Secure storage in the cloud using property preserving encryption

#### **Kenny Paterson**

Information Security Group



ROYAL HOLLOWAY UNIVERSITY

#### Overview

- 1. Application scenarios.
- 2. Deterministic encryption and search.
- 3. OPE/ORE and range queries.
- **4.** Analysing access pattern leakage from range queries.



# Application scenarios

#### **Application Scenarios**

- Data owners wish to **securely** outsource storage to cloud providers whilst preserving capability for users to query data in various ways.
- What kinds of queries?
- What kinds of users?
- What kinds of data?
- What kinds of query?
- What kinds of adversary?
- Meta: Why not just use FHE and be done?

#### Two scenarios, one picture



#### Scenario 1: Searchable File Storage

- Owner has large collection of files, indexed by keywords.
- Owner **encrypts** files and stores these on remote server.
- Owner **encodes** keywords in such a way that keyword searches can still be carried out.
- Encoded keywords also stored on server, as an encoded **index**.
- Owner sends **search token** to server; server uses token and index to find identifiers for matching files.
- Matching file identifiers are returned to owner.

#### Scenario 2: Database Encryption

- Data owner has a large database of records; each record has multiple fields.
- Owner encrypts data in each field in such a way that **standard database queries** can still be carried out.
- **Basic**: simple searches.
  - "Give me all records in which surname = Dubois".
- Advanced: compound searches.
- More advanced: range queries
  - "Give me all records with ages between 21 and 30".
- Finally: arbitrary SQL queries.\*

\*Other db query languages are available.

#### Searchable Encryption

#### Solution for Scenario 1: Searchable Encryption.

- Naïve scheme: owner uses IND-CPA symmetric encryption for files and  $PRF_{k}(kw)$  as encoding of keyword kw.
- Store encrypted files and encoded keywords per file on server.
- Owner sends tok = PRF<sub>K</sub>(kw) to server; server matches tok against encoded keywords; returns matching files.
- Can use an inverted index and file identifiers: server stores database of tuples (*tok*, (*fid*<sub>1</sub>, *fid*<sub>2</sub>,....)).

#### Security Analysis

- Adversarial objectives?
  - Keyword recovery, recovery of file contents,...?
- Adversarial capabilities?
  - "Snapshot", "Honest-but-curious", "Fully malicious".
  - Can/cannot **observe** queries; can/cannot **make** queries; can/cannot **inject** files.
- What about auxiliary information?
  - What if the adversary has a representative data sample or keyword sample?
- Cash *et al*. (CCS15): detailed analysis of different attacks models, leakage profiles, etc. against **SE schemes** in general: **Leakage Abuse Attacks**.
- Fuller *et al.* (S&P17): SoK paper on cryptographically protected database search.

#### Two scenarios, one picture



#### **Deterministic Encryption**

#### Partial solution for Scenario 2: DE

- Simplest possible scheme: owner uses deterministic encryption scheme (KGen, Enc, Dec) to encrypt each column of the database using a per-column key *K*.
- Server can store the encrypted data on server in a traditional database.
- To find matches with value x in a column, send search query for y = Enc<sub>k</sub>(x) to server.
- Server finds matches on y and returns full encrypted records to client.
- Client decrypts returned records using per column keys.
- Use of DE **preserves equality** of plaintexts and allows simple searches.
- (Very similar to naïve SE, with PRF replaced by Enc/Dec).

#### Property Preserving/Revealing Encryption (PPE/PRE)

#### More general solution for Scenario 2: PPE/PRE

- Generalises idea of "equality preserving/revealing" property of DE.
- Main example: Order Preserving/Revealing Encrypion (OPE/ORE).
- OPE: if *x* < *y* then Enc(*x*) < Enc (*y*).
- **ORE**: there exists a (public) efficiently computable function "Order" such that:

#### x < y iff Order(Enc(x), Enc(y)) = 1

- OPE/ORE allows range queries!
- Client who wishes to query on range [*a*,*b*] instead sends query for range [Enc(*a*), Enc(*b*)] to server.



