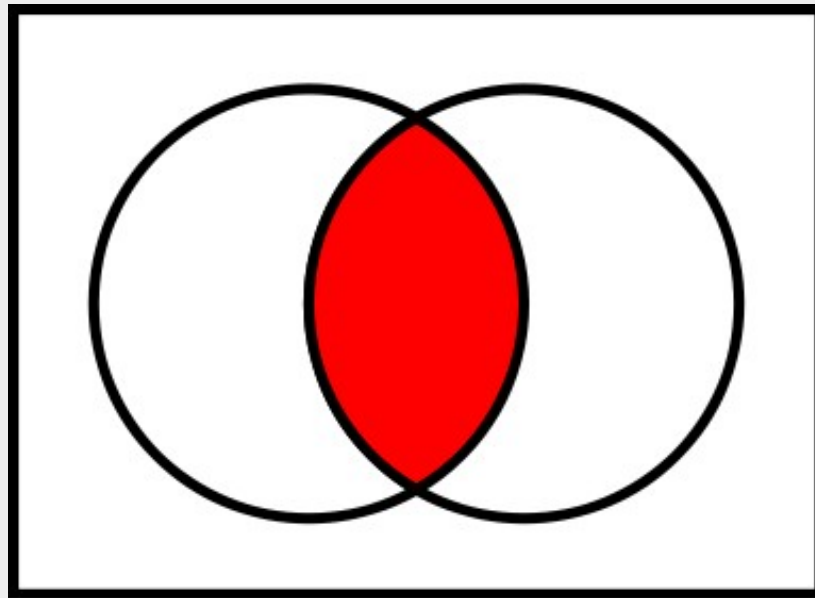# Private Set Intersection

A cryptographic technique that allows two parties to compute the intersection of their sets without revealing anything except the intersection

Matej Kovačič
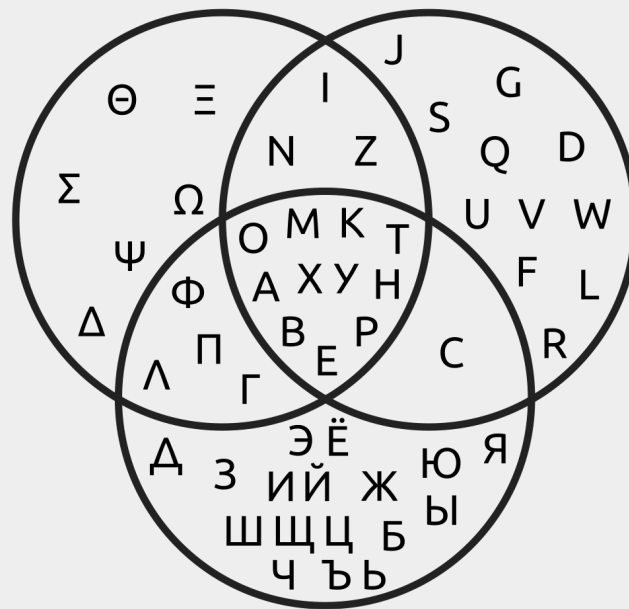
Matjaž Rihtar

# Private Set Intersection

In a private set intersection (PSI) protocol two parties jointly compute the intersection of their private input sets.



A client and a server jointly compute the intersection of their private input sets in a manner that at the end the client learns the intersection and the server learns nothing (one-way PSI) or both learn the intersection (mutual PSI).

# Private Set Intersection

Such protocols are especially useful whenever one or both parties (who do not fully trust each other) must compute an intersection of their respective data sets (i.e. only the required minimum amount of information should be disclosed).



Intersections of the Greek, English and Russian alphabet, considering only the shapes of the letters and ignoring their pronunciation.

# Examples of use

1. A government agency needs to make sure that employees of its industrial contractor have no criminal records. Government does not want to disclose list of all convicted felons, and contractor does not want to disclose their full list of employees. However, both would like to know the intersection, if any.

2. Federal tax authority wants to learn whether any suspected tax evaders have accounts with a certain foreign bank and, if so, obtain their account records. The bank is forbidden of wholesale disclosure of their account holders and the tax authority can not reveal its list of suspects to a foreign bank.

