

Sharing Personal ECG Time-series Data Privately

Luca Bonomi¹, Zeyun Wu², Liyue Fan³

¹Department of Biomedical Informatics, Vanderbilt University Medical Center, Nashville, TN, 37232, USA

²Department of Electrical and Computer Engineering, University of California San Diego, La Jolla, CA, 92093, USA

³Department of Computer Science, University of North Carolina at Charlotte, Charlotte, NC, 28223, USA

Word count excluding title page, abstract, references, figures and tables: 4340

ABSTRACT

Objective Emerging technologies (e.g., wearable devices) have made it possible to collect data directly from individuals (e.g., time-series), providing new insights on the health and well-being of individual patients. Broadening the access to these data would facilitate the integration with existing data sources (e.g., clinical and genomic data) and advance medical research. Compared to traditional health data, these data are collected directly from individuals, are highly unique and provide fine-grained information, posing new privacy challenges. In this work, we study the applicability of a novel privacy model to enable individual-level time-series data sharing while maintain the usability for data analytics.

Methods and materials We propose a privacy-protecting method for sharing individual-level electrocardiography (ECG) time-series data, which leverages dimensional reduction technique and random sampling to achieve provable privacy protection. We show that our solution provides strong privacy protection against an informed adversarial model while enabling useful aggregate-level analysis.

Results We conduct our evaluations on two real-world ECG datasets. Our empirical results show that the privacy risk is significantly reduced after sanitization while the data usability is retained for a variety of clinical tasks (e.g., predictive modeling and clustering).

Discussion Our study investigates the privacy risk in sharing individual-level ECG time-series data. We demonstrate that individual-level data can be highly unique, requiring new privacy solutions to protect data contributors.

Conclusion The results suggest our proposed privacy-protection method provides strong privacy protections while preserving the usefulness of the data.

Keywords: Data Privacy, ECG Data, Time-series, Data Sharing, Predictive Analytics

INTRODUCTION

Advances in mobile sensor technology (e.g., wearable devices, smartphones) are revolutionizing the way in which patient data are collected, enabling high-quality and fine-grained data to be gathered directly from the individual's personal device.¹ Aggregating these data have been shown to provide great opportunities for performing new analytics and advancing healthcare. For example, the use of wearable sensors health data has been shown to be effective in facilitating the diagnosis, prevention, management of chronic diseases, and improving patient care.²⁻⁴ Additionally, recent research studies have shown that data collected from wearable devices could be used in the early detection of asymptomatic and pre-symptomatic cases of COVID-19.^{5,6} Therefore, sensor data are increasingly integrated into large datasets (e.g., NIH All of Us),⁷ with the promise of providing the medical research community with new insights on the health conditions and well-being of individuals. However, sharing these sensor data broadly poses novel privacy challenges. First, the sole removal of personal identifiable information (e.g., SSN) does not provide

adequate privacy protection. Research studies have shown that individual-level data collected from sensor devices could be used as biometric to identify individuals.^{8,9} Specifically, Biel et al.⁹ have demonstrated that carefully selected features from electrocardiography (ECG) data could be used to achieve 100% re-identification rate on a dataset with 20 patients. Second, current data privacy solutions require a trusted site to collect and manage the data. However, individual patients may lack trust in either the data aggregator and data users (e.g., internal and external researchers).¹⁰ To promote data participation and sharing, it is imperative to develop privacy methods that consider the privacy needs of individual data contributors.

Current solutions for individual-level sensor data (e.g., time-series) build on security and data anonymity primitives. Security-based solutions rely on encryption solutions and access control techniques^{11–13}. Despite promising results in some settings,^{14–16} their end-to-end paradigm allows only a small number of authorized health professionals to access the data. Additionally, encryption solutions may still be vulnerable in the presence of an informed adversary (e.g., in genomic applications¹⁷). To broaden data access while protecting data privacy, solutions based on data anonymization have been proposed (e.g., k -anonymity,¹⁸ differential privacy¹⁹). These techniques enable a data curator to sanitize the collected data and share the results with external users (e.g., researchers).^{20–23} However, these privacy solutions build on a traditional central authority assumption, in which a trusted data curator is responsible for collecting, aggregating, and protecting the individual raw health data. As a result, these traditional data anonymization approaches have limited applicability when data are collected directly from individuals who may not trust the data curator (e.g., DTC genomic data²⁴).