# Analysis of Deterministic Encryption

#### Reminder: ECB information leakage



Tux the Penguin, the Linux mascot. Created in 1996 by Larry Ewing with The GIMP. lewing@isc.tamu.edu



ECB-Tux

#### Analysis of Deterministic Encryption

- DE is equality preserving, by design.
- DE therefore preserves frequencies of plaintexts in the ciphertexts, cf. monoalphabetic substitution cipher.
- Naveed-Kamara-Wright (CCS15): let's apply frequency analysis! (al-Kindi, 9<sup>th</sup> century.)
- Assumption 1: attacker has auxiliary information a reasonably accurate estimate for the plaintext distribution.
- Assumption 2: attacker has a snapshot of the encrypted database.



#### Analysis of Deterministic Encryption



## Frequency Analysis is Maximum Likelihood!

- Given a column of ciphertexts **y**, frequency analysis matches:
  - Most frequent item in **y** with most frequent item in aux. dist.
  - Second most frequent item in y with second most frequent item in aux. dist.
  - etc.
- Defines a permutation π mapping plaintexts x to ciphertexts y.
- This procedure is **maximum likelihood**, that is, it maximises the likelihood

 $L(\pi \mid \boldsymbol{y}) := \Pr(\boldsymbol{y} \mid \pi).$ 

• **Proof**: fun exercise, see also eprint 2015/1158.

- Naveed-Kamara-Wright [CCS15] performed an **empirical investigation** of the performance of frequency analysis against DE.
- Using **a large medical dataset**: per-patient data in 12 categories for 200 largest hospitals in the 2009 Nationwide Inpatient Sample (NIS), from the Healthcare Cost and Utilization Project (HCUP), run by the US Agency for Healthcare Research and Quality.
- DE encrypt data per hospital for each category.
- Use 2004 **aggregated** HCUP data as the auxiliary data.
- Run frequency analysis and measure percentage of data items correctly recovered **per hospital**.

Attribute	Num.
	values
Age (AGE)	125
Admission month (AMONTH)	12
Admission source (ASOURCE)	5
Admission type (ATYPE)	6
Patient died (DIED)	2
Sex (FEMALE)	2
Length of stay (LOS)	365
Major diagnostic category (MDC)	25
Primary payer (PAY1)	6
Ethnicity group (RACE)	6
Disease severity (APRDRG_Severity)	4
Mortality risk (APRDRG_Risk_Mortality)	4





21

#### Frequency Analysis Makes Headlines!

Microsoft   TechNet		Search	Q Sign in	
Inside Microsoft Research	Bloa Microsoft Research website Researc	h news center		
				\
Guidelines for U	Jsing the CryptDB	System Securely		
	C 71	5		
Raluca Ada Popa UC Berkelev	Nickolai Zeldovich MIT CSAIL	Hari Balakrishnan MIT CSAIL		
1 Introduction				
This report has two goals. First, we re	eview guidelines for using the	CryptDB system [PRZB11, Pop14	4] se-	
curely by the administrators of databas and elaborated on in [Pop14] but in 1	se applications. These guideling ight of some recent work	nes were already described in [PRZ	CB11] es in-	
correctly, a short document devoted to	summarizing these guideline	s may be useful.		
Second, we explain that the recent usage of CryptDB, in which the auth	study of Naveed, Kamara, and nors violate CryptDB's secur	d Wright [NKW15] represents an u ity guidelines. Hence, the conclu	nsafe Isions	
drawn in that paper regarding CryptDI	B are both unfounded and inco	prrect: had the guidelines been follo	owed,	
none of the cranned attacks would hav	e been possible.		/	
Same Kamara a rasearchar in Mi	endly contraction and a motion with the second second		1000212 2011A TO 11521 I	uu-yeai

none of the claimed attacks would have been possib

### Combatting Frequency Analysis

- We want to smooth out frequency distribution so that frequency analysis becomes ineffective.
  - Performing worse than random guessing of plaintext.
- We also want to preserve ability to efficiently perform search queries on a standard database.
  - Rules out fully randomised/IND-CPA secure encryption.
- What about adding a limited amount of randomness?
- Leads to idea of applying homophonic encoding to produce Frequency Smoothing Encryption (FSE) schemes (Lacharité-Paterson, forthcoming).

# Frequency Smoothing Encryption – Combatting Frequency Analysis



- Homophonic Encoding (HE) consumes **small** amount of randomness.
- Make number of encodings proportional to frequency of *p* for good frequency smoothing.
- DE = Deterministic Encryption.
- Match on  $\{c_1, c_2, c_3, c_4\}$  instead of a single ciphertext.
- Query complexity blow-up by max. number of encodings in worst case.

