



Open access Journal

International Journal of Emerging Trends in Science and Technology

Impact Factor: 2.838

DOI: <http://dx.doi.org/10.18535/ijetst/v3i03.01>

Privacy Preserving Similarity Based Text Retrieval through Blind Storage in Cloud

Authors

V. Thilagavathi*, V.R. Kavitha**

*Computer science and Engineering Prathyusha Engineering College, Chennai

Email: thilaga624@gmail.com

ABSTRACT

In portable distributed computing, a major application is to outsource the versatile information to outside cloud servers for adaptable information stockpiling. To test this issue, in this paper, We build up the Searchable encryption for multi-watchword positioned look over the capacity information. In particular, by considering the extensive number of outsourced the reports in the cloud, we use the significance score and k closest neighbor in methods to build up an effective multi-catchphrase look conspire that can give back the positioned indexed lists in view of the exactness. Security analysis demonstrates that our scheme can achieve confidentiality of documents and index, trapdoor privacy, access unlinkability, and obscuring access pattern of the search user. Finally, using extensive simulations. We leverage an efficient index to further improve the Search efficiency, and adopt the blind storage system to conceal access pattern of the search user.

Keywords used: Searchable encryption, blind storage, access pattern

INTRODUCTION

Portable distributed computing disposes of the equipment impediment of cell phones by investigating the Open and virtualized distributed storage and registering assets, and hence can give considerably more legitimate and versatile portable administrations to clients. The subcontracted data typically contain sensitive privacy information, such as personal photos, emails etc., which would lead to normal confidentiality and privacy violations, if lacking efficient protections. Look over encoded information ought to bolster the accompanying three capacities. To begin with, the searchable encryption associations ought to bolster multi-watchword look, and give the comparable client experience as seeking in Google seek with various catchphrases; single-watchword pursuit is a long way from adequate by just returning exceptionally restricted and erroneous indexed lists. Second, to rapidly distinguish most significant results, the pursuit client would normally lean toward cloud servers to sort the returned list items in an importance based request positioned by the pertinence of the inquiry solicitation to the records. Likewise, demonstrating the positioned pursuit to clients can likewise dismiss the superfluous Network movement by just exchange back the most applicable results from cloud to inquiry clients. What's more, demonstrating the positioned hunt to clients can likewise neglect the pointless Network movement by just exchange back the most important results from cloud to pursuit clients. Third, concerning the pursuit proficiency, since the quantity of the archives contained in a database could be uncommonly extensive, searchable encryption plans ought to be proficient to rapidly react to the hunt demands with minimum deferrals. We give careful security examination to accept that the EMRS can achieve a high

security level including secrecy of reports and file, trapdoor privacy, trapdoor unlink ability, and darkening access example of the pursuit client. The cloud server ought to be forestalled from prying into the outsourced reports and can't find any relationship between the records and catchphrases utilizing the file. To keep her quests from being presented to the cloud server, the cloud server ought to be kept from knowing the accurate catchphrases contained in the trapdoor of the seek client. Trapdoor Unlink ability: The trapdoors ought not to be linkable, which implies the trapdoors ought to be completely diverse regardless of the fact that they contain the same catchphrases. At the end of the day, the trapdoors ought to be randomized cloud server from realizing any extra data about the reports and the file, and to keep seek clients' trapdoors mystery, the EMRS ought to cover all the security prerequisites. We consider the Knowing Background model in the EMRS, which permits the cloud server to know more foundation data of the archives, for example, factual data of the watchwords. A visually impaired capacity framework is based on the cloud server to bolster including, upgrading and erasing records and covering the entrance example of the hunt client from the cloud server. In the visually impaired capacity framework, all reports are separated into fixed-size squares. These pieces are filed by an arrangement of irregular whole numbers produced by an archive related seed. In the perspective of a cloud server, it can just see the squares of encoded records transferred and downloaded. Along these lines, the blind capacity framework releases little data to the cloud Server.

Specifically, the cloud server does not know which squares are of the same archive, even the aggregate number of the archives and the measure of every record. A visually impaired capacity framework is based on the cloud server to bolster A client can unscramble the information just if the qualities implanted in his characteristic keys fulfill the entrance strategy in the cipher text. In CP-ABE, the encrypter holds a definitive power of the access strategy.

