Differentially Private Occupancy Monitoring from WiFi Access Points

Abbas Zaidi, Ritesh Ahuja, Cyrus Shahabi

USC Information Laboratory

University of Southern California

Los Angeles, USA

abbaszai.riteshah.shahabi@usc.edu

Abstract—Accurately monitoring the number of individuals inside a building is vital to limiting COVID-19 transmission. Low adoption of contact tracing apps due to privacy concerns has increased pervasiveness of passive digital tracking alternatives. Large arrays of WiFi access points can conveniently track mobile devices on university and industry campuses. The CrowdMap system employed by the University of Southern California enables such tracking by collecting aggregate statistics from connections to access points around campus. However, since these devices can be used to infer the movement of individuals, there is still a significant risk that even aggregate occupancy statistics will violate the location privacy of individuals. We examine the use of Differential Privacy in reporting statistics from this system as measured using point and range count queries. We propose discretization schemes to model the positions of users given only user connections to WiFi access points. Using this information we are able to release accurate counts of occupants in areas of campus buildings such as labs, hallways, and large discussion halls with minimized risk to individual users' privacy.

Index Terms—differential privacy, WiFi access points

I. INTRODUCTION

Accurately monitoring the number of individuals in indoor spaces is vital to supporting facilities administrators with the challenges of efficiently allocating spaces while ensuring the health & safety of students and staff. For example, in the recent COVID-19 pandemic, contact tracing has been an effective tool in tracking disease spread. Occupancy monitoring in a small business is relatively easy compared to universities, industry campuses, and large businesses whose essential workers must also be afforded protections in the workplace.

Digital contact tracing solutions that utilize GPS and lowenergy Bluetooth technologies in mobile devices to monitor real-time locations or proximity of individuals may be used to monitor occupancy, but raise serious privacy concerns. A location trace can expose users to a wide range of attacks such as unwanted spams/scams or physical danger, and various associated privacy breaches that may disclose sensitive personal details such as one's health status, political, or religious inclinations [1].

At the University of Southern California (USC), we have developed **CrowdMap**, a non-intrusive passive digital tracking platform that utilizes a large network of WiFi routers. Research has shown WiFi to be the most promising alternative to GPS for indoor, context-aware and location-based services [3], [4]. Due to privacy concerns, such a WiFi system collects only

connection logs and their timestamps. This information is generated at thousands of Access Points (APs) across the university's three campuses to estimate occupancy counts. Fig. 2 shows the placement of the APs on a single floor in one of the buildings on campus. The monitoring platform enables users to view aggregate population density counts at each AP over time in different visual and textual formats, such as chloropeth maps, scatter plots, statistic reports, and more.

The granularity of reporting in such traditional WiFi systems is severely limited due to privacy concerns, with only connection events being recorded (as opposed to triangulated users' positions). However, since an individual may make connections to an AP from nearby hallways and classrooms or even through walls, estimating the exact location of individuals is an important challenge. Such reporting is crucial to identifying work and common areas where employees could have close contact (within 6 feet) with others such as the cafeteria, locker rooms, waiting areas, and routes of entry and exit. At the same time, since these devices can be linked to the movement of individuals, there is still a significant danger that aggregate computation will violate the location privacy of individuals. To prevent privacy breaches that may result from uncontrolled, direct release of occupancy statistics, it is important to design and deploy techniques for privacy protection.

In this work we propose a privacy analysis of the CrowdMap system utilizing the powerful Differential Privacy (DP) model. DP [2] is a popular model to achieve privacy, but its mechanisms—via the addition of random noise—can also lead to a severe decrease in data utility. Previous work [7] has shown DP's success in leveraging noise to secure data; our analysis extends DP's use to various queries and data representations.

We first address the use-case of reporting statistics *per AP* which, while of limited use due to reporting granularity, is still useful for learning building-level and floor-level occupancy counts. Since each AP has its own identifier, this data may be viewed as a 1D histogram. We employ state-of-the-art DP mechanisms for histogram publishing for this task. We show that in this setting simple mechanisms are effective even when noise may perturb the true counts significantly, since utility for real-world use is preserved.

