



---

## International Journal of Intellectual Advancements and Research in Engineering Computations

---

### Secure virtual machine based multi keyword data sharing on public cloud

Mrs. B.Deepa<sup>1</sup>, G.V. Kavipriya<sup>2</sup>, K. Lakshmipriya<sup>2</sup>, P. Karthik<sup>2</sup>

<sup>1</sup>Assistant Professor

<sup>2</sup>UG Students

---

#### ABSTRACT

The project solves and defines the difficulty of Multi-keyword Ranked Search over Encrypted cloud data (MRSE) while preserving firm system wise privacy in the cloud computing hypothesis. For the protection of data privacy, sensitive data has to be encrypted before outsourcing, which makes effective data utilization a very challenging task. Even though searchable encryption method allows users to firmly search over encrypted data all the way through the keywords, they carry only search i.e., Boolean. They are not yet enough to meet the utilization of the data successfully because there is instinctively demanded by large number of data files and users located in cloud. Hence, it is required to allow multiple keywords in the search request and return documents in the order of their significance to the keywords. The keyword i.e., Boolean of the search technique only produce the unsorted result. A valuable method proposed for this difficult problem is privacy conserving search over encrypted cloud data. After the data have been encrypted and outsourced by the data owner this method establishes a set of privacy desires to secure cloud data utilization system during splitting the cloud data and storing the chunk data in different servers. Among different multi-keyword ontology, this method chooses the well-organized similarity measure of “coordinate matching” for searching technique. Then according to Top-K query scheme the sorted results are created.

**Index Terms:** Cloud, MRSE, OTP, Product similarity.

---

#### INTRODUCTION

Cloud computing is a casually using the real-time communication network and connect large number of computers that depict different types of computing concepts. Non-ambiguous technical or scientific description in cloud computing has not been accepted. In science, cloud computing is a one kind of the distributed computing network and capability to run a program on many related computers at the similar time. Cloud computing is called as a utility of the computing since it uses pay per use paradigm. In cloud computing, users can also right to use a variety of resources like storage, programs, and application development platforms. Cloud computing is an emerging technology, and it is also called as utility because client are used to store their data in the cloud

server. In cloud server data can also be leaked to hackers therefore encrypted the data before sent to the cloud for data privacy [1].

Cloud computing is transform how businesses uses the information technology. Several trends are opening up the period of cloud computing, which is a development of Internet-based and use of technology of the computer. More controlling processors and increasingly cheaper, collectively with the Software as a Service computing architecture, are transforming data center into pool of computing service on a huge scale.

The escalating network bandwidth and reliable yet stretchy network connections make it even possible that users can now give to high quality services from data and software that dwell exclusively on remote data centers. To protect

privacy of the data and struggle unwanted accesses in the cloud, Cloud Service Providers (CSP) frequently put in force users data security all the way through mechanisms like virtualization and firewalls. However, these techniques do not protect users privacy from the CSP itself since the CSP possesses full organize of the system hardware and minor levels of software stack. Therefore, encryption before outsourcing the data of the cloud; this, however, difficulty in the traditional data utilization service based on plain text keyword investigate. The small solution of downloading all the data and decryption close by is clearly unreasonable, due to the vast amount of bandwidth cost in cloud extent systems. Thus, exploring privacy preserving and effective search service over encrypted cloud data are of dominant importance [2]. Considering the potentially large number of on-demand data users and enormous quantity of outsourced data document in the cloud, this crisis is particularly challenging as it is enormously difficult in system usability, performance and scalability. On the one hand, to meet the effectual data retrieval, the huge amount of documents insists the cloud server to perform result significance ranking, instead of returning results which are not similar. Such ranked of the search system enables the data users to discover quickly the most appropriate information, rather than sorting all the way through each and every match in the content set. On the other hand, to progress the search result exactness as well as to increase the experience of user searching. To provide more exactness to the end Privacy Conserving in Cloud Documents in excess of users result is done by searching, the unlabeled data keywords are incorporated in the index of the server and then searching was done this search results is then characterized, and then they are sorted in their splitting up using Top-k query algorithm. TOP-k selection queries will assist to sort the related data and provide the accurate data to the end user [3].

