

HEalth: Privately Computing on Shared Healthcare Data

Leo de Castro, Erin Hales, Mimeo Xu
peai011@rhul.ac.uk

January 2020

1 Introduction and Motivation

Healthcare in the US is notoriously expensive. Compared to other Organisation for Economic Co-operation and Development (OECD) countries, US healthcare costs are one-third higher or more relative to GDP [7]. According to the Centre for Disease Control and Prevention [6], the average per capita cost of healthcare was \$10,739 for the year 2017. Several Machine Learning (ML) startups aim to improve healthcare by bringing automated expertise to hospitals, in areas such as brain imaging and cancer detection. Additionally, hospitals in the US cannot easily share data. Existing ML solutions often focus on training a model using data which has been obtained through an existing collaboration, which has been pre-processed to the required format by the clients. These solutions work around the data-sharing challenge for training rather than tackling it and as a result produce a static inference model which does not continuously adapt to changes.

It is important to maintain ethical standards for healthcare professionals, especially for regulatory agencies working to protect patients. However, such agencies can often have problems accessing appropriate data to compare between hospitals and doctors, as well as to assess data across many hospitals. This can lead to difficulties in fulfilling a mandate to audit doctors and hospitals. Additionally, there are other healthcare use cases for which assessing data across many hospitals is useful. Such data sharing can be beneficial for managing epidemics and treating rare diseases.

Homomorphic encryption can be brought to bear on these situations by allowing private training, as well as by enabling more secure data sharing and computation.

The initial goal of our work is to allow anomaly detection and discovery in a shared private data-set. In particular we will consider this in the context of healthcare data. To do this we will calculate aggregate statistics across data shared between many hospitals to enable us to consider ‘fairness’ of hospital admissions.

Let us consider hospital admissions across the US. Some summary data is provided in Table 1. The total number of hospital admission across the year 2018 was very large, and if we can harness all of this data we can calculate powerful results.

Total number of US hospitals	6146
Staffed beds in all US hospitals	924,107
Total Admissions in all US Hospitals	36,353,946
Total expenses for all US hospitals	\$1,112,207,387,000

Table 1: Summary data for US hospital admissions in the year 2018, from [2].

We present a scenario where hospitals share records to compute anomalies and audit fairness. The contributions of our idea comprise three main categories. Firstly, fairness auditing at scale without the need to approximate non-linearity. Secondly, enabling continuous sharing of data without the need to refresh keys, meaning that we no longer require pairwise contracts. Finally, we apply an instance of anomaly detection from security management to finding causes and correlations in medical research. This has many applications, for example rare diseases, epidemic management and chronic illness research.

We use ‘fairness’ as an initial calculation goal which would demonstrate the effectiveness of simple statistics. It would then be possible to extend the techniques further to do different calculations for other goals.

2 Our Scenario

Suppose we have a group of hospitals who will combine their data to calculate aggregate statistics. Each hospital has an interest in keeping their data private from the other hospitals, so we introduce a method whereby the hospitals can work together to create a shared secret key to encrypt the shared data on which computation takes place. Then once the private computation has taken place, the hospitals are able to collectively decrypt the results using their shares of the decryption key.

We are calculating these aggregate statistics to compute anomalies and audit fairness in admissions statistics. This is a particular use case we have chosen to focus on, since this will bring a definite benefit to the patients involved. However, this is just an illustrative example and the same techniques could be applied to other scenarios where data is shared between several parties and some private computation takes place on the data.

We now consider the motivation of hospitals to take part in this process. Firstly, this data sharing process will allow hospitals to have access to statistics which encompass far more data than they would have if working independently. Additionally, hospitals will be able to compare themselves to others, obtaining more data for research and improvement.

A potential deployment issue is that different hospitals may store their data differently. So there may be some requirement for participating hospitals to share their data in a particular format. For example, we may need to format the data gathered in the format of Figure 1 as a matrix or a table. It is possible for this data pre-processing to be automated from the existing form. Additionally, as the scheme is used more frequently in a hospital it would increase efficiency if the data was gathered in a uniform way across participating hospitals to remove the need for pre-processing.

3 A Discussion of the Underlying Cryptography

We propose to use the method of multiparty communication for threshold FHE, as introduced in [1]. Asharov et Al’s work is an extension of existing FHE schemes ([5], [4]). In the multiparty communication setting, Key Generation and Decryption become N -party protocols rather than algorithms. Since we have a threshold encryption scheme, we require N out of the M parties to cooperate in order to encrypt and decrypt.

Each of the M participating hospitals would be a party of the threshold scheme, and we would require N hospitals to cooperate in order to decrypt the results of the computations. We must consider the case where a hospital no longer wishes to participate. In this case, the hospital can destroy their share of the key, and other parties will still be able to decrypt the results. The data

```
patient_id: 0
timestamp: 0
age: 37
gender: M
race: Caucasian
medical_history: [severe abdominal pain]
occupancy_at_admission: 70%
reason_for_visit: high fever
.
.
.
decision : admit
```

Figure 1: Example admissions data

of the non-participating hospital will remain in the dataset until that round of computations on the data is complete. It will then be possible to remove the data of one of the hospitals and begin computations again with fresh keys on a new set of data.

4 The initial goal: Fairness

As an initial computation we consider fairness. Here it is possible to compute relatively simple statistics on encrypted data which will have a large impact on healthcare.