Source: Emiliano De Cristofaro, Gene Tsudik. 2009. Practical Private Set Intersection Protocols with Linear Computational and Bandwidth Complexity. Cryptology ePrint Archive, Report 2009/491. <https://eprint.iacr.org/2009/491>.

# Examples of use

3. Two real estate companies would like to identify homeowners who are double-dealing. That means they have signed *exclusive* contracts *with both companies* to assist them in selling their properties. None of them would like to reveal the other party list of their customers.

4. Two national law enforcement bodies want to compare their respective databases of terrorist suspects. National privacy laws prevent them from revealing bulk data, however, by treaty, they are allowed to share information on suspects of common interest.

Source: Emiliano De Cristofaro, Gene Tsudik. 2009. Practical Private Set Intersection Protocols with Linear Computational and Bandwidth Complexity. Cryptology ePrint Archive, Report 2009/491. <https://eprint.iacr.org/2009/491>.

# Why not hashes?

Hashes are supposed to be one-way functions, but this is not always true...

**Our target hash:**
**b2ed3039dd45146d2a1f9f420f8816c5**

First we generate a database of all possible phone numbers...

Finally we look for our target hash in our hash table and so reveal the original phone number.

**Hash list:**
```
57204d13a7d353ab0cd3a0d212bbc425
5cf06d5bf80c49e5156515d932536a36
2b6895ae6a902d00da9e04a4d4269a68
d74e8d6aacdc1747ef626982bd5f11bb
06c9b2b57717d80f08e1773394b6f502
b2ed3039dd45146d2a1f9f420f8816c5
b4d06450af671e0edc403b257315107a
cade9c741dd959d966420a6659238272
2f3498ebf828e6b996a1d29a7b2d877b
bd3ca378488e00055d5b23df1252e443
5bdbd627461fc598bb6a00d0168587fe
90e3415b9b87576ca111bf52e1d07265
```

Then we compute table of hashes for all that phone numbers...

# Why not hashes?



On the Internet we can find lists of precomputed hashes...

# How does PSI work?

| | |
|---|---|
| $a \leftarrow A$ | variable $a$ is chosen uniformly at random from set $A$ |
| $\tau$ | security parameter |
| $n, e, d$ | RSA modulus, public and private exponents |
| $g$ | group generator; exact group depends on context |
| $p, q$ | large primes, where $q = k(p-1)$ for some integer $k$ |
| $H()$ | cryptographic hash function, codomain depends on context |
| $H'()$ | cryptographic hash function $H' : \{0,1\}^* \rightarrow \{0,1\}^\tau$ |
| $\mathcal{C}, \mathcal{S}$ | client's and server's sets, respectively |
| $v, w$ | sizes of $\mathcal{C}$ and $\mathcal{S}$, respectively |
| $i \in [1, v], \; j \in [1, w]$ | indices of elements of $\mathcal{C}$ and $\mathcal{S}$, respectively |
| $c_i, s_j$ | $i$-th and $j$-th elements of $\mathcal{C}$ and $\mathcal{S}$, respectively |
| $hc_i, hs_j$ | $H(c_i)$ and $H(s_j)$, respectively |
| $R_{c:i}, R_{s:j}$ | $i$-th and $j$-th random value generated by client and server, respectively |

Table 1: Notation

- Common input: $n, e, H(), H'()$
- $H()$ is a Full-Domain Hash $H : \{0,1\}^* \rightarrow \mathbb{Z}_n^*$
- Client's input: $\mathcal{C} = \{hc_1, \cdots, hc_v\}$, where: $hc_i = H(c_i)$
- Server's input: $d, \mathcal{S} = \{hs_1, \cdots, hs_w\}$, where: $hs_j = H(s_j)$

**OFF-LINE:**

1. Server:
   $$\forall j, \text{compute: } K_{s:j} = (hs_j)^d \bmod n \text{ and } t_j = H'(K_{s:j})$$

2. Client:
   $$\forall i, \text{compute: } R_{c:i} \leftarrow \mathbb{Z}_n^* \text{ and } y_i = hc_i \cdot (R_{c:i})^e \bmod n$$