RELATED WORK

Searchable encryption is a promising system that gives the hunt administration over the encoded cloud information. It can chiefly be classified into two sorts: Searchable Public-key Encryption (SPE) and Searchable Symmetric Encryption (SSE). Propose the idea of SPE, which bolsters single-watchword hunt over the encoded cloud information. The work is later reached out in to bolster the conjunctive, subset, and range seek inquiries on encoded information. The above proposition require that the indexed lists coordinate all the watchwords in the meantime, and can't return results in a specific request. Propose a positioned seek plan which embraces a veil grid to accomplish cost adequacy. Propose a multi- watchword recovery plan that can give back the top-k pertinent records by utilizing the completely homomorphism encryption. Receive the ascribe based encryption method to accomplish look power in SPE. Propose a positioned. Look plan which receives a cover framework to accomplish cost effectiveness. Propose a multi-catchphrase recovery plan that can give back the top-k applicable archives by utilizing the completely homomorphism encryption. Embrace the credit based encryption strategy to accomplish seek power in SPE. The proposition can accomplish rich functionalities, for example, Multi-watchword and positioned results, yet requires the calculation of significance scores for all archives contained in the database. This operation causes enormous calculation over-burden to the cloud server and is in this manner not suitable for expansive scale datasets.

SYSTEM MODEL

The framework model in the EMRS comprises of three elements: information proprietor, seek clients and cloud server. The information proprietor keeps a substantial gathering of records D to be outsourced to a cloud server in a scrambled structure C. In the framework, the information proprietor sets a catchphrase lexicon W which contains d catchphrases. To empower seek clients to inquiry over the encoded records, the information proprietor assembles the scrambled Record At the point when an inquiry client needs to look over the scrambled records, she first gets the mystery key from the information proprietor. At that point, she picks a conjunctive catchphrase set S which contains l intrigued catchphrases and figures a trapdoor T counting a catchphrase related token stag and the encoded inquiry vector Q. At long last, the pursuit client sends stag, Q, and after accepting stag, Q, and k from the pursuit client, the cloud server utilizes the stag to get to the list z in the visually impaired capacity and figures the significance scores with the encoded question vector Q. At that point, the cloud server sends back descriptors (DSC) of the top-k reports that are most significant to the looked watchwords. The hunt client can utilize these descriptors to get to the visually impaired capacity framework to recover the encoded reports.

An entrance control system, e.g., property based encryption, can be executed to deal with the hunt client's unscrambling capacity. Certiorari number k to the cloud server to ask for the most k applicable results.

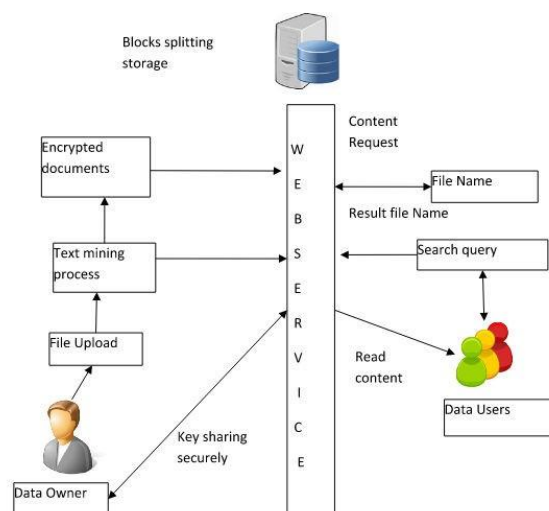
SECURITY REQUIREMENTS

Privacy of Documents and Index

Documents what's more, record ought to be encoded before being outsourced to a cloud server. The cloud server ought to be averted from prying into the outsourced archives and can't conclude any relationship between the archives and watchwords utilizing the list. Trapdoor Privacy: Since the hunt client might want to keep her hunts from being presented to the cloud server, the cloud server ought to be kept from knowing the careful catchphrases contains look client.

Trapdoor confidentiality

Since the pursuit client might want to keep her quests from being presented to the cloud server, the cloud server ought to be kept from knowing the accurate catchphrases contained in the trapdoor of the seek client.



Trapdoor Unlink ability

The trapdoors must not be linkable, which implies the trapdoors ought to be completely diverse regardless of the fact that they contain the same catchphrases. As such, the trapdoors ought to be randomized as opposed to decided. The cloud server can't reason any relationship between two trapdoors.