Next, we focus on *range queries*, since they can be used as building blocks in many processing tasks. Range count queries

can be executed over the space of a floor of a building, thus a 2D query. We show that utility degradation due to DP noise dwarfs in comparison to errors that implicit discretization of the space incur. We use a baseline grid-based approach to uniformly spread the counts at APs in the spatial region around it. We also propose an advanced approach that utilizes a Voronoi partitioning of the floor space into regions where user counts can be uniformly assumed to be spread. We greatly reduce error in reporting counts and hence make possible answering accurate 2D range count queries on WiFi systems that collect only AP connection logs. Once sanitized, the data from CrowdMap (or similar WiFi systems) can be publicly released while guaranteeing the privacy of each individual. Our contributions are as follows:

- We evaluate the state-of-the-art DP mechanisms in 1-Dimensional and 2-Dimensional subsets of the location data reported in CrowdMap.
- We demonstrate the differences between existing DP mechanisms and their impacts in the chosen setting.
- We present competing data representations of user locations to evaluate the privacy-utility trade-off against a manually defined ground truth.

We organize the rest of the paper as follows: Section II describes the CrowdMap system and the data collection process. Section III introduces background on differential privacy. Section IV analyzes 1-Dimensional point queries. Section V illustrates challenges of answering range queries in 2-Dimensional setting. Sections VI and VII present an experimental evaluation on range count queries based on proposed discretization schemes to improve utility. We conclude in Section VIII.

II. CROWDMAP SYSTEM

CrowdMap [5] is a passive tracking system designed at USC's Integrated Media Systems Center (IMSC) to manage COVID-19 transmission on campus. The physical assets in CrowdMap include thousands of access points connected to a centralized controller and a sophisticated network operations system running on dedicated servers. APs are strategically placed throughout USC's campuses to optimize WiFi signals to users, and as a byproduct, are ideal for passively capturing whereabouts of individuals on campus. CrowdMap captures data packets transmitted in existing WiFi traffic and extracts both received signal strength (RSS) and MAC address of each user's mobile device without installing any dedicated apps. Devices periodically transmit data to maintain their association with the access point.

The raw data is comprised of connection information uploaded nightly, each data point consisting of a unique anonymized identifier (MD5-hashed MAC address), AP name, connect time, disconnect time, connection duration, and SSID. As a pre-processing step, raw data is cleaned to filter out non-user connections, as the network management system does not distinguish between stationary devices such as printers and individuals' active devices. Several data cleaning steps, calibrated using on-site analysis, help CrowdMap accurately

estimate the counts of users, as opposed to counts of devices of a user. For example, by simply deleting all connections shorter than 3 minutes or longer than 12 hours, moving individuals and stationary devices can be reconciled. For the privacy-related work our focus will be on USC's main campus, the University Park Campus (UPC), as it is the largest of the three campuses, with daily averages from March 6, 2021 to May 25, 2021 showing 156 buildings, 3902 APs, 580777 connection attempts before cleaning, and 53504 connections after cleaning.

III. DIFFERENTIALLY PRIVATE MONITORING

Differential privacy (DP) is a privacy-preserving technique which guarantees that results will not be affected whether or not any one individual is present in the data [6]. It allows learning of useful information about a population using aggregate statistics while permitting privacy leakage up to a statistically bound value of ε .

An ε -differentially private algorithm [6] is defined as a randomized algorithm M, for neighboring datasets D and D' differing by at most one element and for all $Y \subseteq Range(M)$ for which the following holds:

$$Pr[M(D) \in Y] \le e^{\varepsilon} Pr[M(D') \in Y]$$

To achieve ε -DP, the result obtained by evaluating a function (e.g., a query) f on the input data must be perturbed by adding noise sampled from a random variable Z. The amount of noise required to ensure the mechanism $\mathcal{M}(D)=f(D)+Z$ satisfies a given privacy guarantee depends on how sensitive the function f is to changes in the input, and the specific distribution chosen for Z. The Laplace Mechanism (LPM) [2] is tuned to the sensitivity S_f computed according to the global ℓ_1 -norm as $S_f = \sup_{D \simeq D'} |f(D) - f(D')|_2$ for every pair of neighboring datasets D, D'. LPM adds zero-mean Laplace noise $Z = Laplace(x|b) = \frac{1}{2b}e^{-\frac{|x|}{b}}$, where $b = S_f/\varepsilon$.

We examine the impacts of differential privacy in 1D space with simple point-queries, posing questions such as 'How many individuals are in a particular building?' and in 2D space with more complex range queries, posing queries such as 'How many individuals are in this hallway (specified as a 2D bounding-box) on a particular floor?'. The use of both dimensions is critical in enabling useful applications such as evaluating utilization rate of spaces.