#### Interval-based Homophonic Encoding (IBHE)

- Encoding space = *r*-bit strings / interval [0, 2<sup>*r*</sup>).
- Represent encodings of p having frequency f by an interval of size approximately f x 2<sup>r</sup>.
- Select uniformly at random from interval to encode *p*.
- Needs an encoding table to store an interval for each plaintext item; |p| x 2r bits.
- Also needs a decoding table mapping bits back to plaintexts.

$$p_{o}$$
  $p_{1}$   $p_{2}$   $p_{3}$  ....

- Can prove that as r goes to ∞, no distinguisher can tell apart ciphertexts from uniformly random strings.
- But even for moderate r, IBHE + DE smooths well for all but very skewed data.
- Rapidly limits (generalised) frequency analysis to being worse than a pure guessing attack.
  - Such an attack is always possible for limited domain of plaintexts.
- We used same evaluation framework as Naveed-Kamara-Wright (CCS15).
  - Except that we gave the adversary the **exact, per-hospital distribution** as the auxiliary distribution!





#### Length of stay

- **Warning**: FSE only protects against a basic snapshot attacker.
- Recent work of Grubbs-Ristenpart-Shmatikov (HotOS17) questions legitimacy of snapshot attack model.
- Columns are treated in isolation.
- More powerful adversary could perform frequency analysis on the sets of **responses** to queries.
- Scheme does not protect against an active attacker who can inject his own queries.



# Analysis of OPE/ORE

#### Order Preserving/Revealing Encryption

- OPE: if x < y then Enc(x) < Enc(y).
- **ORE**: there exists a (public) efficiently computable function "Order" such that:

x < y iff Order(Enc(x), Enc(y)) = 1

- OPE/ORE allows range queries.
- Client who wishes to query on range [a,b] instead sends query for range [Enc(a), Enc(b)] to server.

### Order Preserving/Revealing Encryption

- Q: If DE leaks badly, does OPE/ORE leak even more?
- A: Often, yes.
- **Folklore**: if OPE scheme is deterministic and plaintext data is **dense** (every possible plaintext occurs) then a snapshot adversary can learn which plaintext is which.
- Simply order the ciphertexts and then read off the plaintexts.
- **Take-away**: beware of formal security models for OPE/ORE.
- This can sometimes be generalised to the non-dense case...

#### The Scheme of Chenette-Lewi-Weis-Wu (FSE16)

- CLWW (FSE16) presented a clever and practical ORE scheme built using only PRFs.
- CLWW gave a precise characterisation of leakage in a simulation-based security model:
  - Given two ciphertexts Enc(x) and Enc(y), the scheme leaks exactly the first index at which bits of x and y differ (and which is bigger).
- Example: given Enc(x = 1101<sub>2</sub>) and Enc(y = 1001<sub>2</sub>), the scheme would leak that the two plaintexts are equal in MSB (bit o) but that the first one has 1 in bit 1 and the other o in bit 1.
- Leakage is **greater** than in an ideal OPE scheme, which would leak only order.

#### An Attack on the CLWW Scheme

- In the dense case, the folklore analysis applies.
- What about the non-dense case?

- Assumption 1: snapshot attacker.
- Assumption 2: N plaintexts, close to uniformly random on the s MSBs, where N > s2<sup>s</sup>.

Then, with high probability, the attacker can learn the *s* most significant bits of every plaintext.

#### An Attack on the CLWW Scheme

- Assumption 1: snapshot attacker.
- Assumption 2: N plaintexts, close to uniformly random on the s MSBs, where N > s2<sup>s</sup>.

• Second assumption implies that, with high probability, every possible *s*-bit prefix occurs in at least one plaintext:

```
Prob \approx 1 - 2^{-S/2^{(S+1)}}.
```

- Use the CLWW scheme's leakage to order the *N* ciphertexts on the 2<sup>s</sup> distinct *s*-bit prefixes.
- Now read off the *s* most significant bits of each plaintext.

#### Implications of the Attack

- Suppose a company has 10,000 employees with salaries that are 20-bit numbers (between \$0 and \$2<sup>20</sup>-1).
- We can set *s* = 10 (10 x 2<sup>10</sup> ≈ 10,000).
- Attack yields 10 MSBs of *every* salary.
- This is enough to identify each salary up to accuracy of \$1k.