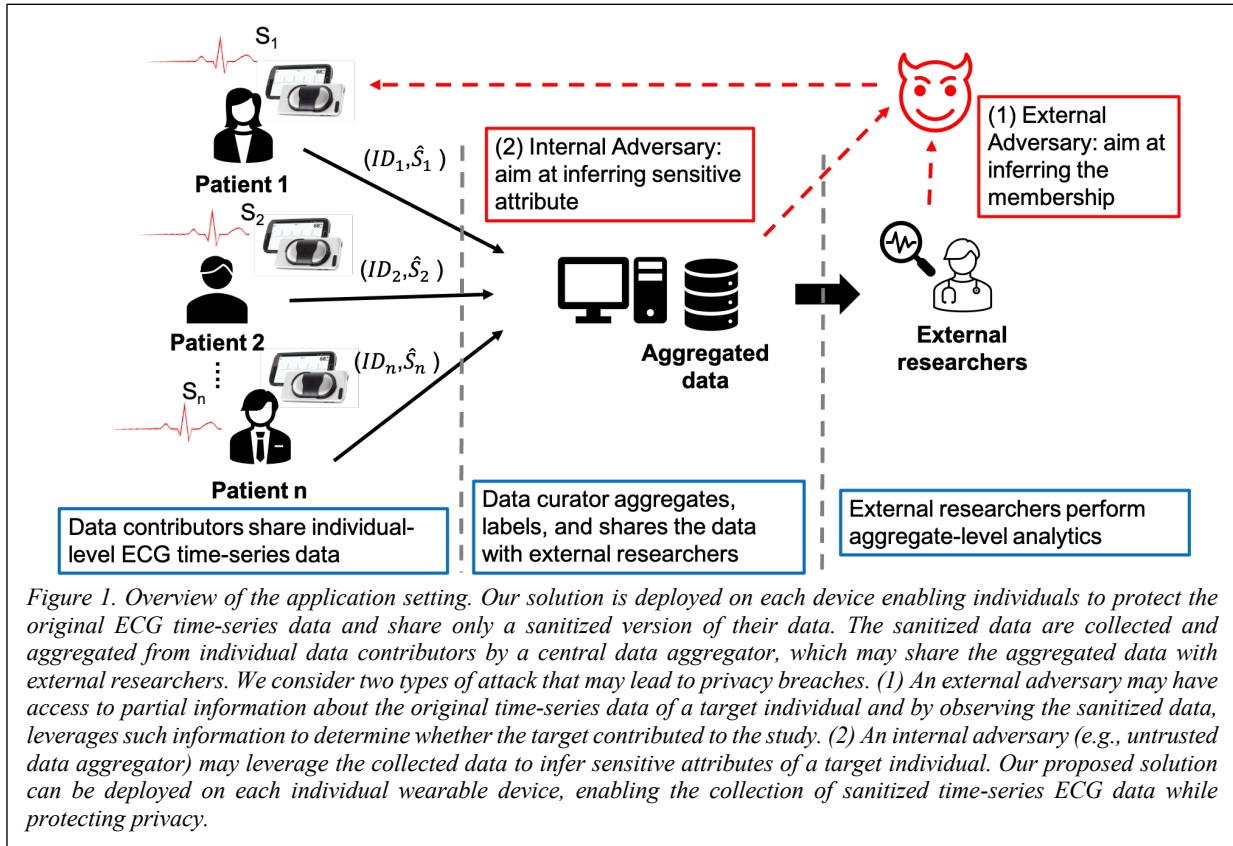
The objective of this work is to assess the privacy risks in sharing ECG time-series data, which are increasingly collected from individual sensor devices, and to develop novel privacy solutions to enable individuals to share their data directly with an untrusted data curator. We show that the ECG time-series data considered in this study are highly unique, enabling an informed adversary to perform accurate inference attacks. As an example, an individual who participates to a research study may inadvertently disclose partial information about her ECG time-series data (e.g., she may share her ECG measurements during a certain day with other mobile applications). Then, an adversary may leverage the disclosed information together with the aggregated data shared by the study to infer the presence of the target individual (e.g., membership inference). To mitigate these privacy risks, we develop a privacy-protecting method that achieves strong privacy protection, enabling individuals to sanitize their data as they are collected for research studies. Our proposed privacy-protecting method builds on the metric privacy model,^{25,26} a generalization of the popular differential privacy notion, which provides provable privacy protection for individual-level data. Our solution can be deployed directly on the sensor devices (e.g., wearables and smartphones) of data contributors, protecting the ECG time-series data as they are collected. We show that our solution is effective in providing privacy protection against an informed adversary, significantly reducing known privacy risks (e.g., membership and attribute inference). Because our solution provides privacy at individual-level, the data aggregation process does not need to rely on a trusted data curator. As a result, individuals have greater control over the shared data compared to traditional central privacy models, which may encourage data participation. Finally, we show that the sanitized aggregate-level data can enable accurate health analytics (e.g., predictive and clustering tasks), demonstrating that the sanitized data preserve the clinical usefulness of the original data.

MATERIALS AND METHODS

Application Setting

In this work, we consider a realistic application setting: ECG data are generated by n individuals and are collected by an untrusted data aggregator, which aggregates these data and shares them with researchers (Figure 1). On the client side, after generating the data each individual applies our privacy method on the original ECG time-series data S to generate a sanitized time-series \hat{S} . These sanitized data are shared with the data aggregator, where a pair (ID, \hat{S}) , represents a pseudo id and the sanitized ECG data contributed by an individual. On the data aggregator side, the collected data may be labeled (e.g., indicating the

presence of certain conditions of interest) and once aggregated they are shared with external researchers to enable clinical predictive tasks (e.g., cardiovascular disease classification).



The application setting is depicted in Figure 1 and presents two major privacy challenges. (1) An external adversary who may have access to known information about a target may leverage the shared data to conduct membership inference attacks with the goal of determining whether the target contributed to the data. (2) The central data aggregator may be untrusted or be compromised, allowing an internal adversary who knows the pseudo-IDs of data contributors to infer some sensitive time-series values (e.g., attribute inference). Mitigating these privacy risks is a challenging problem, which cannot be solved by traditional privacy solutions (e.g., standard differential privacy). Furthermore, there is a tension between protecting individual privacy and the usability of the data: privacy solutions need to preserve usability of the aggregated time-series data to support research studies. Therefore, it is imperative to provide individuals with strong privacy control over the shared individual-level data while preserving the usability at aggregate level.

In the following section, we describe the proposed privacy-protecting solution that sanitizes the individual-level data as they are collected. Specifically, our method allows individuals to share protected time-series ECG data \hat{S} , while the original time-series data S never leave the individual's device. We will show that our solution can provide protection both at individual- and aggregate-level, mitigating attribute and membership inference attacks in the presence of an informed adversary. Additionally, our evaluations show that the data usability at aggregate-level is preserved, enabling researchers to perform accurate aggregate-level analytics (e.g., clustering, and predictive analysis).

Sanitizing Individual Time-series Data

Here, we propose a data sanitization method to provide rigorous privacy for individual-level ECG time-series data. Individual users only share sanitized ECG data with the data aggregator; therefore, the raw time-series data are protected (i.e., never leave individual's wearable devices). Our sanitization method builds on a recent privacy model named *metric privacy*,²⁶ which is a generalization of the differential privacy

model. The metric privacy has shown to provide provable privacy protection for individual-level data.^{27,28} Relevant to our setting, the metric privacy model extends the notion of indistinguishability over arbitrary metric spaces, enabling the design of customizable privacy mechanisms that can be deployed for individual-level data. Specifically, given an arbitrary set of secrets X with a metric d_X , the metric privacy model is defined as follows.

Metric Privacy

A mechanism $M : X \rightarrow P(Z)$ satisfies d_X -privacy, if and only if $\forall x, x' \in X$ the following inequality holds: $K(x)(Z) \leq e^{d_X(x, x')} K(x')(Z) \forall Z \in \mathcal{F}_Z$, where Z is a set of query outcomes, and $P(Z)$ is the set of probability measures over Z .

In practice, the metric d_X can be obtained from a standard metric scaled by a factor ε (privacy parameter/budget). As shown by Chatzikokolakis et al.²⁶, the standard differential privacy protection can be achieved via metric privacy by using the Hamming distance between datasets as metric $d_X = \varepsilon \times d_H$.