Further, we examine DP-algorithms which are either data-independent or data-dependent. A data-independent algorithm maintains a consistent error rate over all possible datasets on a fixed domain since the noise scale for such a mechanism is independent to the scale of the data. Whereas, a data-dependent algorithm adjusts its error rate, while consuming the privacy budget to determine the scale of the data [10]. Of the mechanisms analyzed, *Identity*, *Privelet*, *H*, *HB*, and *Greedy H* are data-independent, while all others are data-dependent [7].

Identity [2] adds noise to each data-point using LPM. Privelet [8] applies a wavelet transform according to input data frequencies before adding logarithmic noise. H [7] is a hierarchical method that uniformly allocates the privacy

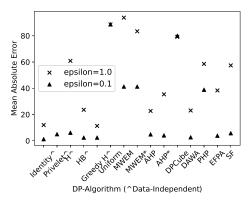


Fig. 1. Mean Absolute Error of DP Algorithms in 1D

budget. HB [9] organizes queries in a tree structure and uniformly distributes the privacy budget over each node, adaptively selecting the optimal branching factor b. Greedy H [10] functions similar to H, forming bins in the generated hierarchical tree in a greedy fashion. Uniform [7] is a data-dependent baseline algorithm that employs a uniformity assumption on the input data. MWEM [11] updates counts with multiplicative weights under an exponential mechanism for a pre-defined number of rounds. AHP [12], or Accurate Histogram Publication, uses a greedy clustering algorithm to form partitions of noisy counts obtained by the Laplace mechanism, setting cells with counts < a pre-defined threshold value to 0. DPCube [13] employs a two-phase partitioning, first using the domain of input data and then using a kd-tree. DAWA [10] uses dynamic programming to compute least cost partitions of noisy counts. QuadTree [17] recursively partitions input data into equal sized quadrants. UGrid [16], or Uniform Grid, partitions input data domain into $m \times m$ equally sized cells and calculates a noisy count for each cell. AGrid [16], or Adaptive Grids, adaptively partitions cells before extracting noisy counts for each. EFPA [14], also known as the Enhanced Fourier Perturbation Algorithm, applies a discrete Fourier transform to input data and adds Laplacian noise to the top kFourier coefficients. SF [15], or StructureFirst, is a histogram based algorithm which determines the structure of input data before applying noise.

IV. 1-DIMENSIONAL RESULTS

Recall that in 1-Dimension we treat each AP as a unique identifier, ignoring its proximity information to other APs. Thus only point queries can be posed against such data, for instance asking questions such as 'How many individuals are in a particular floor?'. Such a query may be answered by aggregating the counts reported at each AP on a floor. Using workloads of such queries, mean absolute error was measured between true counts in the input data and noisy counts produced from application of DP mechanisms with privacy budgets of 0.1 (high privacy regime) and 1.0 (low privacy regime).

Fig. 1 presents the MAE for all state-of-art DP algorithms in 1-Dimension. A Data Independent algorithm such as *Identity* is a basic approach which adds Laplacian noise to the histogram and reports the final counts. Despite the algorithm's simplicity, it performs well on simple point queries. This is as expected since there is limited proximal information that more advanced methods can exploit. For instance, *DPCube* performs a two step partitioning of the histogram, attempting to find regions with similar counts in order to smooth out the random noise. This, while intuitive in instances where nearby regions have similar aggregate counts, results in extreme deviation from input data in our setting and given the small scale of our sample, produces poor results.

We conclude that the signal (the counts at each AP) are large enough to not be diluted by the noise (Laplace mechanism for example has a variance of $2/\varepsilon$). However, 1D point queries have limited use. Therefore we extend our data representation to the 2-Dimensional setting.

V. ESTABLISHING A BASELINE: GROUND TRUTH

The 2-Dimensional setting enables us to estimate how many individuals may be located within a specified range. However, given only the connection logs provided to us from the CrowdMap system, it is difficult to answer accurately how many occupants are within a particular query range. For instance, a query covering the region of a room may return an incorrect answer when the associated AP serving that room is located outside its (and hence the query's) spatial extents.

