Cryptographically Protected Database Search

Benjamin Fuller, Mayank Varia, Arkady Yerukhimovich, Emily Shen, Ariel Hamlin, Vijay Gadepally, Richard Shay, Darby Mitchell, Robert Cunningham

benjamin.fuller@uconn.edu







The Data Economy



The Rise of the Data Economy: Driving Value through Internet of Things Data Monetization

Technology Officers By Albert Opher, Alex Chou, Andrew Onda, and Krishna



Interesting takeaway No. 1: 61% of respondents "acknowledge that big data is now a driver of revenues in its own right and is becoming as valuable to their businesses as their existing products and services."

"Data is the new oil"

– Shivon Zilis, Bloomberg Beta

Sounderrajan

- "Data will become a currency"
- David Kenny, IBM Watson
- "... the fourth industrial revolution is connectivity and data"
- Mukesh Ambani, Reliance

Value implies Risk

The telecommunications company TalkTalk admitted that its data breach last year resulted in

criminals using customer information to commit fraud. This was more bad news for a company













Cryptographically Protected Search

Homomorphic encryption vector-matrix mult: $30s^1$ Multi-party computation: 200,000 AES² blocks/s, does not scale to large data No server protections (encrypt data at rest)

Databases are expected to answer common queries in milliseconds

SELECT count(*), avg(b) FROM t2 WHERE b>=0 AND b<100; SELECT count(*), avg(b) FROM t2 WHERE b>=100 AND b<200; SELECT count(*), avg(b) FROM t2 WHERE b>=200 AND b<300; ... 4994 lines omitted SELECT count(*), avg(b) FROM t2 WHERE b>=499700 AND b<499800; SELECT count(*), avg(b) FROM t2 WHERE b>=499800 AND b<499900; SELECT count(*), avg(b) FROM t2 WHERE b>=499900 AND b<500000;

5000 range queries takes 1s

Return whole dataset encrypted

Use homomorphic encryption or multi-party computation

¹S. Halevi and V. Shoup. (2014) HElib - an implementation of homomorphic

encryption. [Online]. Available: https://github.com/shaih/ Helib Utility of stored data

²M. Keller, E. Orsini, D. Rotaru, P. Scholl, E. Soria-Vazquez, S. Vivek,

Risk of data compromise

"Faster Secure Multi-Party Computation of AES and DES Using Lookup Tables," in ACNS 2017

Cryptographically Protected Search

Includes:

- Symmetric searchable encryption (SSE)
- Property preserving encryption

No server protections (encrypt data at rest)

Return whole dataset encrypted

Use homomorphic encryption or multi-party computation

Utility of stored data



Outline

- Overview of Protected Search
 - Leakage Impacts
- Finding a basis for search results
 - Range queries
 - Compatible approach: Order-Preserving Encryption / CryptDB
 - Custom approach: Partial Order-Preserving Encryption
 - Obliv approach: SisoSPIR
 - Combining queries
- Extending to new database paradigms

Common Language for Leakage

Protected search schemes reveal some information about the query, data set, and result set to *each* party.

Called leakage.

Difficult to compare, phrased to make proofs work, not to compare schemes Define five types of leakage of increasing impact¹:

- 1. Structure
- 2. Identifiers
- 3. Predicates
- 4. Equality
- 5. Order

Some schemes leak:

- 1. At *Initialization* on entire DB
- 2. At Query on relevant records

Hospital Data Set

Birth Month	Length of Stay	Gender	Diagnosis	SSN
February	1	Μ	Flu	000-00-001
April	30	Μ	Cancer	000-00-002
June	3	F	Pneumonia	000-00-003

• Assume:

- Server sees which field queried
- Records are identifiable between queries

Statistical Attack Against Hospital Length of Stay

- Suppose:
 - Queries of form: SELECT * FROM table WHERE length_stay=XXXXX;
 - Observe |records|
 - Create unique id for query



Statistical Attack Against Hospital Length of Stay

Query with highest number of returned records likely represents 4 days



Distribution of length of stay is known, attacker can use prior statistical information

Statistical Attack Against Hospital Length of Stay

Query with highest number of returned records likely represents 4 days



Or statistical prior is inaccurate?