In our proposed mechanism, we will use a metric d_X to achieve indistinguishability between ECG time-series data. Intuitively, metric privacy ensures stronger indistinguishability between similar secrets. In our setting, the indistinguishability is stronger between similar ECG time-series (e.g., running vs cycling heartbeats) and it is weakened for different ECG time-series (e.g., running vs sleeping heartbeats). This relaxed privacy notion enables our method to retain the usefulness of the time-series data. In fact, coarser information is preserved to enable usable analytics, while fine-grained data are made indistinguishable, thus protecting privacy. The greater usability provided by this privacy model has been demonstrated by several research works in applications that rely on individual time-series and location data.^{27,29}

Metric privacy for ECG time series. In our setting, we adopt the metric privacy model to protect the ECG time-series generated by each patient. Given any pair of time-series S_1 and S_2 , our sanitization method bounds the ability of an adversary who observes the sanitized time-series \hat{S} , to determine whether such time-series was originally S_1 or S_2 , thus protecting each input ECG time-series. Specifically, our approach comprises three main steps. First, we perform the sanitization in the Discrete Cosine Transform (DCT) domain, where the original time-series are mapped into vectors representing the first l coefficients of their DCT transform. The DCT coefficients can well capture the high-level structure in time series data, thus providing an accurate estimation of their similarity. Additionally, the DCT domain provides us with a compact representation of the original time-series data, which can benefit both scalability and privacy. Second, in the DTC domain, we consider a metric d_X between the DCT representation of S_1 and S_2 , which is defined as $d_X(S_1, S_2) = \varepsilon \times d_2(DCT_l(S_1), DCT_l(S_2))$, where ε is the privacy parameter and $d_2(DCT_l(S_1), DCT_l(S_2))$ is the Euclidean distance between the first l coefficients in the DCT domain. To achieve privacy in this l -dimensional domain, we use the privacy mechanism proposed in our previous studies to sample a plausible set of DCT coefficients that satisfy d_X -metric privacy.^{29,30} Finally, we reverse the DCT transformation based on the sampled coefficients and generate a sanitized ECG time-series \hat{S} . A detailed description of our sampling process is reported in the Appendix.

In real application settings, an individual may generate a continuous stream of ECG time-series data. For example, a wearable device may record a time-series for each heartbeat, generating a collection of multiple time-series segments associated with the same individual (i.e., a profile for the patient). Because the total number of segments in each patient's profile may be unknown a priori, we protect each individual time-series segment independently in our sanitization method. While a higher privacy cost may be accumulated by this approach, it offers consistent and strong privacy protection to each time-series generated by the individual. This approach is also commonly used in real-world applications (e.g., Apple's privacy safeguard^{31,32}), where the privacy budget is refreshed over after a fixed amount of time (e.g., 1 day). In principle, to provide a bounded overall privacy guarantee, we can divide the overall privacy budget by the total number of time-series generated by the individuals and apply our sanitization method to each time-series.^{26,27}

Privacy Measures

Similar to differential privacy, the metric privacy model provides provable privacy protection against a strong adversary with the privacy parameter ϵ . However, it is important to understand how this theoretical privacy protection is effective in mitigating practical privacy risks. A variety of privacy risk measures, including membership and attribute inference risks, have been proposed for well-studied health data types (e.g., genomics, EHR, tabular data).^{24,33–38} In this work, we adapt those privacy measures to suit ECG time-series data based on realistic adversarial models. Below, we briefly describe our privacy measures, a detailed description is reported in the Appendix.

Data uniqueness. Data uniqueness and privacy protection are closely related. As an example, several privacy methods rely on data manipulation techniques (e.g., generalization) to reduce data uniqueness and achieve privacy (e.g., k -anonymity¹⁸). In our work, we aim at measuring the uniqueness of the recorded ECG time-series both in the original and sanitized data. Our measure of data uniqueness is inspired by the privacy measure for mobility data proposed by De Montjoye et al.³⁹ In our setting, the uniqueness score $u_k(D)$ of the shared dataset D represents the fraction of individuals that can be uniquely identified by a subset of k time-series readings. Higher values of uniqueness score indicate that the time-series are more unique.

Membership Inference. For measuring membership disclosure, we consider an informed adversary who has some prior knowledge about the original time-series of a target and aims at determining whether the target participated in the sanitized data. Practically, an individual may inadvertently disclose partial information collected by her mobile sensor devices, for example, sharing sensor data with other applications or on social media. An adversary may leverage the known information of target together with the aggregate-level data shared by the data aggregator, to determine whether the target contributed to the study. In our attack model, the adversary uses the Dynamic Time Warping Distance⁴⁰ (DTW) to match the known time-series information of the target with the sanitized time-series. Then, by identifying possible matches in the sanitized data (e.g., using a threshold for DTW), the adversary decides whether the target time-series were included in the data. We measure the success of such attack in terms of accuracy, where higher values of accuracy indicate higher success in learning the membership of the target in the data.

Attribute Inference. In attribute inference, we consider an informed adversary who knows a target individual participating in the study and partial information about the target’s time-series data. In this setting, the adversary leverages the sanitized data to infer the value of the remaining time-series of the target. As an example, by inspecting the sanitized data the adversary may reconstruct a time-series in the target profile representing a sensitive heartbeat (e.g., myocardial infarction), thus learning a sensitive condition of the target. In our attack model, the adversary uses a K-NN framework to impute the unknown time-series values. In this framework, the adversary first identifies the top-K similar profiles in the sanitized data according to the known time-series, and then impute the unknown time-series values (e.g., average among the ECG readings). The difference between the imputed and the original values (i.e., inference error) can be used to quantify the adversary’s ability to successfully infer the unknown ECG readings of the target, where lower error values indicate higher accuracy in reconstructing the original sensitive values.

