

Poster: Secure Multi-Party Computation as a Tool for Privacy-Preserving Data Analysis

Samuel Havron
University of Virginia
havron@virginia.edu

Abstract—A *Secure multi-party computation* (MPC) protocol allows two or more parties to compute a function on sensitive input data provided by both parties, without revealing anything about the inputs (other than what can be inferred from the revealed output result). Social scientists often work with private datasets that cannot be shared due to legal restrictions and ownership issues, but many interesting studies could be enabled if MPC allows joint data analyses on private data analysis. We are exploring opportunities for using MPC to enable scientific research that would not otherwise be possible. We have developed tools and libraries that enable joint scientific data analyses with private data, and report on preliminary results using MPC to enable linear regression analyses over private data.

I. INTRODUCTION

Many social scientists and researchers need to perform statistical analyses across large, independently-owned datasets for their work, but they are often met with difficulties in obtaining sensitive data and computing results in a safe manner. For instance, an education researcher may be interested in using statistical methods to analyze the relationship between family income and grade point average for a particular school system or collection of school systems. Obtaining such sensitive data is often difficult to do without trusting one or more parties with some of the private input data, as well as considering ownership and legal restrictions on having clear access to private data. One government agency has information about incomes, but cannot share it without violating regulations (and compromising privacy); a local school district has information about students' grades, but cannot correlate that with their family incomes. Secure multi-party computation (MPC) is a protocol which can be used as a tool for carrying out large-scale scientific data analysis in these situations without compromising the privacy of any party's data, or risking it being exposed.

II. APPROACH

Several approaches to MPC have been developed, including application-specific custom protocols and generic, universal techniques that can be used to compute any function privately. We use universal techniques since it is important that new functions can be developed quickly and without needing new security proofs, and we anticipate needing to incorporate auditing and other functionality into the MPC. The most common universal MPC techniques are based on secret sharing, homomorphic encryption, and garbled boolean

circuits using Yao's protocol. We use garbled circuits, which are generally the most scalable and high performance MPC approach currently known.

Obliv-C (<http://oblivc.org>) is a programming language which allows an application developer to quickly implement scalable, secure MPC protocols, using the languages API or writing specific functionality by extending the language's existing library as well as experimenting with the implementation of library protocols [7]. The language is compiled and built on top of the standard C language, allowing for developers to integrate C tools and libraries with Obliv-C seamlessly. Obliv-C provides an implementation of Yao's garbled circuit protocol for use with semi-honest adversaries, although it can also be used to implement other protocols. Using this language to write applications that can analyze large, privacy-preserving datasets results of building one such application for linear regression analysis are shown in Preliminary Results.

The goal of using Obliv-C as a framework for secure computation is to demonstrate its performance capabilities and ease of use to developers whom have little knowledge of cryptography or circuit structures, but would greatly benefit from using secure MPCs to carry out privacy-preserving dataset analysis (perhaps for social science research). An example of Obliv-C code which calculates and reveals the correlation coefficient of a linear regression are seen in Figure 1. The *obliv* qualifier built into Obliv-C's type system ensures the variable used is encrypted with the garbled circuit scheme [7]; in this case, the variable is the correlation coefficient of linear regression analysis, obtained through some function call to a method within the MPC protocol. The result is revealed to specified parties through the API *RevealOblivInt()*, which works by decrypting the *obliv* value into a normal C struct.

III. PRELIMINARY RESULTS

A. Scalability

A linear regression data analysis program was developed to test the scalability and speed of Obliv-C as a tool for

```
obliv int orsqr = getOblivRSquared();  
revealOblivInt(&io -> rsqr, orsqr, 0);
```

Fig. 1. Code snippet showing *obliv* qualifier and reveal API. Revealing an integer stores the result into a struct which is accessible to specified parties ("0" means all parties in this context).

implementing MPC programs. Testing was done using `c4.large Elastic Compute Cloud` (EC2) nodes from *Amazon Web Services* (AWS) [1], which feature high frequency Intel Xeon E5-2666 v3 (Haswell) processors optimized specifically for EC2, two vCPUs, and 3.75 GiB of DRAM. Two `c4.large` instances were launched and connected through Obliv-C's API for TCP/IP connections via Oblivious Transfer protocol. The instances were both located in the same cloud cluster in Oregon; exploring network latency impositions with secure MPCs is also an area of interest.