- Example generalises to, say, 32-bit salaries that are all zero in the first 12 bit positions.
- Sufficient that data be dense in some positions (and constant in leading positions).

#### Further Research on OPE/ORE Leakage

Several attack recent papers examine the real-world implications of the leakage of OPE/ORE schemes for snapshot attackers:

- **Durak-DuBuisson-Cash (CCS16):** attacks on correlated columns of OPE/ORE-encrypted data, especially longitude/ latitude data.
- **Grubbs-Sekniqi-Bindschaedler-Naveed-Ristenpart (S&P17):** revisit Naveed-Kamara-Wright for OPE/ORE; recast ptxt/ctxt matching problem as min-weight, non-crossing bipartite matching problem, solve it efficiently for many types of data, relies on auxiliary distributions.



# Access Pattern Leakage for Range Queries

#### Analysis of Access Pattern Leakage for Range Queries

**Kellaris-Kollios-Nissim-O'Neill (CCS16):** analysis of access pattern leakage for SE; applicable to OPE/ORE schemes too.

- Honest-but-curious attack setting, stronger than snapshot adversary.
- **Assumption**: adversary can see which database rows are returned in response to any range query.
- For *N*-valued database, complete reconstruction in O(*N*<sup>4</sup>) queries.
- For dense case: O(*N*<sup>2</sup>log*N*) queries suffice.

### Analysis of Access Pattern Leakage for Range Queries

- Adversary in KKNO (CCS16) does not need to directly see the actual ranges queried.
  - In OPE/ORE, adversary would see only ciphertexts Enc(x), Enc(y) corresponding to range endpoints.
  - But in OPE/ORE setting, and in some SE schemes\*, the **rank** also leaks.
  - The rank of a ciphertext is its position in an ordered list of all the ciphertexts.

\*e.g. Arx scheme of Poddar-Boelter-Popa and FH-OPE scheme of Kerschbaum.



#### Exploiting Rank in Analysis of Access Pattern Leakge

- Can we use the rank leakage to improve attack complexity?
- Lacharité-Minaud-Paterson (forthcoming):
  - Yes!
  - And much more besides...

## Exploiting Rank in Analysis of Access Pattern Leakge

Simple motivating example: consider range queries [*a*,*b*] in which *a* is uniformly random.

- Then with probability 1 1/N, after  $2N\log N$  queries, all N possible values for a will have arisen.
  - Follows from standard analysis of the coupon collector problem.
- Easy to identify and order different *a* values based on rank leakage.
- All values of *a* in queries are now known.
- Pick out *N* queries with distinct values of *a*; each such query produces a set of responses *Y*<sub>*a*</sub> (records in database).
- Then the set of records with value *a* is:  $Y_a U_{i>a} Y_i$ .

#### Exploiting Rank in Analysis of Access Pattern Leakge



43

#### Improved Analysis of Access Pattern Leakge

- Our simple example appears to show that rank leakage helps the adversary.
- In fact, we can dispense with rank leakage and obtain an NlogN+O(N) attack in the general "dense" case!
  - Improving on KKNO's O(N<sup>2</sup>logN) attack.

- We also consider the problem of **approximate reconstruction**.
- We can efficiently reconstruct values in records up to an absolute error of  $\varepsilon N$  after seeing only O(N) queries!
  - With a constant of  $2\log(1/\epsilon)$ .



#### Improved Analysis of Access Pattern Leakge

Finally, we study algorithms for **approximate reconstruction** with the assistance of rank and an auxiliary distribution.

- Significant reduction in number of queries required for accurate reconstruction.
- Perform set intersections and then map back to underlying data using rank + auxiliary distribution.
- Experiments with aggregated HCUP data...

#### Approximate Reconstruction with Auxiliary Distribution

Absolute errors of point guesses



46



# Concluding Remarks

### **Concluding Remarks**

- Use DE/OPE/ORE with extreme care if at all.
- We are currently in a propose/break/patch cycle.
  - Despite the provision of security models and proofs.
- Just identifying and proving leakage is not enough; we need to also identify real-world implications of that leakage.
- Does PPE provide added security or a false sense of security?