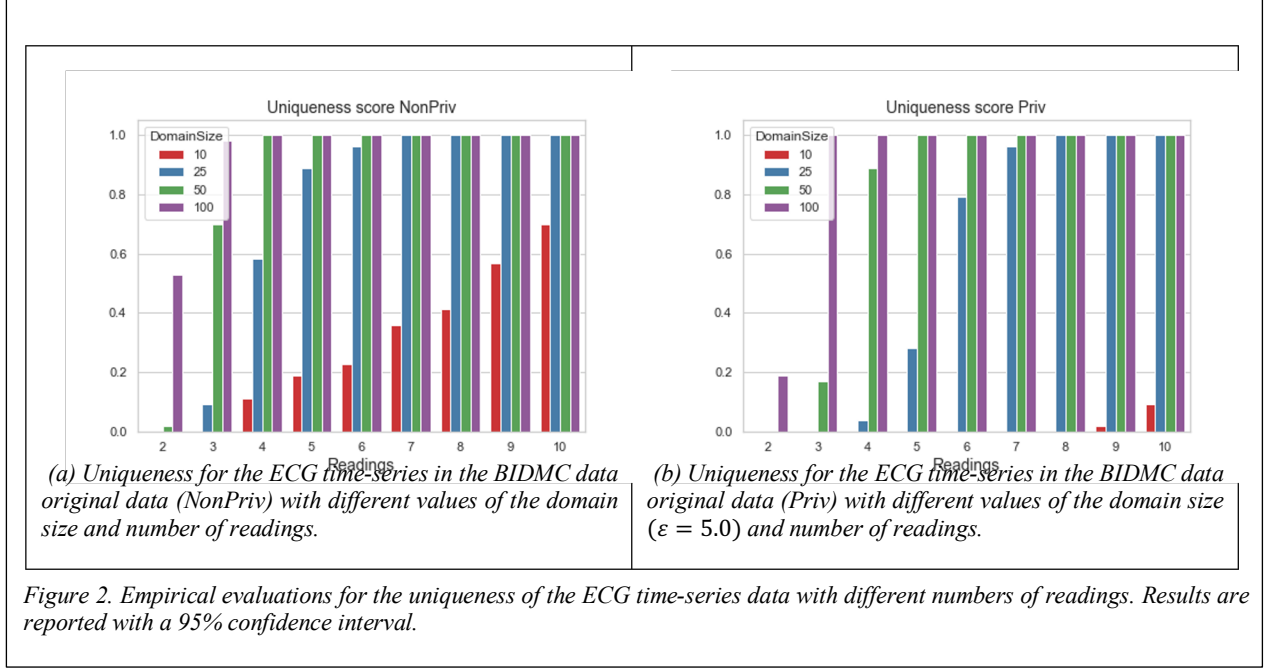
RESULTS

In this section, we describe our evaluations to assess the usability and privacy protection of the sanitized data.

Data Description. In our evaluations, we use two publicly available real-world datasets: ECG 200 and BIDMC. ECG 200 dataset is sampled from the MIT-BIH Supraventricular Arrhythmia Database.⁴¹ The dataset comprises 200 ECG time-series recordings of a single patient, each representing the electrical activity recorded during one heartbeat with 96 values. The time-series are labeled into two classes, representing normal heartbeat and a Myocardial Infarction, respectively.⁴² BIDMC comprises time-series ECG recordings, each of 8-minute duration sampled at 125 Hz, from 53 ICU patients at the Beth Israel

Deaconess Medical Center.^{43,44} In this dataset, we divide the ECG time-series into segments of 0.64 seconds and consider the first 25 segments for each individual's profile.

Algorithm Parameters. We set the number of DCT coefficients to $m = 24$ and vary the privacy parameter ε in the range $[1, 20]$. These parameters were selected according to the guidelines suggested in our previous privacy study for time-series data.²⁹



Privacy Evaluations

Data uniqueness. Figure 2 reports the values of uniqueness for the ECG time-series in the BIDMC dataset. To compute the uniqueness, we map the original real-valued time-series data into a discrete domain. Once discretized, we quantify the uniqueness of the time-series by considering subsequence of consecutive readings that can uniquely identify a time-series in the data. Figure 2(a) shows that for a domain size of 25 symbols, subsequences of 5 readings can uniquely identify more than 89% of the individuals in the BIDMC dataset. We notice that the uniqueness increases as more readings are used. Furthermore, as we decrease the domain size (i.e., time-series are mapped into fewer symbols), the uniqueness decreases. Overall, the time-series in the original BIDMC data are highly unique. Figure 2(b) reports similar results on the time-series sanitized with our privacy method. While our privacy method reduces the data uniqueness compared to the original data, we observe that the time-series are still unique for large domain sizes and higher numbers of readings. Nevertheless, we will show that the sanitized data are protected from membership and attribute inference attacks. In fact, high data uniqueness does not necessarily imply low privacy protection. Specifically, our randomized approach mitigates existing inference attacks by reducing the ability of the adversary to link the known data of the target's with the output sanitized data.

Membership inference risk evaluations. Figure 3 reports the empirical membership inference risk with different values of the privacy parameter. In the membership inference attack, the adversary relies on a threshold value (th) to determine whether the known ECG profile of the target is present in the sanitized data (see detailed methodology in Appendix). To quantify the success of membership inference, we measure the accuracy for different values of the threshold parameter. In the case of the original data (i.e., NonPriv) the accuracy of the adversary decreases as larger values of th are used (i.e., more false positives). Figure 3(a) shows that the accuracy is high even when the adversary knows only one segment, and the values increase with the amount of information available to the attacker (i.e., number of time-series

segments in target's profile). For example, in our evaluations the highest value of accuracy is achieved when 25 time-series are recorded per individual. In Figure 3(b), we assess whether discretizing the time-series data into discrete, coarse representations may provide some privacy protection. From these results, we observe that even with a coarser representation (e.g., domain size 10) the adversary is successful in performing the membership inference attack. In Figure 3(c), we observe that our proposed sanitization method provides robust privacy protection with respect to an adversary with increasing background knowledge about the target profile. In fact, the overall accuracy of the membership inference attack is