Covering Access Pattern of the Search User: Access example is the arrangement of the sought results. In the EMRS, the entrance example ought to be completely covered from the cloud server. Specifically, the cloud server can't take in the aggregate number of the records put away on it nor the extent of the looked report notwithstanding when the inquiry client recovers this report from the cloud server.

BLIND STORAGE SYSTEM

A visually impaired capacity framework is based on the cloud server to bolster including, upgrading and erasing reports and covering the entrance example of the pursuit

client from the cloud server. In the visually impaired capacity framework, all archives are separated into fixed-size squares.

$$\begin{cases} p'_j = p''_j = p^*_j, & \text{if } S_j = 1 \\ p'_j = \frac{1}{2}p^*_j + r, & p''_j = \frac{1}{2}p^*_j - r, \text{ otherwise} \end{cases} \quad (1)$$

These pieces are listed by an arrangement of irregular numbers created by a report related seed. In the perspective of a cloud server, it can just see the squares of encoded report transferred and downloaded. In this manner, the blind capacity framework releases little data to cloud.

$$Enc_{(K_i \oplus \Phi(j))}(H(id_i)||size_i||data) \quad (2)$$

And the rest of the blocks of d_i is as

$$Enc_{(K_i \oplus \Phi(j))}(H(id_i)||data) \quad (3)$$

Algorithm 1 Initialize F

```

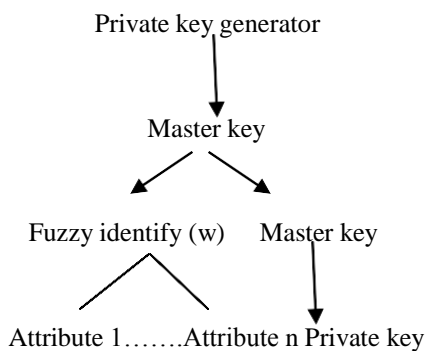
1: for each keyword  $\omega \in W$  do
2:   Set  $t$  an empty list
3:   for each document  $d_i$  containing the keyword  $\omega$  do
4:     Get the associated vector  $P$  of  $d_i$ 
5:     Choose a random number  $x$ 
6:      $Dsc \leftarrow ABE_{v_i}(id_i||K_i||x)$ 
7:     Append the tuple  $(Dsc, P)$  to  $t$ 
8:   end for
9:    $F[\omega] = t$ 
10: end for
11: return  $F$ 
    
```

$$q^* = \{rq, r, t\} \quad (4)$$

$$\begin{cases} q'_j = q''_j = q^*_j, & \text{if } S_j = 0 \\ q'_j = \frac{1}{2}q^*_j + r', & q''_j = \frac{1}{2}q^*_j - r', \text{ otherwise} \end{cases} \quad (5)$$

The search user search the chooses a random integer r' and split the vector q^* into $(d+2)$ dimension vector.

Attribute Based Encryption access policy:



Specifically, the cloud server does not know which squares are of the same report, even the aggregate number of the reports and the span of every record. Besides, all the reports and file can be put away in the

visually impaired capacity framework to accomplish a searchable encryption plan.

PROPOSED SCHEME

Subsequent to the scrambled reports and list z are both put away in the blind capacity framework, we would give the general development of the visually impaired capacity frame specific development of CP- ABE is out of extent of this paper and we just give a straightforward sign here work Development of Blind Storage, Encoded Database Setup, Trapdoor Generation, Efficient what's more, Secure Search, and Retrieve Documents from Blind Storage. At that point, for every record d_i , the information proprietor picks a 192-piece key.

EFFECTIVE AND SECURE SEARCH

After accepting Q , $stag$, and k , the cloud server parses the $stag$ to get an arrangement of whole numbers in the reach. At that point, the cloud server gets to list z in the visually impaired capacity and recovers the pieces listed by the whole numbers to acquire the tuples on these pieces. Note that, these pieces comprise of the pieces and some fake squares. For each recovered encoded significance vector P , property based encryption as an entrance control procedure can be executed to oversee pursuit client's decoding capacity.