The initial goal of aggregating the data and calculating statistics on it is to measure ‘fairness’ of hospital admissions. We will measure fairness within the context of protected characteristics, since these are recorded by hospitals already. We need to use such indicators in order to summarise complex statistical information and ideas of fairness and make them accessible to a non-technical audience. [11]

Since we have available to us many competing frameworks, coming from the different protected characteristics and across different policy domains, it is challenging to conceptualise fairness. For example, what weight do we give to different ‘strands’ of protected characteristics when we evaluate fairness? We aim to consider equality as “an outcome of equal treatment” [11]. In the context of medical data, the concept of ‘treatment’ takes on additional significance.

We conclude that while fairness may vary between different frameworks, if an admission decision lies too far from the norm then it is more likely to be due to unfair practices rather than noise. Over the volume of data that our method allows us to evaluate we will be able to get a better picture of what data stands out.

We will require data submitted by hospitals to be in the same form, for example in Figure 1. We must consider that the features we choose to record should encapsulate what the doctor is seeing and the information they use to make their decisions, as well as being similar to what a doctor would previously have recorded. We can ensure the data submitted by hospitals takes the same format without pre-processing by adapting the interface the data is entered into.

We wish to calculate density-based statistics, and to begin with we calculate simple statistics such

as averages and histograms. In our work so far using histograms we have only calculated one dimensional metrics, but these simple implementations highlight the potential for managing higher dimensional data in a usable manner. In addition, when compared to regression based methods our density-based statistics have a much weaker selection bias since averaging controls for bias.

5 Discussion

This application of threshold HE allows us to calculate healthcare statistics at a scale not seen before. The results of our calculations aim to be readily interpretable, with adaptive decision making rather than a blind prediction or classification. This is a novel use case for the existing research taking place in Threshold FHE [1, 10, 9, 3]. Our application will hopefully be able to leverage developments and optimisations in the field of threshold encryption.

In terms of encrypting and sharing the data, the key benefit of our approach is that it allows hospitals to share data, and provides all participating hospitals with an assurance that the data cannot be unlocked without all parties cooperating.

Hospitals have a strong privacy incentive to engage with the calculations, since uploaded data is encrypted and cannot be decrypted by one of the hospitals. This is because we require participation of many parties to decrypt, since each hospital only has access to their own share of the decryption key.

Another particular benefit of our approach is that as long as the hospitals retain their keys, additional algorithms could be developed to compute collective statistics, including implementing additional algorithms [8] on historical data without revealing secrets.

As it stands, our scheme does not consider malicious users. Future work would analyse how the system would behave if users were malicious or honest but curious, and how such malicious users could work together to compromise security. Additionally, future work could explore what would happen if hospitals were to submit some false or corrupted data, or an entirely false dataset. It would be useful to see how much false or corrupted data the system could tolerate.

References

- [1] Gilad Asharov, Abhishek Jain, Adriana López-Alt, Eran Tromer, Vinod Vaikuntanathan, and Daniel Wichs. Multipart computation with low communication, computation and interaction via threshold fhe. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, pages 483–501. Springer, 2012.
- [2] American Hospital Association. *Fast Facts on U.S. Hospitals, 2020*, (accessed January 14, 2020). <https://www.aha.org/statistics/fast-facts-us-hospitals>.
- [3] Dan Boneh, Rosario Gennaro, Steven Goldfeder, Aayush Jain, Sam Kim, Peter MR Rasmussen, and Amit Sahai. Threshold cryptosystems from threshold fully homomorphic encryption. In *Annual International Cryptology Conference*, pages 565–596. Springer, 2018.
- [4] Zvika Brakerski, Craig Gentry, and Vinod Vaikuntanathan. (leveled) fully homomorphic encryption without bootstrapping. *ACM Transactions on Computation Theory (TOCT)*, 6(3):13, 2014.
- [5] Junfeng Fan and Frederik Vercauteren. Somewhat practical fully homomorphic encryption. *IACR Cryptology ePrint Archive*, 2012:144, 2012.
- [6] Centers for Disease Control and Prevention. *Health Expenditures, 2017* (accessed January 14, 2020). <https://www.cdc.gov/nchs/fastats/health-expenditures.htm>.
- [7] Organisation for Economic Co-operation and Development. *Health expenditure and financing, 2020* (accessed January 14, 2020). <https://stats.oecd.org/Index.aspx?DataSetCode=SHA>.

- [8] Susan Graham, Deborah Estrin, Eric Horvitz, Isaac Kohane, Elizabeth Mynatt, and Ida Sim. Information technology research challenges for healthcare: From discovery to delivery. *ACM SIGHIT Record*, 1(1):4–9, 2011.
- [9] Aayush Jain, Peter MR Rasmussen, and Amit Sahai. Threshold fully homomorphic encryption. *IACR Cryptology ePrint Archive*, 2017:257, 2017.
- [10] Berry Schoenmakers. *Threshold Homomorphic Cryptosystems*, pages 1293–1294. Springer US, Boston, MA, 2011.
- [11] Sylvia Walby and Jo Armstrong. Developing key indicators of ‘fairness’: Competing frameworks, multiple strands and ten domains – an array of statistics. *Social Policy and Society*, 10(2):205–218, 2011.