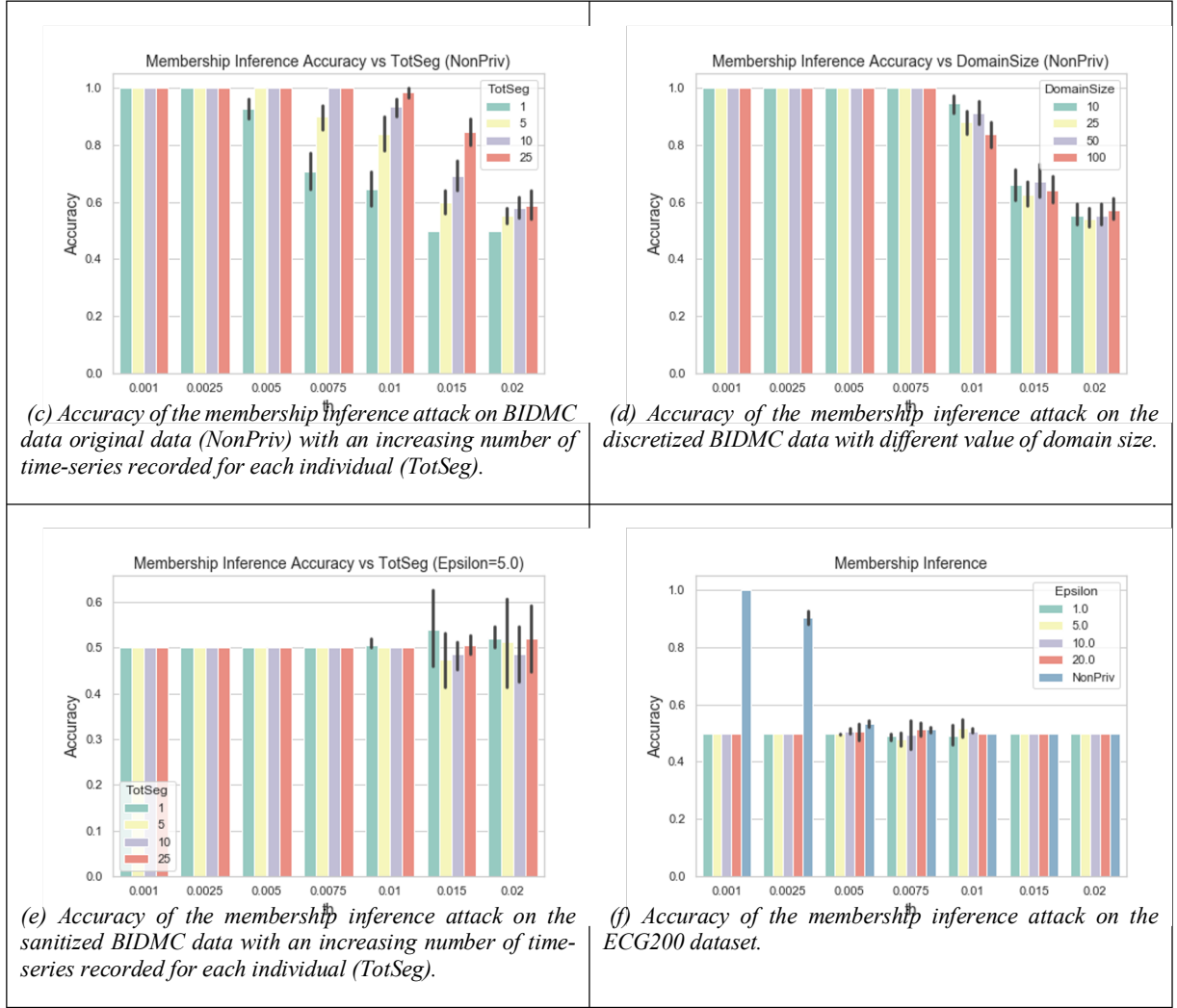
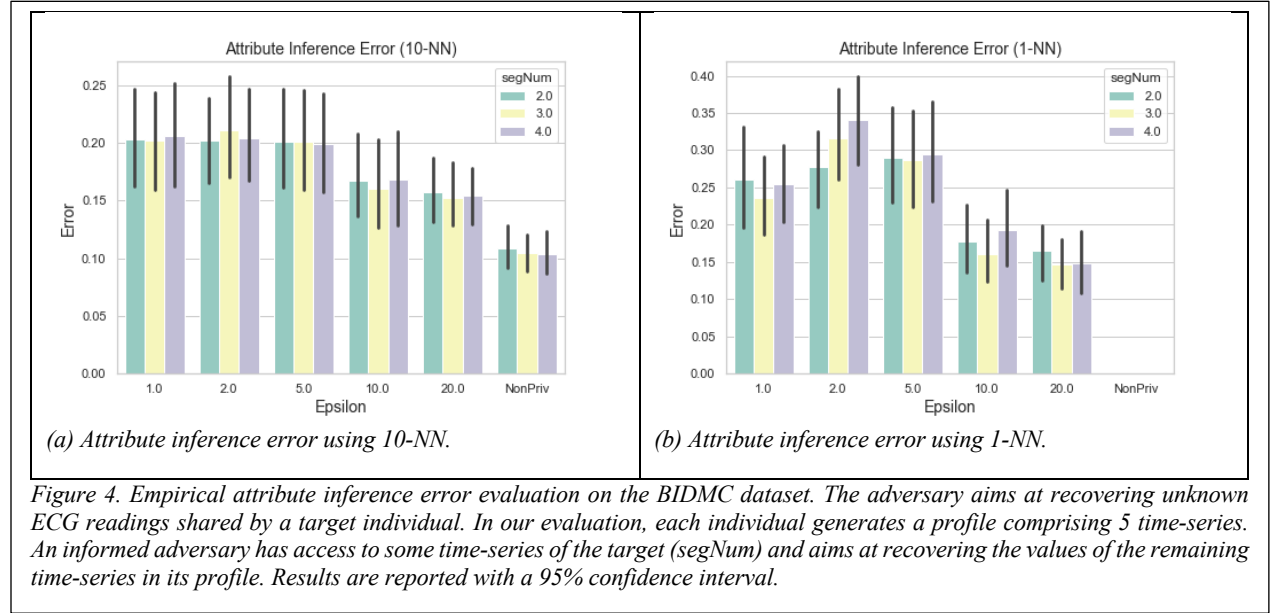


Figure 3. Empirical membership inference risk evaluations. The success of the adversary in inferring the membership of a target individual in the sanitized data is measured in terms of accuracy. We report the privacy measure both on the original data (NonPriv), and the sanitized data with different values of the privacy parameter epsilon. Results are reported with a 95% confidence interval.

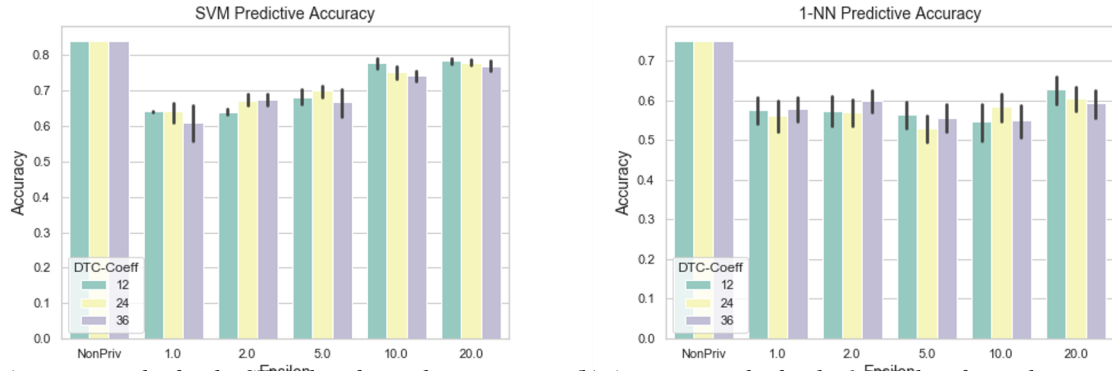
always below 60%. Additionally, our privacy-protecting method is robust against changes in the privacy parameter. For example, the accuracy of the attack is always less than 70%, with the maximum accuracy achieved by the attacker with the largest value of the privacy parameter (see Figure 7 in Appendix). Our proposed sanitization method significantly reduces the accuracy of the membership inference attack also on the ECG200 dataset. For example, in Figure 3(d), the adversary can exactly infer the membership $th \leq 0.001$ on the original data, while after sanitization, the adversary can only achieve 50% accuracy (i.e., random guess) for all the values of th .