$$\begin{aligned} Score_i &= P \cdot Q \\ &= \{M_1^T \cdot p', M_2^T \cdot p''\} \cdot \{M_1^{-1} \cdot q', M_2^{-1} \cdot q''\} \\ &= p' \cdot q' + p'' \cdot q'' \\ &= p^* \cdot q^* \\ &= (p, \epsilon, 1) \cdot (rq, r, t) \\ &= r(pq + \epsilon) + t \end{aligned}$$

COMPUTATION OVERHEAD:

We evaluate the performance of the EMRS through simulations and compare the time cost with Cao's [6]. We apply a real dataset National Science Foundation Research Awards Abstracts 1990-2003 [7], by randomly selecting some Documents. Then, we conduct real-world experiments on a 2.8Hz- processor, computing machine to evaluate the performance of index construction and search phases. Moreover, We implement the trapdoor generation on a 1.2GHz smart phone. We would show the simulation experiments of the EMRS, and demonstrate that the computation overhead of index construction and trapdoor generation are almost the same compared with that of Cao's [6]. Then we would compare the execution time of search phase with Cao's [6] and show that the EMRS achieves better search efficiency.

RECOVER DOCUMENTS FROM BLIND STORAGE

In the event that the pursuit client's qualities fulfill the entrance arrangement of the archive, the hunt client can unscramble the descriptor utilizing her mystery credit keys to get the archive id idi and the related symmetric key Ki . Produce a sufficiently long piece number through the capacity utilizing the seed fi , parse it as an arrangement of numbers in the reach and pick the first integers numbers as the set Sof . Recover the pieces ordered by these numbers from the scrambled database D through visually impaired capacity framework. What's more, this mix can promote cover access example of the inquiry client since the cloud server even does not know the quantity of records that the inquiry client requires.

$$P = \frac{1}{2^{\eta r} 2^{\mu i} 2^{\mu \eta q}} = \frac{1}{2^{\eta r + \mu i + \mu \eta q}} \quad (7)$$

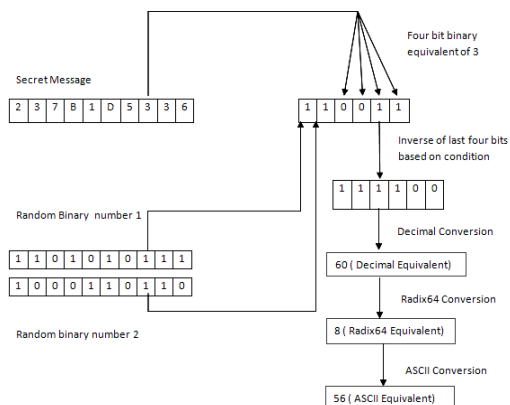
$$P' = \frac{1}{2^{2\alpha + s i z e_{\omega} + n \eta}} \quad (8)$$

SECURITY ANALYSIS

We give investigation of the EMRS as far as condentiality of records and file, trapdoor protection, trapdoor unlink ability and disguising access example of the hunt client.

SECRECY OF DOCUMENTS AND INDEX

The records are encoded by the customary symmetric cryptography method before being outsourced to the cloud server. Without a right key, the hunt client and cloud server can't unscramble the archives. Concerning file confidentiality, the pertinence vector for every archive is encoded utilizing the mystery key $M1$, $M2$, and S . What's more, the descriptors of the archives are scrambled utilizing CP-ABE system. Therefore, the cloud server can just utilize the record z to recover the encoded importance vectors without knowing any extra data, for example, the relationship between the archives furthermore, the watchwords. Also, just the inquiry client with right quality keys can decode the descriptor $ABEi (idijjKijjx)$ to get the report id and the related symmetric key. In this way, the confidentiality of archives and record can be well secured.



TRAPDOOR PRIVACY

At the point when a hunt client produces her trapdoor including the watchword related token $stag$ and encoded inquiry vector Q , she arbitrarily picks two numbers r and t . At that point, for the question vector q , the hunt client augments it as $(rq; r; t)$ and scrambles the question vector utilizing the mystery key $M1;M2$ and S . Consequently, the inquiry vectors can be very surprising regardless of the fact that they contain same catchphrases. Also, we utilize the safe capacity 9 and 0 to offer the inquiry client some assistance with computing catchphrase related token $stag$ utilizing the mystery key $K9$. Without the mystery key $M1;M2; S$ and $K9$, the cloud server can't pry into the trapdoor. What's more, the pursuit client can add sham numbers to the set Sf to cover what she is genuinely scanning for. In this manner, the watchword data in the trapdoor is completely covered from the cloud server in the EMRS and trapdoor security is all around ensured.