In order to evaluate how successful a particular DP mechanism is in answering such queries, we must have some existing knowledge of where said user is located. For this setting, we pose range queries on a single floor: the 1st floor of the USC Herman Ostrow Dental School building, shown in Fig. 2. Our ground truth is approximated as bounding boxes, seen in Fig. 2, around each access point using knowledge of WiFi signal strength in relation to physical obstructions (walls in a building) and on-site WiFi receiving devices. These boxes serve to estimate a region in which a user connected to a particular AP is assumed to be uniformly located within. Aggregate counts at each AP, previously used for 1-Dimensional point queries, are uniformly spread in each bounding box, providing each connection with a coordinate location used for evaluating accuracy of query responses.

VI. 2-DIMENSIONAL RESULTS

A. Query Size

In the 2-Dimensional space we are now also presented with another parameter to tune: query size. For clarity of presentation, we normalize our input floorplan to a -5 x 5 grid. A query consists of location coordinates (x,y) within our normalized space and a query size q. Using this information, each query is denoted by the polygon described by (x-q,y-q),(x-q,y+q),(x+q,y+q),(x+q,y-q). For example, a large query with q=2.0 covers a region of more than 10% of our input floor space. Hence we focus on query

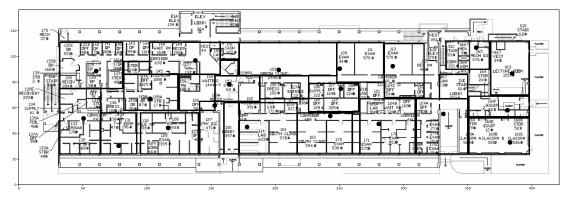


Fig. 2. Ground Truth

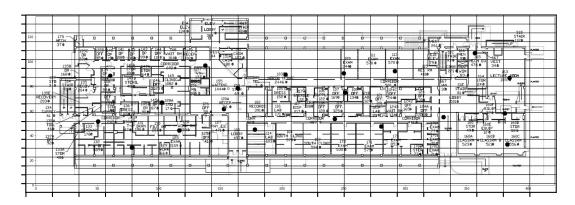


Fig. 3. Gridding Approach

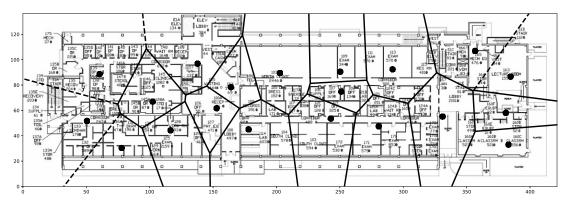


Fig. 4. Voronoi Approach

sizes of q = 0.25, 0.50, 1.0 to simulate realistic queries that a query issuer might ask in practice.

B. Grid discretization approach

Querying in 2D naturally gives way to a grid partitioning of our floorplan. In this approach we split our floorplan into $n \times n$ cells of equal size, shown in Fig. 3, and employ a uniformity assumption on each cell, assuming that counts at a particular AP are uniformly spread over the entire space of any cells that contain said AP. A query overlapping a cell that contains an AP will return the total count at said AP. Further, a query

overlapping 0 AP-containing cells will have a true response of 0. Choice of granularity n is highly specific to the floor plan, the number and position of APs, and the query sizes of interest to a query issuer. This makes automating this parameter very difficult.

We then utilize differentially private mechanisms on work-loads of randomly generated range queries over our grid-partitioned floorplan. The utilization of range queries instead of the previously used point queries is critical in that we are now able to identify coordinates of individual connections to each AP, similar to identifying user's exact locations.

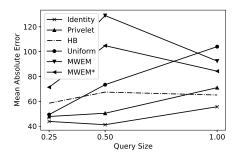


Fig. 5. 2D Gridding Results at $\varepsilon = 0.1$

Results shown in Fig. 5, and for the remainder of this paper, are a subset of the mechanisms evaluated in Fig. 1. Poorly performing algorithms, like Greedy H and DPCube are left out for clarity. Considering previously determined algorithmic performance in related works such as those in [7], we are able to identify the variance between relative inserted noise. We observe that error values are significantly larger than the noise inserted by the privacy mechanisms at the corresponding privacy levels. Therefore, we hypothesize that the grid discretization itself introduces significant errors and hence fails to convey any useful information in the range counts. Also, the employed uniformity assumption fails to accurately account for spread of users surrounding a particular AP, such as connections that may truly reside in a different physical room than that of the AP. Ultimately, the large incurred error is a result of the poor representation mechanism and accompanying uniformity assumption.