Attribute inference risk evaluations. Figure 4 reports the average attribute inference error for the BIDMC dataset. In our attack model, the adversary uses the K-NN based method to impute the unknown time-series shared by the target individual in the sanitized data. In this setting, we fix the number of time-series segments per profile to 5, and we assess the impact of the adversary’s background knowledge on the attribute inference error by increasing the background knowledge of the adversary (i.e., varying the number of time-series segments known by the adversary from 2 to 4). Additionally, we vary the number of neighboring profiles used in the imputation process. We observe that using the original data, the average inference error (i.e., error in imputing the unknown readings) is roughly 10% for $K=10$, while for $K=1$ the adversary can correctly recover the unknown time-series of the target. Our proposed privacy method increases the inference error inflicted by the adversary for both $K=10$ and $K=1$. As an example, for $\epsilon = 5.0$, the average attribute inference error is roughly 20% and 28% for $K=10$ and $K=1$, respectively. The average error gently decreases as the privacy protection is reduced. Overall, we observe that our proposed sanitization method significantly mitigates attribute inference attack on the sanitized data.



Usability of the shared ECG data

Classification results. To evaluate the usability of the sanitized data, we consider a ECG classification task using two classifiers: Support Vector Machine (SVM) and the 1-nearest neighbors (1-NN). For the NN method, we use the DTW to measure the distance between time-series data. These classifiers are simple machine learning models that are widely used in healthcare application settings, including time-series classification.⁴⁵ In these evaluations, we use the ECG200 dataset, in which the recorded time-series data are labeled in two classes: normal heartbeat and Myocardial Infarction. The overall data are divided into 80% training and 20% testing. Our goal is to assess the impact of the privacy protection on the predictive accuracy of these classifiers. To evaluate the impact of privacy, we train the classifiers on the sanitized training set with increasing values of the privacy parameter. Then, we test the classifiers on the non-private test set. **Error! Reference source not found.** reports the predictive accuracy results on the test set with the classifiers trained on the original vs. sanitized data. Among the two classifiers, SVM outperforms 1-NN with DTW. From **Error! Reference source not found.** (a), we observe that the predictive accuracy for the SVM classifier approaches the results on the non-private data as epsilon increases. As an example, with $\epsilon \geq 10.0$, the average predictive accuracy is above 75%. Regarding the results for the 1-NN classifier, the predictive accuracy results in **Error! Reference source not found.** (b) are robust against changes in the privacy parameter, achieving predictive accuracy results above 60% for all settings.



(a) Accuracy results for the SVM classifier with increasing values of privacy parameter (ϵ) and different number of DTC coefficients. (b) Accuracy results for the 1-NN classifier with increasing values of privacy parameter (ϵ) and different number of DTC coefficients.

Figure 5. Utility evaluation in ECG classification. We report the accuracy for the classifiers on the non-private data and the accuracy results on the sanitized data with increasing values of the privacy parameter. Results are reported with a 95% confidence interval.

Clustering results. On the BIDMC dataset, we use an unsupervised clustering approach to identify groups of individuals with similar ECG time-series profiles and assess whether the privacy method proposed in this work may affect the clustering results. Specifically, we consider a hierarchical clustering approach where the DTW distance is used to quantify the similarity between ECG profiles. In principle any unsupervised clustering approach could be used, and hierarchical clustering is a preferred method to eliminate randomness (e.g., in comparison to k-means). Using the hierarchical clustering approach, we identify 2 distinct clusters among individuals in the original data. Information about the age and gender distribution in the identified clusters is summarized in Table 1. Then, we perform clustering on the sanitized data produced with different values of the privacy parameter. To visualize these high dimensional clusters, we use t-SNE technique,⁴⁶ in which the results are projected into a two-dimensional space (**Error! Reference source not found.** (a)-(b)). Additionally, we report all-pairs distance between ECG profiles in the original data and sanitized data (**Error! Reference source not found.** (c)-(d)). As our privacy method performs random sampling in the DTC, it may incur utility loss for clustering sanitized profiles. We observe that for large epsilon values (e.g., $\epsilon = 20$), the pair-wise distance between ECG time-series is well preserved, enabling the hierarchical clustering method to generate clusters that well resemble those identified in the original data.

Table 1. Information about gender and age for the patients in hierarchical clustering with the DTW distance between their ECG time-series profiles.

	Label 0 (n=12)	Label 1 (n=41)
Female	25%	56.1%
Male	75%	43.9%
Age < 65	50%	43.9%
Age \geq 65	50%	56.1%