TRAPDOOR UNLINK ABILITY

Trapdoor unlink ability is defined as that the cloud server cannot deduce associations between any two trapdoors. Even though the cloud server cannot decrypt the trapdoors, any association between two trapdoors may lead to the leakage of the search user's privacy. We consider whether the two trapdoors including $stag$ and the encrypted query vector Q can be linked to each other or to the keywords

DISGUIISING ACCESS PATTERN OF THE SEARCH USER

The entrance design implies the succession of the sought results, the look client specifically acquires the related records from the cloud server, which might uncover the relationship between the hunt demand and the records to the cloud server. In the EMRS by adjusting the visually impaired capacity framework, access example is very much covered from the cloud server. Subsequent to the headers of the squares are encoded with the piece number j what's more, every descriptor has an irregular cushioning,

$$P_{err}(\gamma, \alpha, \kappa) \leq \max_{n \geq \frac{\kappa}{\alpha}} \sum_{i=0}^{n-1} \binom{\lceil \alpha n \rceil}{i} \left(\frac{\gamma-1}{\gamma}\right)^i \left(\frac{1}{\gamma}\right)^{\lceil \alpha n \rceil - i} \quad (9)$$

INDEX CONSTRUCTION

Considering a extensive integer of information and look clients in a cloud location, searchable encryption procedure should to authorize security saving multi- watchword seek and return archives in a request of high demand. Cash's scheme chains multi-keyword search, but cannot return results in a specific order of the relevance score. Naveed's scheme equipment the blind storage space system to protect the access pattern but it only supports single-keyword search and returns undifferentiated results Upon receiving $stag$, the cloud server can use $stag$ to access blind storage and recover the encrypted consequence vector on the blocks indexed by the $stag$. These blocks consist of blocks of permit containing the $stag$ -related keyword and some replica blocks. Thus, the EMRS can considerably reduce the number of documents which are significant to the searched keywords. Then, the cloud server only needs to compute the internal product of two $(d+2)$ -dimension vectors for the computing relevance scores for all documents as that in Cao's scheme.

Finally, we adopt the index z via the blind storage in the EMRS to improve search efficiency and conceal the access pattern of the search user. For each keyword W , we need to build the list z of tuples $(Abe(id_{ij}K_{ij}x); P)$ of documents that contain the keyword and upload it using the B. Build function. So the computation complexity to build the index z is $O(\%d)$, where $\%$ represents the average number of tuples contained in the list z and is no more than the number of document m . Since the access pattern is not considered in most schemes, we are not going to give the specific comparison of the implementation of the blind storage.

Upon receiving $stag$, the cloud server can use $stag$ to access blind storage and retrieve the encrypted relevance vector on the blocks indexed by the $stag$. These blocks consist of blocks of documents containing the $stag$ -related keyword and some dummy blocks the computation cost of search phase is mainly affected by the number of documents in the dataset and the size of the keyword dictionary. In our experiments, we implement the index on the memory to avoid the time-cost. Although the time costs of search operation are linearly increasing in both schemes, the increase rate of the EMRS is less than half of that in Cao's scheme.

Each index number is nb -bit long. In the EMRS, We modify the way the search user computes the sequence Sf that indexes the blocks by adding some dummy integers to Sf to conceal what the search user is searching for. The communication comparison is shown in As we can see, even though the EMRS requires a little more communication overhead, the EMRS can achieve more functionalities compared with [2], [5] as shown in and better search efficiency compared with [6] as When the system is once setup, including generating encrypted documents and index, the communication overhead is mainly influenced by the search phase. In this section, we would compare the communication overhead among the EMRS, Cash's scheme [8], Cao's scheme [6] and Naveed's scheme [7] when searching over the cloud server. Since most existing schemes of SSE only consider obtaining a sequence of results rather than the related documents, the Comparison here would not involve the communication of retrieving the documents. In Cao's scheme [6], the search user needs to compute the trapdoor and send it to the cloud server. Then it can obtain the searched results. The communication overhead in Cao's is $2(d C2)_q$, where d represents the size of the keyword dictionary and each dimension of the encrypted query vector is q -bit long. According to Cash's scheme [7], when a search user needs to continuously compute the x token until the cloud server sends stop, which indicates that the total number of the x tokens is linear to $\%$, the number of documents containing the keyword related to the $stag$.