**ON-LINE:**

3. Client ⟶ Server: $\{y_1, .., y_v\}$

4. Server:
   $$\forall i, \text{compute: } y_i' = (y_i)^d \bmod n$$

5. Server ⟶ Client: $\{y_1', ..., y_v'\}, \{t_1, .., t_w\}$

6. Client:
   $$\forall i, \text{compute: } K_{c:i} = y_i'/R_{c:i} \text{ and } t_i' = H'(K_{c:i})$$
   OUTPUT: $\{t_1', .., t_v'\} \cap \{t_1, .., t_w\}$

Figure 4: Blind RSA-based PSI Protocol with linear complexity

PSI exchange, which is a functionality in secure multi-party computation (also known as privacy-preserving computation) **is much more secure than exchange with hashes**, because it uses encryption.

Source: Emiliano De Cristofaro, Gene Tsudik. 2009. Practical Private Set Intersection Protocols with Linear Computational and Bandwidth Complexity. Cryptology ePrint Archive, Report 2009/491. <https://eprint.iacr.org/2009/491>.

# How does PSI work?

**Server** takes its data, computes hash (SHA-256), convert this to integer, makes RSA decode (mathematical operation which is a part of RSA decode procedure), and computes new hash (SHA-256).

$$data \rightarrow hash(s\_data) \rightarrow s\_number1 \rightarrow (s\_number1)^d \bmod n \rightarrow s\_number2 \rightarrow hash(s\_number2) \rightarrow s\_number3$$

**Client** computes hashes (SHA-256) for all data.

$$c\_data \rightarrow hash(c\_data) \rightarrow c\_number$$

Then client for each data element generates random number (c_rand).

Now client has: c_data :: c_number :: c_rand

# How does PSI work?

Then client multiplies each data element (c_data) with RSA encoded random number. Random number c_rand is RSA encoded with server's public key. {RSA(c_rand) * c_data mod n}. → This result is then sent to server.

For each received data server then computes RSA decode. Decoded client data and server's final hash of its data (s_number3) are sent to client. ←

Finally client takes (his) RSA decoded client data from server and divides them with c_rand number and computes a hash. This hash is then compared with server's data (s_number3).

If both hashes are the same, the data are the same, so there is intersection between the data on server and a client.

# Transliteration

Another problem is that different languages has different writing system.

Transliteration is utilized when a word or phrase must be "transmitted" to a language with a different writing system.

For the we took ICAO (The International Civil Aviation Organization) transliteration tables, however we improved them a little, so our system now uses "extended ICAO" transliteration standard.

Example (lines containing # are ICAO original standard, lines containing ## are our "extended ICAO"):
- 0x00DD: 'Y',   # Y with acute
- 0x00DE: 'TH',  ## Thorn (Iceland)
- 0x00DF: 'ss',  # sharp s (Germany)
- 0x0136: 'K',   # K with cedilla
- 0x0137: 'k',   ## k with cedilla
- 0x0138: 'k',   ## kra
- 0x0139: 'L',   # L with acute

# Transliteration

Our "extended ICAO" supports all small and some additional upper case and additional cyrillic characters.

Our procedure is that we first perform transliteration and then convert characters to uppercase.

Example of a problem with ICAO standard (one character could be transliterated to differend latin characters – depending of an original language):

- 0x0413: 'G',   # GHE (or G, except Belorussian and Serbian H)

Suggestions:

- If you are using small caps, you can leave them – our system will automatically transcode them to uppercase.

- If you have cyrillic characters in original, leave them in cyrillic.

# Transliteration

Example of transliteration:

Matej Kovačič → MATEJ KOVACIC

Matjaž Rihtar → MATJAZ RIHTAR

# Permutations

Asumption is that one person can have from 2 to 5 names (first name + up to 4 additional names or surnames).

This means from 2 to 394 possible permutations.

Example: "Kyle Reese Sergeant, born 1. 1. 2003".

