

# 資訊安全 Final Project

5105056013 吳嘉偉、5105056029 江青霞、5105056019 廖健智

2017/06/25

## 1 Introduction

Many online services let users query large public datasets: some examples include restaurant sites, product catalogs, stock quotes, and searching for directions on maps. In these services, any user can query the data, and the datasets themselves are not sensitive. However, web services can infer a great deal of identifiable and sensitive user information from these queries, such as her current location, political affiliation, sexual orientation, income, etc.

or mirror site. For a given user query, all the providers have to run it on the same view of the data. A user splits her query into shares, using the Splinter client, and submits each share to a different provider. The user can select any providers of her choice that host the dataset. The providers use their shares to execute the user's query over the cleartext public data, using the Splinter provider library. As long as one provider is honest (does not collude with others), the user's sensitive information in the original query remains private. When the user receives the responses from the providers, she combines them to obtain the final answer to her original query.

## 2 Proposed Method

### 2.1 Splinter

#### 2.1.1 Architecture

There are two main principals in Splinter: the user and the providers. Each provider hosts a copy of the data. Providers can retrieve this data from a public repository

#### 2.1.2 Security Goals

The goal of Splinter is to hide sensitive parameters in a user's query. Specifically, Splinter lets users run parametrized queries, where both the parameters and query results are hidden from providers.

Splinter supports a subset of the SQL language, hides the information represented by the questions marks and the query's re-

sults, but the column names being selected and filtered are not hidden.

### 2.1.3 Threat Model

Splinter keeps the parameters in the user's query hidden as long as at least one of the user-chosen providers does not collude with others. Splinter also assumes these providers are honest but curious: a provider can observe the interactions between itself and the client, but Splinter does not protect against providers returning incorrect results or maliciously modifying the dataset.

It assumes that the user communicates with each provider through a secure channel (e.g., using SSL), and that the user's Splinter client is uncompromised. The cryptographic assumptions are standard. They only assume the existence of one-way functions in our two-provider implementation. In our implementation for multiple providers, the security of Paillier encryption is also assumed.

## 2.2 Function Secret Sharing

Function Secret Sharing lets a client divide a function  $f$  into function shares  $f_1, f_2, \dots, f_k$  so that multiple parties can help evaluate  $f$  without learning certain of its parameters.

## 3 Case Studies

### 3.1 Restaurant review site

We implement a restaurant review site using the Yelp academic dataset. The original dataset contains information for local businesses in 10 cities, but we duplicate the dataset 4 times so that it would approximately represent local businesses in 40 cities. We use the following columns in the data to perform many of the queries expressible on Yelp: name, stars, review count, category, neighborhood and location. For location-based queries, e.g., restaurants within 5 miles of a user's current location, multiple interval conditions on the longitude and latitude would typically be used. To run these queries faster, we quantize the locations of each restaurant into overlapping hexagons of different radii (e.g., 1, 2 and 5 miles), following the scheme from [1]. We precompute which hexagons each restaurant is in and expose these as additional columns in the data (e.g., hex1mi and hex2mi). This allows the location queries to use '=' predicates instead of intervals. For this dataset, we present results for the following three queries: Q1: SELECT COUNT(\*) WHERE category="Thai" Q2: SELECT TOP 10 restaurant WHERE category="Mexican" AND (hex2mi=1 OR hex2mi=2 OR hex2mi=3) ORDER BY stars Q3: SELECT restaurant, MAX(stars) WHERE category="Mexican" OR category="Chinese" OR category="Indian" OR category="Greek" OR category="Thai" OR

category="Japanese" GROUP BY category Q1 is a count on the number of Thai restaurants. Q2 returns the top 10 Mexican restaurants within a 2 mile radius of a user-specified location by querying three hexagons. We assume that the provider caches the intermediate table for the Top 10 query as described in Section 5.2.3 because it is a common query. Finally, Q3 returns the best rated restaurant from a subset of categories. This requires more communication than other queries because it performs a MAX with many disjoint conditions, as described in Section 5.2.2. Although most queries will probably not have this many disjoint conditions, we test this query to show that Splinter's protocol for this case is also practical.

### 3.2 Flight search

Implement a flight search service similar to Kayak, using a public flight dataset. The columns are flight number, origin, destination, month, delay, and price. To find a flight, we search by origin- destination pairs.

### 3.3 Map routing

We implement a private map routing service, using real traffic map data from for New York City. However, implementing map routing in Splinter is difficult because the providers can perform only a restricted set of operations. The challenge is

to find a shortest path algorithm compatible with Splinter. Fortunately, extensive work has been done to optimize map routing. One algorithm compatible with Splinter is transit node routing (TNR), which has been shown to work well in practice. In TNR, the provider divides up a map into grids, which contain at least one transit node, i.e. a transit node that is part of a "fast" path. There is also a separate table that has the shortest paths between all pairs of transit nodes, which represent a smaller subset of the map. To execute a shortest path query for a given source and destination, the user can use FSS to download the paths in her source and destination grid. She locally finds the shortest path to the source transit node and destination transit node. Finally, she queries the provider for the shortest path between the two transit nodes. We used the source code from and identified the 1333 transit nodes. We divided the map into 5000 grids, and calculated the shortest path for all transit node pairs. The grid table has 5000 rows representing the edges and nodes in a grid, and the transit node table has about 800,000 rows representing the number of shortest paths for all transit node pairs.

