

# Efficient Similarity Search over Encrypted Data

then [Enc<sub>id</sub>(B<sub>k</sub>), Enc<sub>payload</sub>(V<sub>Bk</sub>)  $\in$  I.

## Introduction

- Sensitive information is increasingly outsourced into the cloud in encrypted form.
- Encrypted storage complicates the computation on the data such as similarity search.
- Available searchable encryption schemes have been designed for only exact query matching.
- Secure multi-party computation based sophisticated cryptographic techniques enable similarity search but they are not practical.

## Motivation

- Efficient similarity search have been extensively studied for plain data and one of the best solution is locality sensitive hashing (LSH).
- Appealing properties of LSH can be used in the context of encrypted data via a secure index.
- Secure index should provide some guarantees to ensure confidentiality of the sensitive data.

## **Security Definition**

- Access Pattern (Ap): Identifiers of data items that are in the result set of a specific query.
- Similarity Pattern (Sp) : Relative similarity between distinct queries.
- History (H<sub>n</sub>:) Data collection and query set for n consecutive queries.
- Trace(Y): Maximum amount of information that data owner allows the leak to an adversary (encrypted data collection (C), Ap and Sp).
- View(V) : Information that is accessible to the adversary (C, query trapdoors).

 $Pr[Dist(V(H_n)) = 1] - Pr[Dist(S(\gamma(H_n)) = 1] < \frac{1}{poly(s)}$ 

Mehmet Kuzu, Mohammad Islam, Murat Kantarcioglu Department of Computer Science, University of Texas at Dallas

## **Secure Search Scheme** Similarity Searchable Encryption Keygen(s): Output a private key K according to Secure LSH Index Alice **Shared Secret** security parameter s. Data Identifiers • Enc(K, $D_i$ ): Encrypt $D_i$ with key K to form $C_i$ . collection encryption • $Dec(K, D_i)$ Decrypt $C_i$ with key K to obtain $D_i$ . index construction BuildIndex(K, D) Extract feature set for each function $D_i$ ( $F_i$ ) in the data collection and form secure index I with the extracted features. Trapdoor Construction: Trapdoor for any Trapdoor(K, f) Generate a trapdoor for a feature $f_i(T_{f_i})$ is formed as follows: specific feature f with key K and output T. $T_{f_i} = \{Enc_{id}(g_1(\rho(f_i))), ..., Enc_{id}(g_\lambda(\rho(f_i)))\}$ FuzzySearch(I, T) Perform search on I with Search: Retrieve encrypted bucket contents trapdoor of feature f (T) and output encrypted collection C with high probability: decrypt bit vectors and rank items: $C_j \in C$ if $\exists (f_i \in F_j) [dist(f_i, f) \le \alpha]$ $V_{B_k} = Dec_{K_{payload}}(\sigma_{V_{B_k}}), \ \pi_{B_k} \in T_{f_i} \&\& [\pi_{B_k}, \sigma_{V_{B_k}}] \in I$ $C_j \notin C \quad if \quad \forall (f_i \in F_j) \; [dist(f_i, f) \ge \beta]$ $score(id(D_j)) = |Z_j|$ s.t. $V_{B_k} \in Z_j \ iff \ V_{B_k}[id(D_j)] = 1$ **Secure LSH Index** buckets between compared features. Feature Extraction : Extract features that characterizes sensitive data item D<sub>i</sub>. Metric Space Translation: Map features to vectors in a metric space to apply LSH. **Error Aware Keyword Search** Bucket Index Construction: Build a locality sensitive function $g : (g_1, g_2, ..., g_\lambda)$ from a family and apply g on feature vectors to form LSH buckets. Content of any bucket B<sub>k</sub> Elegant locality sensitive family has been is a bit vector $(V_{B_{\mu}})$ : is $[r_1, r_2, 1-r_1, 1-r_2]$ sensitive. $V_{B_k}[id(D_z)] = 1$ if $g_i(\vec{f_j}) = B_k$ for $g_i \in g, f_j \in D_z$ $[\mathbf{r}_1, \mathbf{r}_2, \mathbf{p}_1, \mathbf{p}_2]$ sensitive family can be formed $V_{B_k}[id(D_z)] = 0$ otherwisefrom MinHash via AND-OR construction: Bucket Index Encryption: Let Enc<sub>id</sub> be a $p_1 = 1 - (1 - (1 - r_1)^k)^{\lambda}, \quad p_2 = 1 - (1 - (1 - r_2)^k)^{\lambda}$ pseudorandom permutation, Enc<sub>pavload</sub> be a secure encryption scheme, $B_k$ be a bucket with content $V_{B_{\nu}}$ and I be the secure index,

This work has been published in IEEE ICDE 2012 Conference Proceedings.



- K<sub>coll</sub>: Secret key of data
- K<sub>id</sub>, K<sub>pavload</sub>: Secret keys of
- p: metric space translation
- g: locality sensitive function

corresponding to components of the trapdoor,

Ranking depends on the number of common

Data Decryption: Retrieve encrypted data items with highest scores and decrypt them with K<sub>coll</sub>.

Bloom filter encoding enables efficient space translation for approximate string matching.

designed for Jaccard distance (MinHash) that

 $A = \rho(john) = \{1, 2, 4, 5, 7, 10\}$ B = ρ(johhn) = {1, 2, 4, 5, 6, 7, 10}  $J_{d}(A,B) = 1 - 6/7 = 0.14$ 

## **Experimental Evaluation**



# Conclusion

A sample corpus of 5000 emails is constructed from publicly available Enron e-mail dataset.

Words in e-mails are embedded into 500 bit Bloom filters with 15 hash functions.

Common typos are introduced into the queries. (0.45, 0.8, 0.85, 0.01)-sensitive family is formed from MinHash to tolerate typos. LSH buckets are encrypted with AES in CTR mode.

We proposed LSH based secure index and search scheme to enable fast similarity search over encrypted data.

We provided a rigorous security definition and proved the security of the scheme to ensure confidentiality.

Efficiency of the proposed scheme is verified with empirical analysis.