All possible permutations:



SERGEANT REESE KYLE 2003-01-01
REESE SERGEANT KYLE 2003-01-01
SERGEANT REESE 2003-01-01
REESE SERGEANT 2003-01-01
REESE KYLE SERGEANT 2003-01-01
SERGEANT KYLE 2003-01-01
KYLE SERGEANT 2003-01-01
SERGEANT KYLE REESE 2003-01-01
KYLE REESE 2003-01-01
KYLE REESE SERGEANT 2003-01-01
REESE KYLE 2003-01-01
KYLE SERGEANT REESE 2003-01-01

# Permutations

We can find even cases when one or more names are similar and not exactly the same. For example:

- country A has a person "SERGEANT REESE **KYLE** 2003-01-01"

- country B has a person "SERGEANT REESE **KYLLE** 2003-01-01"

PSI exchange will find a match, because one of the permutations in both countries is "SERGEANT REESE 2003-01-01" (name KYLE or KYLLE is neglected in this permutation).

However, please note, that birthday of both persons should always be the same. Otherwise, there is no match.

# Technical requirements

PSI application is written in programming language Python.

Installation script provides all installation packages (Python 2.7, required Python modules and the application itself).

Application runs on operating system Windows 10 (minimum Windows 7), but it could be ported to other platforms as well.

Other needed software:

- for importing the data from Excell files, Microsoft Office 2010 or more recent version is required;

- web browser (suggested is the latest version of Chrome or Firefox).

# Technical requirements

Data could be imported via CSV or Excell (XLSX) file with predefined set of fields.

# Hardware setup

PSI exchange is taking place in closed local network. It could also be implemented in VPN network.

Client 1 (country A)

Dispatch server
(Central PSI Monitor)

Dispatch server monitors state of clients and starts PSI procedure

router
(closed local network)

Client 2 (country B)

Clients contain data. PSI exchange is done directly between clients

Client 3 (country C)

# Data flow



**Client 1 (country A):**
- user selects country;
- user uploads data (XLSX/CSV);
- client transliterates data;
- when user confirms (s)he is ready, client prepares itself for PSI excange and reports to dispatch server;
- client is now waiting for dispatch server to start PSI process;
- when PSI process is done, client reports it to dispatch server.

**Data exchanging through PSI procedure is encrypted.**

1. Client reports it is ready for PSI process

2. Central PSI Monitor instructs client to start PSI process

3. Client performs PSI exchange with all other clients

**Dispatch server (Central PSI Monitor):**
- records **status** of clients (which client is assigned to which country, is transliteration on client done, is client ready for PSI exchange);
- can **start PSI** process **between clients** and can **monitor** PSI exchange progress;
- **does not see client data**!

PSI exchange is taking place **only among clients**. Central PSI Monitor (dispatch server) **is not involved** in PSI exchange.

Client 2 (country B)

Client 3 (country C)

# PSI application

Application consists of a three parts. First part is Central PSI Monitor, which monitors statuses of clients and starts PSI procedure.

Client application is divided into two parts. One is used to exchange the data about the persons and the other to exchange so called electronic identificator data (phone numbers, licence plates, e-mail addreses social network ID's,...).



Each country can monitor its own status.

Central PSI Monitor at dispatch server is monitoring country's status and PSI progress.

PSI exchange is taking place **directly between two countries** (peer-to-peer).

# PSI application



Step 1: Select your country.

# PSI application



Step 2: Select your file with data. You can upload Excel or CSV file.

# PSI application



Step 3: File is uploaded to the application. Transliteration is done and permutations are computed.

# PSI application



Step 4: You can rewiev the data.

# PSI application



Step 5: Click to "Main page" and then "Start PSI monitor".

