

# Have You Been Pwned? the Aftermath of Data Leakage

AMRIT KUMAR & CÉDRIC LAURADOUX

Privatics team  
INRIA Rhône-Alpes

27th February 2014

# Recent Leakages

Victim	Data	Amount (million)
Adobe	username, email, encrypted pwd, pwd hint	150
LinkedIn	hashed password	6.5
Snapchat	user name, telephone number, city	4.6
Last.fm	hashed password	2.5
eHarmony	hashed password	1.5
Apple	UDID, device name, device type, etc.	1
coming soon	...	

How may a user know if his personal data has leaked?

# Alerting Websites : Good

If data leaked, user needn't download huge leaked database.

Example : [lastpass.com/adobe](http://lastpass.com/adobe)

Was My Adobe Account Hacked?

If you would like to find out if your Adobe account was one of the 150 million that were leaked, you can use the below tool:

Enter your email address here:

Your Adobe account was one of the ones that was compromised.

Your **email address** and encrypted Adobe **password** were found in the list of stolen Adobe accounts.

Did you know that 250 other people used the same password as you did for their Adobe account? Hackers have their password hints and can use them to guess your password too!

We have not sent another email to **toto@gmail.com** as there have already been several prior emails sent with instructions on how to obtain your Adobe password hint.

**We strongly urge you to follow our recommendations and Immediately change your Adobe and related passwords!!**

# Alerting Websites : Phishing

If data isn't leaked, website learns a new data (**possibly your password**).

Example. [didigetgawkered.com](http://didigetgawkered.com)

Enter your username or email address to see if you were affected by the [Gawker hack](#):

roger

Did I Get Gawkered?

Powered by Duo Security



Sharing is caring! Pass this along to a friend:

Powered by Duo Security



# Alerting Websites Cont.

Websites	Type	Database(s)	https	Statement	Answer	Descrip.
<a href="http://adobe.cynic.al">adobe.cynic.al</a>	S	ADOBE	✗	✗	✓	✓
<a href="http://bit.ly/1by3hd9">bit.ly/1by3hd9</a>	S		✓	✓	✓	✗
<a href="http://lucbie.com/credgrep">lucbie.com/credgrep</a>	S		✗	✗	✗	✗
<a href="http://adobe.breach.il.ly">adobe.breach.il.ly</a>	S		✗	✗	✗	✗
<a href="http://snapcheck.org">snapcheck.org</a>	S	SNAPCHAT	✗	✗	✓	✗
<a href="http://findmysnap.com">findmysnap.com</a>	S		✗	✗	✓	✗
<a href="http://lookup.gibsonsec.org">lookup.gibsonsec.org</a>	S		✗	✗	✓	✗
<a href="http://didigetgawkered.com">didigetgawkered.com</a>	S	GAWKER	✗	✗	✗	✗
<a href="http://lastpass.com">lastpass.com</a>	A	6	✓	✓	✓	✗
<a href="http://haveibeenpwned.com">haveibeenpwned.com</a>	A	9	✓	✓	✓	✓
<a href="http://bit.ly/1fj0SqV">bit.ly/1fj0SqV</a>	A	12	✓	✗	✗	✗
<a href="http://dazzlepod.com/disclosure">dazzlepod.com/disclosure</a>	A	28	✗	✓	✓	✗
<a href="http://bit.ly/1aJubEh">bit.ly/1aJubEh</a>	A	3346/194	✓	✗	✓	✗
<a href="http://hacknotifier.com">hacknotifier.com</a>	H	Unknown	✗	✓	✓	✗
<a href="http://pwnedlist.com/query">pwnedlist.com/query</a>	H	Unknown	✓	✗/✓	✓	✓
<a href="http://sicherheitstest.bsi.de">sicherheitstest.bsi.de</a>	H	Botnets	✓	✓	✓	✗

# Problem

## Scenario

- Server hosts a public (leaked) database.
- Client wishes to know if his data belongs to the database.

## Problem Statement

Can the client

- **privately** query a database,
- and remain **more efficient** than asking for the database itself?

# Private Information Retrieval

“Private Information Retrieval”, by Chor et al, FOCS'95.

## Scenario

- Server has a database  $\mathcal{DB} = \{x_1, \dots, x_n\} \in \{0, 1\}^n$ .
- Client has an index  $1 \leq i \leq n$  and wishes to retrieve  $x_i$ .