C. Voronoi Approach

We propose to split the floorplan into a Voronoi partitioning to improve estimates of user positions. By setting each AP as the site used to construct the surrounding polygons, shown in Fig. 4. In contrast to the Grid-based approach, we uniformly spread the counts at each AP over the entire space of the corresponding Voronoi polygon, assigning each count a coordinate location. This allows us to reasonably infer user locations in relation to whichever AP they are connected to. In this way we can also more accurately represent the spread of WiFi signals generated at each AP subject to physical constraints such as walls, which a baseline Gridding approach fails to achieve.

We evaluate this approach by examining DP mechanisms on workloads of randomly generated range queries over our Voronoi diagram. Results of the best performing DP mechanisms in Fig. 6 show that larger query sizes incur larger error. We hypothesize that this is due to incrementally incurred DP noise of queries covering multiple Voronoi cells. However, the Voronoi Approach performs significantly better than the Gridding Approach. We formalize this notion further by comparing results of both approaches.

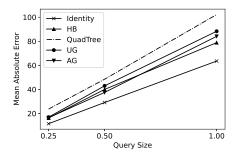


Fig. 6. 2D Voronoi Results at $\varepsilon = 0.1$

VII. 2-DIMENSIONAL METHOD COMPARISON

Fig. 7 and Fig. 8 show significant differences in results produced by each approach at the same DP privacy parameter $\varepsilon=0.1$. Unsurprisingly, as epsilon increases in Fig. 11 and Fig. 12 (meaning less noise is inserted to the counts of each AP), each DP mechanism improves in accuracy. Notably the Voronoi approach produces much less error at each measured value of ε , primarily due to its better ability to spread counts in the vicinity of each AP. Hence we examine the differences between the approaches used in spreading aggregate counts at each AP across their respective floorplans.

As stated previously, the two approaches have fundamental differences in their discretization of the floorplan. We focus our attention on queries of smaller sizes (such as in Figure 7) as they align more closely with queries we may pose in reallife scenarios. Assuming that one AP covers one particular room on a given floor, identifying locations of users connected to that AP would constitute a small query size. In this setting, our Voronoi approach clearly outperforms the baseline grid partitioning by reducing the implicit errors in discretization. We believe this to be a result of the Voronoi partitioning's ability to much more accurately mimic WiFi signal strength. Slightly larger query sizes of Q = 0.50 shown in Fig. 8 show similar results. But for very large query sizes (analogous to large auditoriums) we see that both approaches produce similar results, as in Fig. 9. This is due to larger query sizes covering multiple Voronoi cells and Grid cells, respectively. Hence, noise error is incrementally accumulated in each case.

In Fig. 10 we evaluate two of the best performing algorithms in both approaches: *Identity* and *AG*. Results displayed show that the DP mechanisms in the Voronoi approach function much better than in the Gridding approach especially at smaller query sizes. As previously noted, lower query sizes more closely align to real-world applications of such a system. Hence, the proposed Voronoi partitioning is much more suited to this application as it is capable of mimicking WiFi spread in a highly accurate manner around access points.

VIII. CONCLUSION

In a 1-Dimensional setting, we see that simple DP mechanisms are sufficiently accurate in protecting user privacy, provided that the use-cases only require simple point queries.

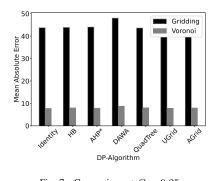


Fig. 7. Comparison at Q = 0.25

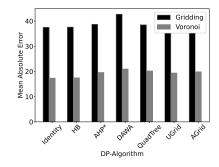


Fig. 8. Comparison at Q = 0.50

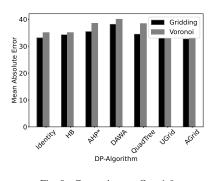


Fig. 9. Comparison at Q = 1.0

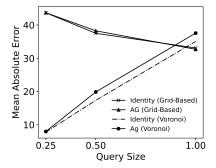


Fig. 10. Identity and AG By Approach at $\varepsilon=1.0$

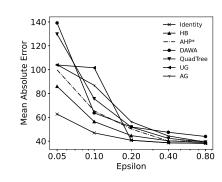


Fig. 11. Avg MAE Over All Query Sizes in Gridding Approach

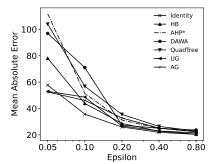


Fig. 12. Avg MAE Over All Query Sizes in Voronoi Approach

However, when a query issuer may want to know counts in a 2-Dimensional spatial region, whether that be an entire room, hallway, or floor, we show that discretization schemes are more important to accurately answering such queries than the noise introducing DP mechanisms.