# PSI application



Step 6: Now application generates RSA keys. After that, you can click to "Connect". Your client will be connected to Central PSI Monitor. Then you can click "Ready" button. This will inform the Central PSI Monitor that your country is ready for PSI exchange (this will be indicated by message: {"STATUS": "OK, country ready"}).
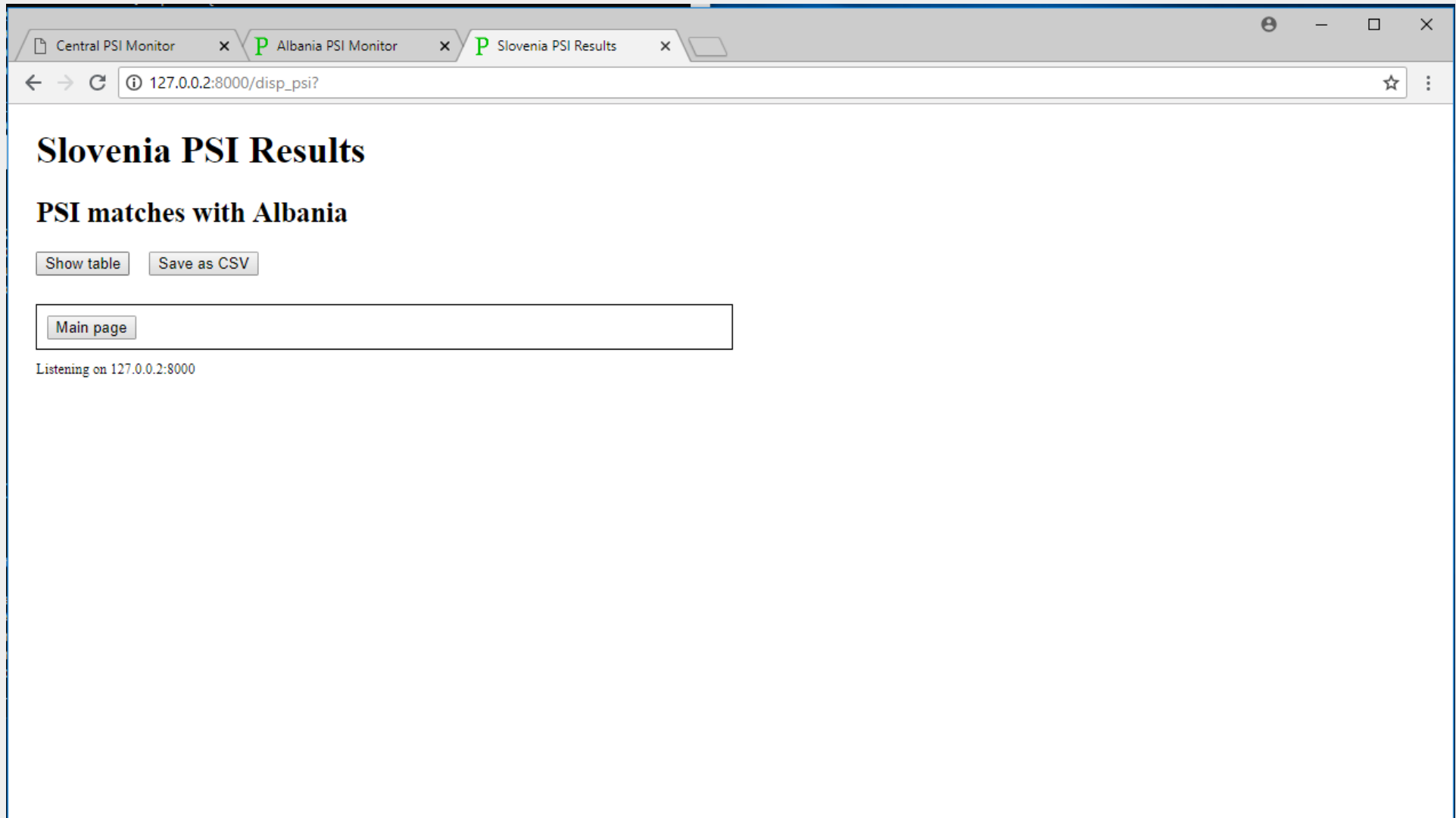
# PSI application



Step 7: When all clients are ready, operator of the Central PSI Monitor can start PSI exchange among countries. Central PSI Monitor will then indicate which countries are doing PSI exchange and when the procedure is finished.

# PSI application



Step 8: Now you can return to main screen and check PSI results.

# PSI application



Step 9: You can see with which clients (countries) you have common results.

# PSI application



Step 10: You can review the results (common records in both databases) and save the results in CSV file.

# Discussion



matej.kovacic@telefoncek.si
matjaz@eunet.si