One node instance provided independent (x) data points, while the other provided dependent (y) data points; data points used were 32-bit integers, using fixed-point mathematics to convert raw data values into scaled integers, as Obliv-C does not currently support floating point numbers. The time needed to execute the MPC between instances appears to scale linearly with the size of the data input; 100K data points finished execution in 12.7 minutes on average, 500K completed in 63.7 minutes, and 1 million data points finished execution in just over 127 minutes on average. Given the relative cost of using a `c4.large` instance (as of writing, \$0.105 per hour), executing this program over larger inputs, such as 10 million data points, will only incur about \$4.45 between two instances in the estimated runtime of 21 hours.

Artificial data was generated for testing the scalability of input size, and considerations to automated data match-ups between two separate datasets were not implemented; the artificial data was presumed to already be matched and sorted properly. To provide a clear example of the utility of Obliv-C for analyzing sensitive datasets, additional data for computation was obtained from the public New York State Department of Health dataset of Hospital Inpatient Discharges from 2011 [2]. Comparisons between fields such as "Length of Days stayed" and "Total Costs" over approximately 2.6 million data points are currently being tested. This dataset is particularly amenable to analysis, as all data is already matched properly and each data value is a comparable number.

B. Discussion

All numbers are averages over 5 executions between `c4.large` instances. The datasets used in Preliminary Results assume that data is matched, sorted, and comprises of comparable values. Using private set intersection to assist matching data with a unique identifier and automatically filtering out incompatible data points is considered for future work. The proposed health dataset is publicly available and in a single database. A more realistic situation would be for the independent and dependent data points to be each residing in separate datasets, owned by different organizations, and being tied to legal restrictions for sharing the data.

IV. RELATED WORK

A similar approach in taking distinct federal datasets and comparing them is seen in Dan Bogdanov *et al.*, where correlations between working hours and failure to graduate on time in Estonia was investigated, matching over 10 million tax records

and 500K education records [4]. This analysis utilized the researchers' own framework for secure computation, *ShareMind*, a database and analytics system which almost exclusively uses three-parties to carry out computations, and uses arithmetic manipulations to implement secure MPC, rather than boolean circuit evaluation [3].

Another privacy-preserving approach to analyzing millions of records uses a combination of homomorphic encryption (for linear computations) and Yao circuits (for non-linear computations) in order to compute ridge regression [6]. This approach is designed for many users to send data to a central server called the *Evaluator*, in contrast to the primarily two-party model presented in the Obliv-C language. Yet another approach used for secure multiple linear regression relies on protocols based on homomorphic secret sharing, and data which is partitioned across several databases [5].

V. CONCLUSION

Using MPCs for scientific analysis of large datasets is promising for social scientists and researchers whom would otherwise need to reveal some of one party's input to another party in order to analyze their data or be unable to perform analysis due to legal restrictions and ownership issues. The Obliv-C language is particularly well-suited for implementing such scalable, sensitive data analysis between two or more parties, and shows the practicality of using secure MPCs for processing large datasets. Future work invites a closer examination of automatic data-matching between separate datasets with private set intersection, improving fixed-point integer conversion for decimal data values used in computation, and other privacy-preserving applications.

REFERENCES

- [1] <https://aws.amazon.com/ec2/instance-types/>.
- [2] <https://health.data.ny.gov/>.
- [3] Dan Bogdanov, S. Laur, and J. Willemson, *Sharemind: A framework for fast privacy-preserving computations*, In 13th European Symposium on Research in Computer Security (ESORICS), volume 5283 of LNCS, pages 192206. Springer, 2008. <http://sharemind.cyber.ee>.
- [4] Dan Bogdanov, Liina Kamm, Baldur Kubo, Reimo Rebane, Ville Sökk, Riivo Talviste, *Students and Taxes: a Privacy-Preserving Social Study Using Secure Computation*, Cryptology ePrint Archive, Report 2015/1159, 2015, <http://eprint.iacr.org/>.
- [5] Rob Hall, Stephen E. Fienberg, and Yuval Nardi, *Secure multiple linear regression based on homomorphic encryption*, J. Official Statistics, vol. 27, no. 4, 2011.
- [6] Valeria Nikolaenko, Udi Weinsberg, Stratis Ioannidis, Marc Joye, Dan Boneh, Nina Taft, *Privacy-Preserving Ridge Regression on Hundreds of Millions of Records*, in 2013 IEEE Symposium on Security and Privacy (S&P 2013), pp. 334-348, IEEE Computer Society, 2013, <http://www.technicolorbayarea.com/papers/2013/NWUJBT13garbled.pdf>.
- [7] Samee Zahur and David Evans, *Obliv-C: A Language for Extensible Data-Oblivious Computation*, Cryptology ePrint Archive, Report 2015/1153, 2015, <http://eprint.iacr.org/>.