Sagl, M. Loidl, and E. Beinat. A visual analytics approach for extracting spatio-temporal urban mobility information from mobile network traffic. *ISPRS International Journal of Geo-Information*, 2012.
- [37] E. Shi, T. H. Chan, E. Rieffel, R. Chow, and D. Song. Privacy-preserving aggregation of time-series data. In *NDSS*, 2011.
- [38] R. Silva, S. M. Kang, and E. M. Airoldi. Predicting traffic volumes and estimating the effects of shocks in massive transportation systems. *National Academy of Sciences*, 2015.
- [39] A. Stathopoulos and M. G. Karlaftis. A multivariate state space approach for urban traffic flow modeling and prediction. *Transportation Research Part C: Emerging Technologies*, 2003.
- [40] B. Sun, F. Yu, K. Wu, and V. Leung. Mobility-based anomaly detection in cellular mobile networks. In *ACM Workshop on Wireless Security*, 2004.
- [41] D. Thom, H. Bosch, S. Koch, M. Wörner, and T. Ertl. Spatiotemporal anomaly detection through visual analysis of geolocated twitter messages. In *IEEE Pacific Visualization Symposium*, 2012.
- [42] G. Yavas, D. Katsaros, Ö. Ulusoy, and Y. Manolopoulos. A data mining approach for location prediction in mobile environments. *Data & Knowledge Engineering*, 2005.
- [43] H. Zang and J. Bolot. Anonymization of location data does not work: A large-scale measurement study. In *MobiCom*, 2011.
- [44] Y. Zheng, H. Zhang, and Y. Yu. Detecting collective anomalies from multiple spatio-temporal datasets across different domains. *ACM SIGSPATIAL*, 2015.
- [45] C. Zhong, M. Batty, E. Manley, J. Wang, Z. Wang, F. Chen, and G. Schmitt. Variability in regularity: Mining temporal mobility patterns in london, singapore and beijing using smart-card data. *PloS One*, 2016.