## Privacy

- Information-theoretic privacy : no information on  $i$  should be revealed. Chor et al. '98, Devet et al. '12.
- Computational privacy :  $i$  cannot be retrieved by a PPT server. Kushilevitz and Ostrovsky '97, Cachin et al. '99, Chang '04, Lipmaa '05, Gentry and Ramzan '05.

Trivial PIR : Send the database, communication complexity  $n$ .

## PIR Variants

- Single server/multiple non-communicating servers.
- Private block retrieval (PBR) :

$$DB = \{x_1, \dots, x_n\} \in \{0, 1\}^{\ell \cdot n}, 1 \leq i \leq n.$$

- PIR by Keywords :

$$DB = \{x_1, \dots, x_n\} \in \{0, 1\}^{\ell \cdot n}, kw \in \{0, 1\}^{\ell}.$$

“Private Information Retrieval by Keywords ”, Chor et al. 1997

- ▶ Store database in hash table.
- ▶ Invoke PBR on digest.



# PIR in Our Context ?

## Is PIR applicable ?

- Original PIR :  
**NO**, user doesn't have physical address.
- PIR by keywords :  
**NO**, keywords statically defined  $\Rightarrow$  prevents non-membership query.
- **IDEA** : Use **Bloom filter** to represent data  
Allows to apply original PIR + low memory footprint.

# Bloom Filters

“Space/Time Trade-offs in Hash Coding with Allowable Errors”, Burton Bloom, Communications of the ACM, 1970.

## Tools

- $S$  : a set of  $n$  elements.
- $k$  hash functions with range  $0, 1, \dots, m - 1$ .  
hash functions  $h_1, \dots, h_k$  can be truncated SHA-1.
- $b$  : array of bits of size  $m$ , initialized to 0.

## Working Mechanism

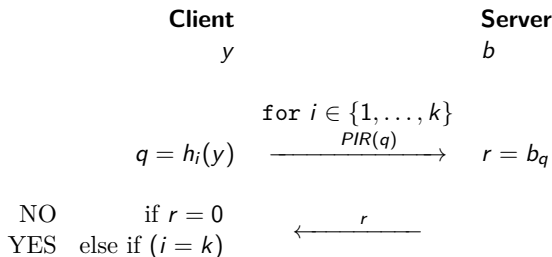
- Insertion : find  $k$  hash-indices and set these bits to 1.
- Query : find  $k$  hash-indices and check if all these bits are set to 1.

## False positive

Truncation  $\Rightarrow$  Collisions  $\Rightarrow$  False positive  $f$ .

## PIR with Bloom Filter

- Server builds the Bloom filter  $b$  using  $k$  hash functions  $\{h_1, h_2, \dots, h_k\}$
- Client for a data  $y$  generates  $\{h_1(y), h_2(y), \dots, h_k(y)\}$ .
- For each  $1 \leq i \leq k$ , Client invokes a single-server PIR and retrieves  $b[h_i(y)]$ .
- If  $b[h_i(y)] = 0$  for any  $i$ , then  $y$  is not in the database.



# Possible Solutions

- $k$  invocations of single server PIR to Bloom filter.  
⇒ computational privacy + communication complexity sublinear/square logarithmic.
- Send the complete filter.  
⇒ complete privacy + low communication complexity than trivial PIR.

Which is the most time efficient ?

# Performance Evaluation

Building Bloom filter ( $f = 2^{-128}$ ), SHA-1.

Database	Size	$n$	$m$ (MB)	Build time (mins)	# CPU
SNAPCHAT	49 MB	4609621	102	6	1
LINKEDIN	259 MB	6458019	142	10	1
ADOBE	3.3 GB	153004872	412	198	1
ADOBE	3.3 GB	153004872	412	50	4

$$f = 2^{-64}$$

Database	$m$ (MB)	Build time (mins)
SNAPCHAT	52	2
LINKEDIN	72	3
ADOBE	206	72

$$f = 2^{-32}$$

Database	$m$ (MB)	Build time (mins)
SNAPCHAT	26	1
LINKEDIN	36	1.5
ADOBE	102	30

# Conclusion

- Data leakages, the panorama.
- Privately alerting a user : PIR with Bloom filter vs. filter transmission.
- Compression techniques on Bloom filter.
- Experimental evaluation of techniques.

# Thank you !