## RELATED WORK

Organizations, companies store more and more valuable information is on cloud to protect their data from virus, hacking. The benefits of the new

computing model include but are not limited to: relief of the trouble for storage administration, data access, and avoidance of high expenditure on hardware mechanism, software, etc. Ranked search improves system usability by normal matching files in a ranked order regarding certain relevance criteria (e.g., keyword frequency). As directly outsourcing relevance scores will drips a lot of sensitive information against the keyword privacy, We proposed asymmetric encryption with ranking result of queried data which will give only expected data.

## EXISTING SYSTEM

The cloud server hosts intermediary information storage spaces and service can also be regained. Since information may holds susceptible information, defending information cannot be fully entrusted by cloud servers. For this reason, files that can be should be encrypted. Any information type that is leaking that would have an effect on data solitude is regard as insufferable. To get mutually the effectual data recovery need, the huge quantity of documents stress the cloud server instead of recurring not similar results and to carry out result impact ranking. Such data users approaches ranked search system to discover the most appropriate information rapidly, rather than onerously sorting all the way through every match in the content set. Ranked of search remove the needless network traffic by distribution reverse only the most of the valid data, which is also very attractive in the “pay- as-you-use” cloud perception. For segregation support, such grade of the process, however, cannot be seeped out any information related to the keyword. To get good search consequence correctness as well as to improve the user sharp skill, such ranking system essentially hold up multiple keywords search, too rude results were produced by single keyword of search [4].

## Draw Backs of Existing System

- Accurate data we are not getting.
- Users will not get adequate required ranking functionality.
- Data that can be sharing will not be safe.



## RANKED SEARCH

The multi-keyword search method checks whether queried keywords exist in a document or not. If the user searches for a single or more keywords, there will possibly be many correct matches where some of them may not be useful for the user at all. Therefore, it is difficult to decide as to which documents are the most relevant. I add ranking capability to the system by adding extra index information for frequently occurring keywords in a file. With ranking, the user can retrieve only the top  $\tau$  matches where  $\tau$  is chosen by the user. In order to rank the documents, a ranking function is required, which assigns relevancy scores to each document matching to a given search query. One of the most widely used metrics in information retrieval is the term frequency. Term frequency is defined as the number of times a keyword appears in a document. Instead of using term frequency itself, we assign relevancy levels based on the term frequencies of keywords. I assume that there are  $\eta$  levels of ranking in our proposed method for some integer  $\eta \geq 1$ . For each document, each level stores an index for frequent keywords of that document in a cumulative way in descending order. This basically means that its level index includes all keywords in the  $(i + 1)$ th level, and the keywords that have term frequency for the  $i$ th level. The higher the level, the higher the term frequency of the keywords is. For instance, if  $\eta = 3$ , level 1 index includes keywords that occur at least once in the document while levels 2 and 3 include keywords that occur at least, say 5 times and 10 times<sup>4</sup>, respectively. There are several variations for relevancy score calculations, and we use a very basic method. The relevancy score of a document is calculated as the number representing the highest level search index that the query index matches. All the keywords that exist in a document is included in the first level search index of that document. The other higher level indices include the frequent keywords that also occur in its previous level, but this time they have to occur the number of times, which are at least the term frequency of the corresponding level. The highest level includes only the keywords that have the highest term frequency. In the oblivious search phase, the server starts comparing

the user query against the first level indices of each document. The matching documents found as a result of the comparison in the first level is then compared with the search indices in the other levels according to the Algorithm 1. In this method, some information may be lost due to the ranking method employed here. Rank of two documents will be the same if one involves all the queried keywords infrequently and the other involves all the queried keywords frequently except one infrequent one. The rank of the document is identified with the least frequent keyword of the query. We tested the correctness of our ranking method by comparing with a commonly used formula for relevance score calculation, which is given in the following:

The number of levels and the term frequency of each level can be chosen in any convenient way.

Algorithm Ranked Search

```
{
for all documents Ri do
{
Compare (level1 index of Ri , query
index) j = 1 while match do
{
increment j
Compare (level j indices of Ri, query index)
end while
}
rank of Ri = highest level that match
with query index end for
}
```

$$\text{Score}(W, R) = \sum_{t \in W} \frac{1}{|R|} (1 + \ln fR, t) \cdot \ln(1 + \frac{1}{ft})$$