IX. ACKNOWLEDGEMENT

This research has been funded in part by NIH award R01LM014026 and NSF grants IIS-1910950, CNS-2027794, CNS-2125530 and IIS-2128661, and an unrestricted cash gift from Microsoft Research. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the sponsors such as the NIH and NSF.

REFERENCES

- P. Kalnis, G. Ghinita, K. Mouratidis and D. Papadias, "Preventing Location-Based Identity Inference in Anonymous Spatial Queries," in IEEE Transactions on Knowledge and Data Engineering, vol. 19, no. 12, pp. 1719-1733, Dec. 2007, doi: 10.1109/TKDE.2007.190662.
- [2] Dwork, C., McSherry, F., Nissim, K., Smith, A. (2006). Calibrating Noise to Sensitivity in Private Data Analysis. In: Halevi, S., Rabin, T. (eds) Theory of Cryptography. TCC 2006. Lecture Notes in Computer Science, vol 3876. Springer, Berlin, Heidelberg.
- [3] H. Zou, Y. Zhou, J. Yang, W. Gu, L. Xie and C. Spanos, "FreeDetector: Device-Free Occupancy Detection with Commodity WiFi," 2017 IEEE International Conference on Sensing, Communication and Networking (SECON Workshops), 2017, pp. 1-5, doi: 10.1109/SEC-ONW.2017.8011040.

- [4] H. Zou, Y. Zhou, J. Yang, C. J. Spanos, "Device-free occupancy detection and crowd counting in smart buildings with WiFienabled IoT," Energy and Buildings, 2018, pp. 309-322, doi: 10.1016/j.enbuild.2018.06.040
- [5] S. Min, R. Ahuja, Y. Liu, A. Zaidi, C. Phu, L. Nocera, and C. Shahabi, "CrowdMap: Spatiotemporal Visualization of Anonymous Occupancy Data for Pandemic Response," ACM SIGSPATIAL, 2021.
- [6] C. Dwork and A. Roth, "The algorithmic foundations of Differential Privacy," Foundations and Trends® in Theoretical Computer Science, vol. 9, no. 3-4, pp. 211–407, 2013.
- [7] M. Hay, A. Machanavajjhala, G. Miklau, Y. Chen, and D. Zhang, "Principled evaluation of differentially private algorithms using DPBench," ACM SIGMOD, 2016.
- [8] X. Xiao, G. Wang, and J. Gehrke, "Differential privacy via wavelet transforms," ICDE 2010.
- [9] W. Qardaji, W. Yang, and N. Li, "Understanding hierarchical methods for differentially private histograms," PVLDB, 2013.
- [10] C. Li, M. Hay, G. Miklau, and Y. Wang, "A data- and workload-aware algorithm for range queries under Differential Privacy," PVLDB 2014.
- [11] M. Hardt, K. Ligett, F. McSherry, "A simple and practical algorithm for differentially private data release," NeurIPS 2012
- [12] X. Zhang, R. Chen, J. Xu, X. Meng, and Y. Xie, "Towards accurate histogram publication under Differential Privacy," Proceedings of the 2014 SIAM International Conference on Data Mining, 2014.
- [13] Y. Xiao, L. Xiong, L. Fan, S. Goryczka, "DPCube: Differentially Private Histogram Release through Multidimensional Partitioning," 2012
- [14] G. Acs, C. Castelluccia, and R. Chen, "Differentially private histogram publishing through lossy compression," 2012 IEEE 12th International Conference on Data Mining, 2012.
- [15] J. Xu, Z. Zhang, X. Xiao, Y. Yang, and G. Yu, "Differentially private histogram publication," 2012 IEEE 28th International Conference on Data Engineering, 2012.
- [16] W. Qardaji, Weining Yang, and Ninghui Li, "Differentially private grids for geospatial data," 2013 IEEE 29th International Conference on Data Engineering (ICDE), 2013.
- [17] G. Cormode, C. Procopiuc, D. Srivastava, E. Shen, and T. Yu, "Differentially private spatial decompositions," IEEE ICDE, 2012.