Attacks exploit correlation between fields, use techniques from optimization

Why systematize?

No server protections (encrypt data at rest)

Return whole dataset encrypted

Use homomorphic encryption or multi-party computation

Utility of stored data

Approaches to Protected Databases

Define five types of leakage of increasing impact⁴:

- 1. Structure
- 2. Identifiers
- 3. Predicates
- 4. Equality
- 5. Order

Distinguish between schemes that leak this information at *Initialization* and at *Query*

⁴Partially based on S.Kamara, "Structured encryption and leakage suppression," presented at Encryption for Secure Search and Other Algorithms, Bertinoro, Italy, June 2015.

Find three approaches to protected databases:

- 1. Legacy:
 - Leak at Initialization
 - Inherit DB advances
- 2. Custom:
 - Leak during Query
- 3. Obliv:
 - Leak only structure
 - Require multiple servers to be efficient

Approaches to Protected Databases

- Developed⁹:
 - a database instrumentation platform
 - data and query generator
- Used in prior work^{10, 11}



⁹https://github.com/mit-ll/sparta

¹⁰V. Pappas et al. "Blind Seer: A Private Scalable DBMS," S&P 2014
¹¹D. Cash et al. "Dynamic Searchable Encryption in Very-Large Databases: Data Structures and Implementation," NDSS 2014

Find three approaches to protected databases:

- 1. Legacy:
 - Leak at Initialization
 - Inherit DB advances
- 2. Custom:
 - Leak during Query
- 3. Obliv:
 - Leak only structure
 - Require multiple servers to be efficient

How to compare functionality?

- Natural approach: what fraction of a unprotected database language is supported?
- Current systems implement *base* queries using cryptography, extend from these base queries:
 - Keyword Equality
 - Range
 - Boolean Combination
 - Other (graph alg and substring)

Find three approaches to protected databases:

- 1. Legacy:
 - Leak at Initialization
 - Inherit DB advances
- 2. Custom:
 - Leak during Query
- 3. Obliv:
 - Leak only structure
 - Require multiple servers to be efficient

Outline

- Overview of Protected Search
- Leakage Impacts
 - Finding a basis for search results
 - Range queries
 - Order-Preserving Encryption
 - Partial Order-Preserving Encryption
 - SisoSPIR
 - Combining queries
- Extending to new database paradigms

Order-Preserving Encryption

- Enc that preserves plaintext order:
 - If $m_1 \le m_2$ then $Enc(m_1) \le Enc(m_2)$

- 1. Encrypt query Enc(a), Enc(b)
- 2. Let server use standard search mechanism
- 3. Return encrypted records





Database server

Leakage Attacks of OPE

- Data is sorted, does not protect dense data
- Strongest leakage attack applies to OPE
- Technique used in many commercial product

	Required S leakage		Required conditions	attack	Attack efficacy						
Attacker goal	Init	Query	Ability	Prior	Runtime	Sensitivity	Keyword				
			to inject	knowledge		to prior	universe				
			uata			Kilowieuge	lesieu				
4	0	0		•		?	0				
aer	0	0	~	0	0	0	0				
aso o	0	O	—	O	•	?	0				
A	0	•	—	•	lacksquare	•					
Morri	0	O	~	•	lacksquare	0	\bullet				
<u>u</u>	0	•	—		lacksquare	•					
A	0	•	_	\bullet	•	•	\bullet				
-0 ²	•	—	_	O	0	?	0				
Aer -	•	—	~	0	0	?	•				
x2			_	4	0	?					
$\mathfrak{I}_{\mathfrak{G}}$		—			0	0					

Row corresponding to OPE

Partial Order Preserving Encoding¹³



- Client sends data to server encrypted and unsorted
- Client and Server work together to create partially sorted tree
 - Client performs all comparisons
 - Server is able to build tree based on client comparisons
- Stronger security than Order-Preserving Encryption if tree is only partially built