A user issues 2 grid queries and one transit node query. The two grid queries are issued together in one message, so there are a total of 2 network round trips. and transit node table. One observation is that the multi-party version is slightly faster than the two party version because it is faster at processing the grid query. The two-party version of FSS requires using GMP oper-

ations, which is slower than integer operations used in the multi-party version, the two-party version requires much less bandwidth.

## 4 Related Work

### 4.1 PIR(Private Information Retrieval ) systems

Splinter is most closely related to systems that use Private Information Retrieval(PIR) to query a database privately. In PIR, a user queries for the  $i$ th record in the database, and the database does not learn the queried index  $i$  or the result. Much work has been done on improving PIR protocols. Work has also been done to extend PIR to return multiple records, but it is computationally expensive. Our work is most closely related to the system in, which implements a parametrized SQL-like query model similar to Splinter using PIR. However, because this system uses PIR, it has up to  $10\times$  more round trips and much higher response times for similar queries.

Popcorn[4] is a media delivery service that uses PIR to hide user consumption habits from the provider and content distributor. However, Popcorn is optimized for streaming media databases, like Netflix, which have a small number (about 8000) of large records. The systems above have a weaker security model: all the providers

need to be honest. Splinter only requires one honest provider, and it is more practical because it extends Function Secret Sharing (FSS) [5, 6], which lets it execute complex operations such as sums in one round trip instead of only extracting one data record at a time.

### 4.2 Garbled circuits

Systems such as Embark [7], BlindBox [8], and private shortest path computation systems [9] use garbled circuits [10, 11] to perform private computation on a single untrusted server. Even with improvements in practicality [12], these techniques still have high computation and bandwidth costs for queries on large datasets because a new garbled circuit has to be generated for each query. (Reusable garbled circuits [13] are not yet practical.) For example, the recent map routing system by Wu et al. [9] uses garbled circuits and has  $100\times$  higher response time and  $10\times$  higher bandwidth cost than Splinter.

### 4.3 Encrypted data systems

Systems that compute on encrypted data, such as CryptDB [14], Mylar [15], SPORC [16], Depot [17], and SUNDR [18], all try to protect private data against a server compromise, which is a different problem than what Splinter tries to solve. CryptDB is most similar to Splinter because it allows for SQL-like queries over encrypted

data. However, all these systems protect against a single, potentially compromised server where the user is storing data privately, but they do not hide data access patterns. In contrast, Splinter hides data access patterns and a user’s query parameters but is only designed to operate on a public dataset that is hosted at multiple providers.

#### 4.4 ORAM(Oblivious RAM) systems

Splinter is also related to systems that use Oblivious RAM [19, 20]. ORAM allows a user to read and write data on an untrusted server without revealing her data access patterns to the server. However, ORAM cannot be easily applied into the Splinter setting. One main requirement of ORAM is that the user can only read data that she has written. In Splinter, the provider hosts a public dataset, not created by any specific user, and many users need to access the same dataset.

## 5 Experimental Result

Can Splinter be used practically for real applications?

Built and evaluated clones of three applications on Splinter: restaurant reviews, flight data, and map routing, using real datasets. Can support realistic applica-

tions including the search features of Yelp and flight search sites, and data structures required for map routing. Achieves end-to-end latencies below 1.6 sec for queries in these applications on realistic data. Use up to 10x fewer round trips than prior systems and have lower response times.

## 6 Conclusion

Splinter is a new private query system that protects sensitive parameters in SQL-like queries while scaling to realistic applications. Splinter uses and extends a recent cryptography primitive, Function Secret Sharing (FSS), allowing it to achieve up to an order of magnitude better performance compared to previous private query systems. [1] Develop protocols to execute complex queries with low computation and bandwidth. As a proof of concept, it has evaluated Splinter with three sample applications—a Yelp clone, map routing, and flight search—and showed that Splinter has low response times from 50 ms to 1.6 seconds with low hosting costs.

## 7 Reference

[1] Splinter: Practical Private Queries on Public Data Frank Wang, Catherine Yun, Shafi Goldwasser, Vinod Vaikuntanathan, Matei Zaharia<sup>†</sup> MIT CSAIL, <sup>†</sup>Stanford InfoLab