The size of the returned results in the EMRS is mainly affected by the choice of the security parameters, n and k . The larger these two numbers are, the higher security guarantee the scheme provides, Cao et al. [6] propose a privacy-preserving multi-keyword search scheme that supports ranked results by adopting secure k -nearest neighbors (kNN) technique in searchable encryption. The proposal can achieve rich functionalities such as multi-

keyword and ranked results, but requires the computation of relevance scores for all documents contained in the database. This operation incurs huge computation overload to the cloud server and is therefore not suitable for large-scale datasets.

SCOPE OF RETURNED OUTCOMES

The size of the resumed fallouts in the EMRS is mostly affected by the select of the security parameters, n and k . The larger these two numbers are, the sophisticated safety assurance the system provides, The size of resumed outcomes for every document can be a^*size_w blocks, which comprise the blocks of examined document and replica blocks. Moreover, the search user can require many documents at one time and thus can circumvent demanding dummy blocks. The EMRS affords stability factors for search users to fulfil their altered requirements on communication and computation cost, as well as confidentiality.

CONCLUSION

In this paper, we have projected a multi-keyword ranked search pattern to facilitate exact, effective and protected search over encrypted mobile cloud data. Security analysis have demonstrated that proposed scheme can effectively achieve confidentiality of brochures and catalogue, access secrecy, access unlinkability, and obscuring access pattern of the search user. General concert evaluations have shown that the proposed scheme can achieve better efficiency in terms of the functionality and computation in the clouds compared with surviving ones. For the prospect work, we will consider on the validation and access control issues in searchable encryption performance.

REFERENCES

1. H.Liang,L.X.Cai,D.Huang,X.Shen,andD.Peng"An SMDP based service model for interdomain resource allocation in mobile cloud networks",IEEE Trans.veh.technol., vol.61, no.5,pp.2222-2232,jun 2012.
2. M. M. E. A. Mahmoud and X. Shen, "A cloud-based scheme for protecting source-location privacy against hotspot-locating attack in wireless sensor networks,"*IEEE Trans. Parallel Distribution System*, vol. 23, no. 10, pp. 1805_1818, Oct. 2012.
3. Q. Shen, X. Liang, X. Shen, X. Lin, and H. Y. Luo, "Exploiting geo distributed clouds for a e-health monitoring system with minimum service delay and privacy preservation," *IEEE J. Biomed. Health Information*, vol. 18, no. 2, pp. 430_439, Mar. 2014
4. H. T. Dinh, C. Lee, D. Niyato, and P. Wang, "A survey of mobile cloud computing: Architecture, applications, and approaches," *Wireless Communication Mobile Computing*, vol. 13, no. 18, pp. 1587_1611, Dec. 2013.
5. H. Li, Y. Dai, L. Tian, and H. Yang, "Identity-based authentication for cloud computing," in *Cloud Computing*. Berlin, Germany: Springer-Verlag, 2009, pp. 157_166.
6. W. Sun, et al., "Privacy-preserving multi-keyword text search in the cloud supporting similarity-based ranking," in Proc. 8th ACM SIGSAC Symp.In., Comput. Commun. Secur., 2013, pp. 71_82.

7. B.Wang, S.Yu,W. Lou, andY. T. Hou, "Privacy-preserving multi-keyword fuzzy search over encrypted data in the cloud," in Proc. IEEE INFOCOM, Apr./May 2014, pp. 2112_2120.
8. E. Stefanov,C.Papamanthou, and E. Shi, "Practical dynamic searchable encryption with small leakage,"in Proc. NDSS, Feb. 2014.
9. Y. Yang, H. Li, W. Liu, H. Yang, and M. Wen, "Secure dynamic searchable symmetric encryption with constant document update cost,"in Proc.GLOBECOM, Anaheim, CA, USA, 2014.
10. D. Cash, S. Jarecki, C. Jutla, H. Krawczyk, M.- C. Ro³u, and M. Steiner, "Highly-scalable searchable symmetric encryption with support for Boolean queries," in Proc. CRYPTO, 2013, pp. 353_373.