Here  $W$ ,  $fR, t$ ,  $ft$ , and  $M$  denote the set of searched key- words, the term frequency of term  $t$  in file  $R$ , the number of files that contain term  $t$ , and the number of files in the database, respectively.  $|R|$  is the length of the file  $R$ . We use a synthetic database to compare the two ranking methods. We assume that there are 1000 files of equal lengths in the database and 3 keywords are searched for. We further assume that the number of files containing the queried keywords ( $ft$ ) is equal to 200 and only 20 of those contain all 3 keywords. Term frequencies of the keywords in 20 matching files are chosen uniformly randomly between 1 and 15 and  $\eta = 5$  in our proposed ranking method. In 40% of the time, the top match for a given

relevance score, is also the top match for our proposed ranking method, and 100% of the time in the top 3 matches of our ranking method. Additionally, in 80% of the time, at least 4 of the top 5 matches for the given relevance score is also in the top 5 of our proposed ranking method. Note that since we assume  $f_t$  is the same for all  $t \in W$ , changing  $f_t$  has no effect on the performance of both methods. As can be observed from these experimental figures, while the performance of the proposed method is unacceptable levels, the choice of the method parameters (especially  $\eta$  and term frequency of each level) depends very much on the characteristics of the database and the documents. While this new method necessitates an additional  $r$ -bit storage per level for a document, it reduces the communication overhead of the user since matches with low rank documents will not be retrieved unless the user requests. Considering  $\eta$  search indices are stored instead of a single search indexes per document, storage overhead for indexing mechanism increases  $\eta$  times due to ranking. This additional cost is not a burden for the server since the index sizes are usually negligibly small compared to actual document sizes.

Algorithm Used

## RSA ALGORITHM

This algorithm is used to encrypt n decrypt file contents. It is an asymmetric algorithm. The RSA algorithm involves three steps: key generation, encryption and decryption. Key generation RSA involves a public key and a private key.

The public key can be known to everyone and is used for encrypting messages. Messages encrypted with the public key can only be decrypted using the private key. The keys for the RSA algorithm are generated the following way

- Choose two distinct prime numbers  $a$  and  $b$ .
- Compute  $n = ab$ .  $n$  is used as the modulus for both the public and private keys
- Compute  $\phi(n) = (a-1)(b-1)$ , where  $\phi$  is Euler's totient function.
- Choose an integer  $e$  such that  $1 < e < \phi(n)$  and greatest common divisor of  $(e, \phi(n)) = 1$ ; i.e.,  $e$  and  $\phi(n)$  are co prime.  $e$  is released as the public key exponent having a short bit-length .

## K-Nearest Neighbor

K-nearest neighbor search identifies the top  $k$  nearest neighbors to the query. This technique is commonly used in predictive analytics to estimate or classify a point based on the consensus of its neighbors. K-nearest neighbor graphs are graphs in which every point is connected to its  $k$  nearest neighbors. The basic idea of our new algorithm: The value of  $d_{max}$  is decreased keeping step with the ongoing exact evaluation of the object similarity distance for the candidates. At the end of the step by step refinement,  $d_{max}$  reaches the optimal query range  $E_d$  and prevents the method from producing more candidates than necessary thus fulfilling the  $r$ -optimality criterion.

- Nearest Neighbor Search ( $q, k$ ) // optimal algorithm
- Initialize ranking = index. increm-ranking ( $F(q), df$ )
- Initialize result = new sorted-list (key, object)
- Initialize  $d_{max} = w$
- While  $o = \text{ranking. getNext}$  and  $d(o, q) \leq d_{max}$ , do
- If  $d(o, q) < d_{max}$  then result. insert ( $d(o, q), o$ )
- If result. length  $\geq k$  then  $d_{max} = \text{result}[k].key$
- Remove all entries from result where key  $> d_{max}$
- End while Report all entries from result where key  $\leq d_{max}$

## EXPECTED RESULTS

### Data Encryption and decryption Result

When RSA algorithm is applied on the data then we get encrypted data and that encrypted data are store on the cloud. User can access the data after downloading and decrypting file. For encryption and decryption keys are provided

### Ranking Result

When any User request for the data then Ranking is done on requested data using k-nearest neighbor algorithm. For Ranking co- ordinate matching principle is used. After, ranking user gets the expected results of the query.

## Alert System Results

If any unauthorized User tries to access or updating the data on cloud, then alert will be generated in the form of mail and messages. The alert intimates the authorized user.

## CONCLUSION

Thus we proposed the problem of multiple-keyword ranked search over encrypted cloud data,

and construct a variety of security requirements. From various multi-keyword concepts, we choose the efficient principle of coordinate matching. We first propose secure inner data computation. Also, we achieve effective ranking result using k-nearest neighbor technique. This system is currently work on single cloud Provide better security in multi-user systems.