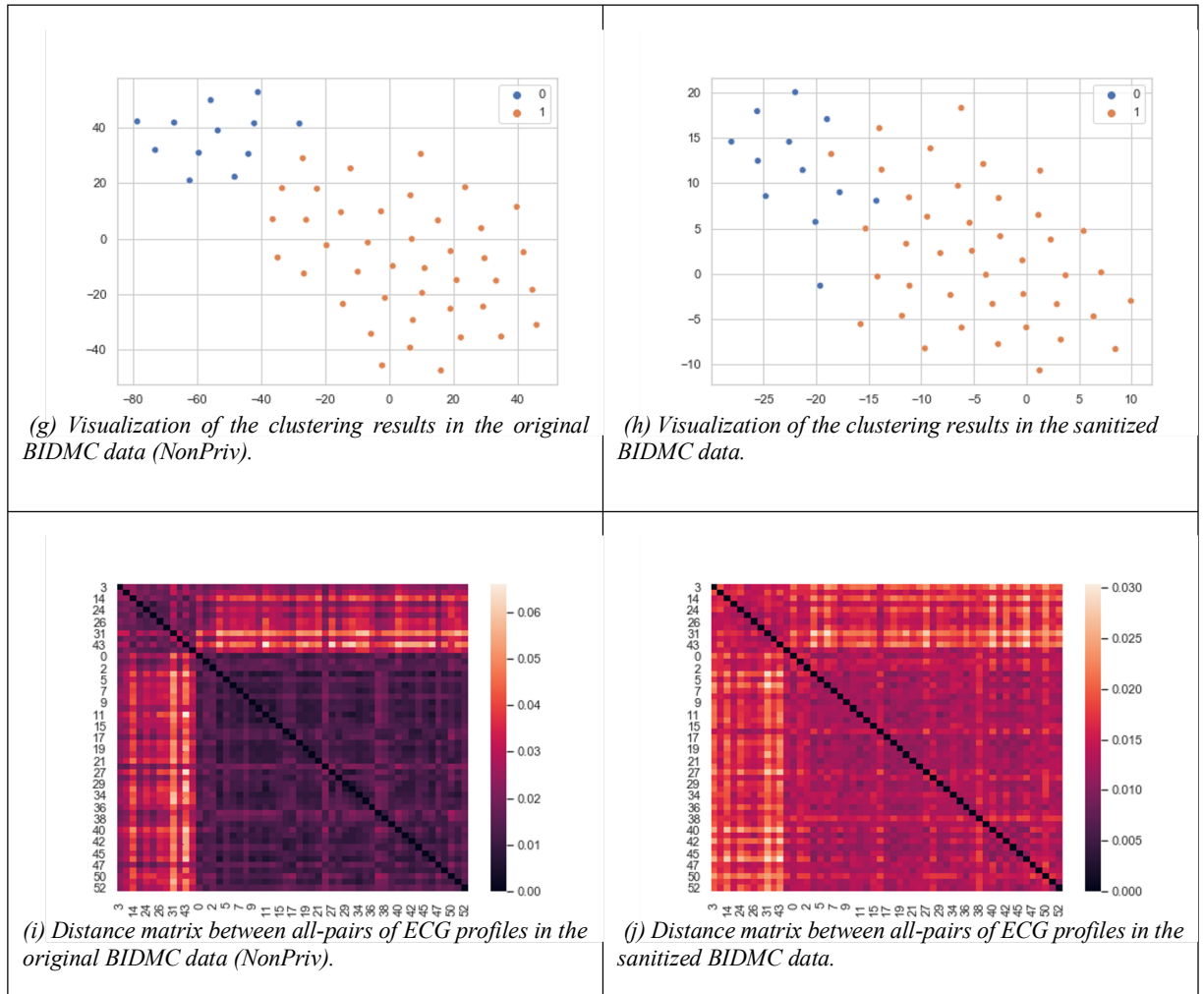


Figure 6. Distance-based clustering visualization on the BIDMC dataset. Using a hierarchical clustering approach, we identified two clusters (labeled 0 and 1) among all individuals. The clustering results for original data and sanitized data are visualized by t-SNE and are reported in (a) and (b), respectively. We observe that the general structure in clustering is preserved in the sanitized data. The all-pairs distance between ECG profiles both on the original and sanitized data are reported in (c) and (d), respectively. While the absolute distances between profiles are reduced in the sanitized data, the structure of distance matrix is preserved, which can be used to separate the clusters. The privacy parameter for obtaining the sanitized results was $\epsilon = 20$.

DISCUSSION

Our evaluations provide important insights on the design of privacy-protecting methods that can be deployed to protect individual-level health data.

We observed that ECG time-series data are highly unique in the datasets considered in this study. Our results show that a coarser data granularity can help reducing the uniqueness in the data, but it may not prevent membership inference attacks. For example, deterministically reducing the domain size had limited effects on the attack accuracy of the informed adversary considered in this work (Figure 3(b)). Our proposed sanitization method uses random sampling to introduce uncertainty in the output time-series data, reducing

the ability of the adversary to link the target's data with the sanitized output, thus reducing the success of membership inference.

In the differential privacy model, large values of the privacy parameter (epsilon) indicate weaker privacy protection. As a result, small privacy parameters (e.g., $\epsilon \leq 1.0$) may be used to provide strong privacy guarantees, at a high cost of data usability. In our evaluations, we show that our solution can successfully mitigate exiting inference attacks even for large values of the privacy parameter, enabling useful analytics (e.g., clustering, and predictive tasks) otherwise impossible with small privacy parameters. Future research efforts in studying the gap between theoretical privacy guarantees and practical adversarial models could provide useful insights in the design of privacy mechanisms, in order to find the right balance between privacy and usability.

Recent advances in technology have made it possible to collect fine-grained data directly from individuals (e.g., DTC genomics, wearable devices). Although this manuscript addresses some of the privacy issues in individual-level data sharing, more research in the areas of ethics, human-computer interaction, and education is needed to advance the design of individual-level privacy solutions.

CONCLUSION

In this work, we studied the applicability of the formal metric privacy model to provide privacy protection for individual-level ECG time-series data. Our evaluations demonstrated that our sanitization approach can provide strong privacy protection against powerful privacy attacks, and the aggregate ECG data can be used to develop predictive models and fine-grained analysis for cardiovascular diseases. Overall, our privacy study provides important insights on the development of privacy-protecting pipelines for collecting individual-level data and making them available for secondary use, which could facilitate emerging health applications (e.g., telemedicine and personalized medicine).

CONTRIBUTIONS

LB provided the motivation for this work, developed the method, contributed most of the writing, and conducted the experiments. ZW contributed with the usability evaluations and provided helpful comments. LF provided the motivation for this work, detailed edits, and critical suggestions.

COMPETING INTERESTS

None.

FUNDING

This work was supported by the National Human Genome Research Institute grant K99HG010493, National Institute of General Medical Sciences grant R01GM118609, and in part by the National Science Foundation Award number 2040727. LF is supported in part by the National Science Foundation CNS-1949217, CNS-1951430, and a UNC Charlotte FRG award.

REFERENCES

1. Dunn J, Runge R, Snyder M. Wearables and the medical revolution. *Per Med*. 2018;15(5):429-448.
2. Oresko JJ, Jin Z, Cheng J, et al. A wearable smartphone-based platform for real-time cardiovascular disease