- [2] E. Boyle, N. Gilboa, and Y. Ishai. Function secret sharing. In Proceedings of the 34th Annual International Conference on the Theory and Applications of Cryptographic Techniques (EUROCRYPT), pages 337-367. Sofia, Bulgaria, Apr. 2015.
- [3] N. Gilboa and Y. Ishai. Distributed point functions and their applications. In Proceedings of the 33rd Annual International Conference on the Theory and Applications of Cryptographic Techniques (EUROCRYPT), pages 640 - 658. Copenhagen, Denmark, May 2014.
- [4] T.Gupta, N.Crooks, S.T.Setty, L.Alvisi, and M.Walfish. Scalable and private media consumption with Popcorn. In Proceedings of the 13th Symposium on Networked Systems Design and Implementation (NSDI), pages 91 - 107, Santa Clara, CA, Mar. 2016.
- [5] E. Boyle, N. Gilboa, and Y. Ishai. Function secret sharing. In Proceedings of the 34th Annual International Conference on the Theory and Applications of Cryptographic Techniques (EUROCRYPT), pages 337-367. Sofia, Bulgaria, Apr. 2015.
- [6] N. Gilboa and Y. Ishai. Distributed point functions and their applications. In Proceedings of the 33rd Annual International Conference on the Theory and Applications of Cryptographic Techniques (EUROCRYPT), pages 640 - 658. Copenhagen, Denmark, May 2014.
- [7] C. Lan, J. Sherry, R. A. Popa, S. Ratnasamy, and Z. Liu. Embark: Securely outsourcing middleboxes to the cloud. In Proceedings of the 13th Symposium on Networked Systems Design and Implementation (NSDI), pages 255 - 273, Santa Clara, CA, Mar. 2016.
- [8] J. Sherry, C. Lan, R. A. Popa, and S. Ratnasamy. Blind-box: Deep packet inspection over encrypted traffic. In Proceedings of the 2015 ACM SIGCOMM, pages 213-226, London, United Kingdom, Aug. 2015.
- [9] D. J. Wu, J. Zimmerman, J. Planul, and J. C. Mitchell. Privacy-preserving shortest path computation. In Proceedings of the 2016 Annual Network and Distributed System Security Symposium, San Diego, CA, Feb. 2016.
- [10] M. Bellare, V. T. Hoang, and P. Rogaway. Foundations of garbled circuits. In Proceedings of the 19th ACM Conference on Computer and Communications Security (CCS), pages 784 - 796, Raleigh, NC, Oct. 2012.
- [11] S. Goldwasser. Multiparty computations: past and present. In Proceedings of the 16th Annual ACM Symposium on Principles of Distributed Computing (PDC), pages 1-6, 1997.
- [12] M. Bellare, V. T. Hoang, S. Keelveedhi, and P. Rogaway. Efficient garbling from a fixed-key blockcipher. In Proceedings of the 34th IEEE Symposium on Security and Privacy, pages 478 - 492, San Francisco,

CA, May 2013.

2010.

- [13] S. Goldwasser, Y. Kalai, R. A. Popa, V. Vaikuntanathan, and N. Zeldovich. Reusable garbled circuits and succinct functional encryption. In Proceedings of the 45th Annual ACM Symposium on Theory of Computing (STOC), pages 555-564, Palo Alto, CA, June 2013.
- [14] R.A.Popa, C.M.S.Redfield, N.Zeldovich, and H.Balakrishnan. CryptDB: Protecting confidentiality with encrypted query processing. In Proceedings of the 23rd ACM Symposium on Operating Systems Principles (SOSP), pages 85 - 100, Cascais, Portugal, Oct. 2011.
- [15] R. A. Popa, E. Stark, J. Helfer, S. Valdez, N. Zeldovich, M. F. Kaashoek, and H. Balakrishnan. Building web applications on top of encrypted data using Mylar. In Proceedings of the 11th Symposium on Networked Systems Design and Implementation (NSDI), pages 157-172, Seattle, WA, Apr. 2014.
- [16] A.J.Feldman, W.P.Zeller, M.J.Freedman, and E.W.Felten. SPORC: Group collaboration using untrusted cloud resources. In Proceedings of the 9th Symposium on Operating Systems Design and Implementation (OSDI), Vancouver, Canada, Oct. 2010.
- [17] P. Mahajan, S. Setty, S. Lee, A. Clement, L. Alvisi, M. Dahlin, and M. Wal-fish. Depot: Cloud storage with minimal trust. In Proceedings of the 9th Symposium on Operating Systems Design and Implementation (OSDI), Vancouver, Canada, Oct. 2010.
- [18] J. Li, M. Krohn, D. Mazières, and D. Shasha. Secure untrusted data repository (SUNDR). In Proceedings of the 6th Symposium on Operating Systems Design and Implementation (OSDI), pages 91 - 106, San Francisco, CA, Dec. 2004.
- [19] J. R. Lorch, B. Parno, J. Mickens, M. Raykova, and J. Schiffman. Shroud: Ensuring private access to largescale data in the data center. In Proceedings of the 11th USENIX Conference on File and Storage Technologies (FAST), pages 199-213, San Jose, CA, Feb. 2013.
- [20] E. Stefanov, M. van Dijk, E. Shi, C. Fletcher, L. Ren, X. Yu, and S. Devadas. Path ORAM: An extremely simple oblivious RAM protocol. In Proceedings of the 20th ACM Conference on Computer and Communications Security (CCS), Berlin, Germany, Nov. 2013.