## REFERENCES

- [1]. L. Vaquero, L. Rodero-Merino, J. Caceres, and M. Lindner, "A break in the clouds: towards a cloud definition," *ACM SIGCOMM Computer Communication Review*, vol. 39, no. 1, pp. 50–55, 2008.
- [2]. Huang, "Mobile cloud computing," *IEEE COMSOC Multimedia Communications Technical Committee (MMTC) E-Letter*, 2011.
- [3]. J. Oberheide, K. Veeraraghavan, E. Cooke, J. Flinn, and F. Jahanian, "Virtualized in-cloud security services for mobile devices," in *Proceedings of the First Workshop on Virtualization in Mobile Computing*. ACM, 2008, 31–35.
- [4]. J. Oberheide and F. Jahanian, "When mobile is harder than fixed (and vice versa): demystifying security challenges in mobile environments," in *Proceedings of the Eleventh Workshop on Mobile Computing Systems* ACM, 2010, 43–48.
- [5]. A. Moffat, T. C. Bell et al., *Managing gigabytes: compressing and indexing documents and images*. Morgan Kaufmann Pub, 1999.
- [6]. Song, D. Wagner, and A. Perrig, "Practical techniques for searches on encrypted data," in *Security and Privacy, 2000. S&P 2000. Proceedings. 2000 IEEE Symposium on*. IEEE, 2000, 44–55.
- [7]. Boneh, G. Di Crescenzo, R. Ostrovsky, and G. Persiano, "Public key encryption with keyword search," in *Advances in Cryptology- Eurocrypt*. Springer, 2004, 506–522.
- [8]. R. Curtmola, J. Garay, S. Kamara, and R. Ostrovsky, "Searchable symmetric encryption: improved definitions and efficient constructions," in *Proceedings of the 13th ACM conference on Computer and communications security*. ACM, 2006, 79–88.
- [9]. Y. Chang and M. Mitzenmacher, "Privacy preserving keyword searches on remote encrypted data," in *Applied Cryptography and Network Security*. Springer, 2005, 391–421.
- [10]. S. Zerr, D. Olmedilla, W. Nejdl, and W. Siberski, "Zerber+ r: Topk retrieval from a confidential index," in *Proceedings of the 12th International Conference on Extending Database Technology: Advances in Database Technology*. ACM, 2009, 439–449.
- [11]. C. Wang, N. Cao, K. Ren, and W. Lou, "Enabling secure and efficient ranked keyword search over outsourced cloud data," *Parallel and Distributed Systems, IEEE Transactions on*, 23(8), 2012, 1467–1479.
- [12]. N. Cao, C. Wang, M. Li, K. Ren, and W. Lou, "Privacy-preserving multi-keyword ranked search over encrypted cloud data," *Parallel and Distributed Systems, IEEE Transactions on*, 25(1), 2014, 222–233.
- [13]. J. Zobel and A. Moffat, "Inverted files for text search engines," *ACM Computing Surveys (CSUR)*, 38, (2), 6, 2006.
- [14]. R. Agrawal, J. Kiernan, R. Srikant, and Y. Xu, "Order preserving encryption for numeric data," in *Proceedings of the 2004 ACM SIGMOD international conference on Management of data*. ACM, 2004, 563–574.

- [15]. D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent dirichlet allocation," the Journal of machine Learning research, 3, 2003, 993–1022.
- [16]. J. Ramos, "Using tf-idf to determine word relevance in document queries," Technical report, Department of Computer Science, Rutgers University, 2003.
- [17]. Q. Chai and G. Gong, "Verifiable symmetric searchable encryption for semi-honest-but-curious cloud servers," in Communications (ICC), 2012 IEEE International Conference on. IEEE, 2012, 917– 922.
- [18]. S. Kamara and K. Lauter, "Cryptographic cloud storage," in Financial Cryptography and Data Security. Springer, 2010, 136– 149.
- [19]. M. Li, S. Yu, K. Ren, W. Lou, and Y. T. Hou, "Toward privacy assured and searchable cloud data storage services," Network, IEEE, vol. 27(4), 2013, 56–62.
- [20]. N. Cao, C. Wang, M. Li, K. Ren, and W. Lou, "Privacy-preserving multi-keyword ranked search over encrypted cloud data," in INFOCOM, 2011 Proceedings IEEE. IEEE, 2011, 829–837.