- detection via electrocardiogram processing. *IEEE Trans Inf Technol Biomed.* 2010;14(3):734-740.
3. Sim I. Mobile devices and health. *N Engl J Med.* 2019;381(10):956-968.
4. Uddin M, Syed-Abdul S. Data analytics and applications of the wearable sensors in healthcare: An overview. *Sensors.* 2020;20(5):1379.
5. Ates HC, Yetisen AK, Güder F, Dincer C. Wearable devices for the detection of COVID-19. *Nat Electron.* 2021;4(1):13-14.
6. Quer G, Radin JM, Gadaleta M, et al. Wearable sensor data and self-reported symptoms for COVID-19 detection. *Nat Med.* 2021;27(1):73-77.
7. Investigators A of URP. The “All of Us” Research Program. *N Engl J Med.* 2019;381(7):668-676.
8. Irvine JM, Israel SA, Wiederhold MD, Wiederhold BK. A new biometric: human identification from circulatory function. In: *Joint Statistical Meetings of the American Statistical Association, San Francisco.* ; 2003:1957-1963.
9. Biel L, Pettersson O, Philipson L, Wide P. ECG analysis: a new approach in human identification. *IEEE Trans Instrum Meas.* 2001;50(3):808-812.
10. Kim J, Kim H, Bell E, et al. Patient Perspectives About Decisions to Share Medical Data and Biospecimens for Research. *JAMA Netw Open.* 2019;2(8):e199550--e199550. doi:10.1001/jamanetworkopen.2019.9550
11. Sufi F, Khalil I. Enforcing secured ecg transmission for realtime telemonitoring: A joint encoding, compression, encryption mechanism. *Secur Commun Networks.* 2008;1(5):389-405.
12. Sufi F, Khalil I, Hu J. ECG-based authentication. In: *Handbook of Information and Communication Security.* Springer; 2010:309-331.
13. Layouni M, Verslype K, Sandikkaya MT, De Decker B, Vangheluwe H. Privacy-preserving telemonitoring for ehealth. In: *IFIP Annual Conference on Data and Applications Security and Privacy.* Springer; 2009:95-110.
14. Poon CCY, Zhang Y-T, Bao S-D. A novel biometrics method to secure wireless body area sensor networks for telemedicine and m-health. *IEEE Commun Mag.* 2006;44(4):73-81.
15. Pandey S, Voorsluys W, Niu S, Khandoker A, Buyya R. An autonomic cloud environment for hosting ECG data analysis services. *Futur Gener Comput Syst.* 2012;28(1):147-154.
16. Bhalerao S, Ansari IA, Kumar A, Jain DK. A reversible and multipurpose ECG data hiding technique for telemedicine applications. *Pattern Recognit Lett.* 2019;125:463-473.
17. Goodrich MT. The mastermind attack on genomic data. In: *Security and Privacy, 2009 30th IEEE Symposium On.* IEEE; 2009:204-218.
18. Sweeney L. k-anonymity: A model for protecting privacy. *Int J Uncertainty, Fuzziness Knowledge-Based Syst.* 2002;10(05):557-570.
19. Dwork C. Differential privacy. *Int Colloq Autom Lang Program.* 2006;4052(d):1-12.
20. Papadimitriou S, Li F, Kollios G, Yu PS. Time series compressibility and privacy. In: *Proceedings of the 33rd International Conference on Very Large Data Bases.* ; 2007:459-470.
21. Fan L, Xiong L. An Adaptive Approach to Real-Time Aggregate Monitoring With Differential Privacy. *IEEE Trans Knowl Data Eng.* 2014;26(9):2094-2106.
22. Xiao X, Wang G, Gehrke J. Differential Privacy via Wavelet Transforms. *IEEE Trans Knowl Data Eng.* 2011;23(8):1200-1214. doi:10.1109/TKDE.2010.247
23. Beaulieu-Jones BK, Wu ZS, Williams C, et al. Privacy-preserving generative deep neural networks support clinical data sharing. *Circ Cardiovasc Qual Outcomes.* 2019;12(7):e005122.
24. Bonomi L, Huang Y, Ohno-Machado L. Privacy challenges and research opportunities for genomic data sharing. *Nat Genet.* 2020;1-9.
25. Alvim M, Chatzikokolakis K, Palamidessi C, Pazii A. Local differential privacy on metric spaces: optimizing the trade-off with utility. In: *2018 IEEE 31st Computer Security Foundations Symposium (CSF).* IEEE; 2018:262-267.
26. Chatzikokolakis K, Andrés ME, Bordenabe NE, Palamidessi C. Broadening the scope of differential privacy using metrics. In: *International Symposium on Privacy Enhancing Technologies Symposium.* Springer; 2013:82-102.
27. Andrés ME, Bordenabe NE, Chatzikokolakis K, Palamidessi C. Geo-indistinguishability: Differential privacy for location-based systems. In: *Proceedings of the 2013 ACM SIGSAC Conference on Computer & Communications Security.* ACM; 2013:901-914.
28. Xiang Z, Ding B, He X, Zhou J. Linear and range counting under metric-based local differential privacy. In: *2020 IEEE International Symposium on Information Theory (ISIT).* IEEE; 2020:908-913.
29. Fan L, Bonomi L. Time Series Sanitization with Metric-based Privacy. In: *IEEE Big Data Congress.* ;

- 2018:264-267.
30. Fan L. Practical image obfuscation with provable privacy. In: *2019 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE; 2019:784-789.
31. Thakurta AG, Vyrros AH, Vaishampayan US, et al. Emoji frequency detection and deep link frequency. July 2017.
32. Tang J, Korolova A, Bai X, Wang X, Wang X. Privacy Loss in Apple's Implementation of Differential Privacy on MacOS 10.12. *arXiv e-prints*. 2017:arXiv-1709.
33. Erlich Y, Narayanan A. Routes for breaching and protecting genetic privacy. *Nat Rev Genet*. 2014;15(6):409-421. doi:10.1038/nrg3723
34. El Emam K, Jonker E, Arbuckle L, Malin B. A systematic review of re-identification attacks on health data. *PLoS One*. 2011;6(12):e28071.
35. Fernandes AC, Cloete D, Broadbent MTM, et al. Development and evaluation of a de-identification procedure for a case register sourced from mental health electronic records. *BMC Med Inform Decis Mak*. 2013;13(1):1-14.
36. El Emam K, Rodgers S, Malin B. Anonymising and sharing individual patient data. *bmj*. 2015;350.
37. Choi E, Biswal S, Malin B, Duke J, Stewart WF, Sun J. Generating multi-label discrete patient records using generative adversarial networks. In: *Machine Learning for Healthcare Conference*. PMLR; 2017:286-305.
38. Bonomi L, Jiang X, Ohno-Machado L. Protecting patient privacy in survival analyses. *J Am Med Informatics Assoc*. 2020;27(3):366-375.
39. De Montjoye Y-A, Hidalgo CA, Verleysen M, Blondel VD. Unique in the crowd: The privacy bounds of human mobility. *Sci Rep*. 2013;3:1376.
40. Keogh E, Ratanamahatana CA. Exact indexing of dynamic time warping. *Knowl Inf Syst*. 2005;7(3):358-386.
41. MIT-BIH Supraventricular Arrhythmia Database. <https://physionet.org/content/svdb/1.0.0/>. Published 1999. Accessed September 1, 2021.
42. Olszewski RT. *Generalized Feature Extraction for Structural Pattern Recognition in Time-Series Data*. Carnegie Mellon University; 2001.
43. Pimentel MAF, Johnson AEW, Charlton PH, et al. Toward a robust estimation of respiratory rate from pulse oximeters. *IEEE Trans Biomed Eng*. 2016;64(8):1914-1923.
44. Goldberger AL, Amaral LAN, Glass L, et al. PhysioBank, PhysioToolkit, and PhysioNet : Components of a New Research Resource for Complex Physiologic Signals. *Circulation*. 2000;101(23):e215-e220. doi:10.1161/01.CIR.101.23.e215
45. Kate RJ. Using dynamic time warping distances as features for improved time series classification. *Data Min Knowl Discov*. 2016;30(2):283-312.
46. Van der Maaten L, Hinton G. Visualizing data using t-SNE. *J Mach Learn Res*. 2008;9(11).