¹³D. Roche, D. Apon, S. Choi, A. Yerukhimovich "POPE: Partial Order Preserving Encoding" CCS 2016

SisoSPIR¹⁴ – Obliv Approach to Range



B+ trees are used in many unprotected databases Variable number of children per node

Idea of approach: use crypto to hide all information in traversing B+ tree

Requires multiple servers for practical efficiency

¹⁴Y. Ishai, E. Kushilevitz, S. Lu and R. Ostrovsky, "Private Large-Scale Databases with Distributed Searchable Symmetric Encryption," CT-RSA 2015



ciphertexts

Query Combination

- Techniques to *combine base* queries:
 - Range \rightarrow Equality, search for [a, a]
 - Boolean \rightarrow Range, using set covers
 - Range \rightarrow Substring, by inserting each prefix
- Most combination techniques are less efficient and have more leakage than equivalent base query
- Allow for rapid expansion of query functionality

Approaches to Protected Databases

- Natural approach: what fraction of a unprotected database language is supported?
- Systems implement base queries w/ crypto, extend from these base queries:
 - Keyword Equality
 - Range
 - Boolean Combination
 - Other (graph alg and substring)

SQL has a well defined mathematical set-theory basis of operations¹⁴:

- Union: AUB
- Difference: A \ B
- Join: A x B
- Projection: Take some dimensions of results
- Selection: Take rows satisfying some condition

¹⁴E. Codd, "A relational model of data for large shared data banks," *Communications of the ACM*, 1970

Unprotected DB Development



Keeping up with database diversification

Common unprotected databases have a mathematical basis of operations:

- For SQL: Union, Difference, Join, Projection, Selection
- For Array-Store: Construct, Find, Array (+, x), Element-wise x
- For Graph: Linear algebra over matrices

Cryptographers and DB designers should work together to:

- Identify base queries that are likely to be useful across DB paradigms
- 2. Understand critical functions of emerging databases
- 3. Quickly fill gaps using *combiners*

Questions? https://arxiv.org/abs/1703.02014

Backups

DB Paradigm	Basis Operation	Crypto Base Operation?					
NoSQL – Key Value Store	Construct	Yes					
	Find	Yes – Mature range search with variety of techniques					
	Array (+, x)	Some – Addition possible using partially homomorphic techniques					
	Element Wise x	Some – Using partially homomorphic techniques					

Main gap is support for very high insert rates above 1M records per second

DB Paradigm	Basis Operation	Crypto Base Operation?					
Graph Databases– Linear Algebra	Construct	Yes					
Linear	Find	Yes – Mature range search with variety of techniques					
	Matrix (+, x)	Some – Have private algorithms for matrix mult./add.					
	Element Wise x	Some – Using homomorphic operations					

Current matrix operations operate on full structure, need algorithms for sparse matrices (most graph algorithms)

Current systems

Questions? https://arxiv.org/abs/1703.02014

- Currently mature systems with peer-reviewed descriptions
- All systems use the basis and combination approach to get rich functionality

	Equality	Boolean	Keyword	Range	Substring	Wildcard	Sum	Join	Update	Approach	# of parties	Code available	Multi-client	User auth.	Access control	Query policy	Leakage	Performance
System		Supported Operations							Prop	erties			Feat	ures				
CryptDB [15]			0							Logocy	2					0		
		-	\sim	-	0	\cup	-	•	•	Legacy	2	•	•	•	•	\cup		
Arx [14]	•	0	0	•	0	0	•	•	•	Custom	2	0	0	0	0	0	Ō	Ō
Arx [14] BLIND SEER [16], [17]	•	0	0	•	0	0	•	•	•	Custom Custom	2 2 3	0	• •	0	0 0	0	0	0
Arx [14] BLIND SEER [16], [17] OSPIR-OXT [18]–[21], [103], [104]	•		0 •	• • •		0	• • • •	• • • •	•	Custom Custom Custom	2 2 3 3	0 0 0	• •	0 0 0	0 0 0	0 •	0 0 0	