# A multi-keyword parallel ciphertext retrieval scheme based on inverted index under the robot distributed system

Jiyue Wang[1,*], Xi Zhang[1] and Yonggang Zhu[1]

[1]School of Mechanical Engineering, Zhengzhou University of Science and Technology, Zhengzhou, China
Email: wjiyue@126.com;xdwangxd@163.com;zhuyg11@126.com

## Abstract

The traditional ciphertext retrieval scheme has some problems, such as low retrieval performance, lack of single keyword retrieval mode and limitation of single machine resources in traditional single server architecture. At the same time, for searchable encryption, it needs to balance the data security and retrieval efficiency. In this paper, a multi-keyword parallel ciphertext retrieval system based on inverted index is proposed. The system adopts different index encryption methods to improve the performance of ciphertext retrieval. Through the segmentation of ciphertext inverted index, the block retrieval of inverted index is realized, which overcomes the limitation of single machine resources and improves the retrieval efficiency. By combining the characteristics of distribution, the traditional single-machine retrieval architecture is extended and multi-keyword parallel retrieval is realized. The experimental results show that compared with SSE-1 scheme, the proposed scheme can improve the efficiency of retrieval, update and other operations on the premise of ensuring the security of ciphertext data, achieve multi-keyword retrieval, and dynamically expand the distributed architecture of the system. Finally, it can improve the system load capacity.

*Corresponding author. Email: wjiyue@126.com

## 1. Introduction

With the development of information technology, the scale of data generated and used by people in daily life and work is increasing. In order to efficiently screen out meaningful data from massive data, document retrieval technology has been widely used. In recent years, the explosive growth of data scale leads to the fact that the local computing and storage resources of users can no longer meet the storage and management needs of huge data volume. Due to the convenience, speed and flexibility of cloud services, more and more users choose to migrate local data to the cloud for storage and management [1] to save the cost of local data management. However, due to the openness, distribution and other characteristics of cloud services, data is stored in the cloud without the user's physical control, which makes the security and privacy of data increasingly prominent. Security of big data has been paid more attention [2-3]. Data encrypted storage can ensure the security and privacy of data to some extent. However, encrypted data makes retrieval operation very difficult [5-7]. In order to realize data retrieval operation on the basis of ensuring data security and user privacy, Searchable Encryption (SE) is generated.

As one of the current privacy protection mechanisms, searchable encryption has been widely studied and used. For searchable encryption of unstructured data, SE scheme [8] based on construction index can effectively support keyword retrieval, thus becoming the main construction strategy of SE scheme. In an open system, an index of a document is generated by a user or a trusted third party and the ciphertext data and ciphertext index are uploaded to the server. The user retrieves the target document by retrieving the secure index. The whole retrieval process is conducted without decrypting the document data and index files. Therefore, the index-based searchable encryption scheme not only ensures the security of ciphertext data, but also makes use of the index to realize the efficient searchability of data. In recent years, inverted index, as one of the most efficient index structures, has been widely used in plaintext data retrieval because of its efficient retrieval characteristics [9,10]. Meanwhile, inverted index structure can save disk space, improve retrieval efficiency, and support incremental update and delete. Therefore, searchable encryption scheme based on inverted index has become one of the main searchable encryption schemes.

Although the searchable encryption scheme based on inverted index has certain advantages, it still has some obvious disadvantages:

1) On the one hand, there are limitations in the SSE-1 scheme. In the scheme, the encryption of index pointer results in much computation in the retrieval process, which affects the retrieval performance.

2) On the other hand, there are limitations in stand-alone resources. When the size of index files increases with the amount of data, the size of indexable data is limited by single machine memory resources, and the retrieval efficiency is limited by single machine computing resources. Finally, they are the limitations for the inverted index structure. For multi-keyword retrieval request, serial retrieval of inverted index structure seriously affects the efficiency of multi-keyword retrieval.

In view of the above shortcomings, this paper adopts different index encryption methods, to extend the traditional single-server model architecture. It designs and implements a searchable encryption scheme based on inverted index and multi-keyword parallel retrieval. In this paper, the distributed management and retrieval architecture of inverted index is designed. Compared with the traditional stand-alone model architecture, this new architecture model can improve the efficient utilization of stand-alone resources, and has good scalability. Thus it can avoid the limitation of retrievable data scale. Meanwhile, the scheme of this paper can realize the parallel retrieval of multiple keywords to some extent, and make up for the deficiency of inverted index structure in multi-keyword retrieval process.

## 2. Related works

Searchable encryption technology not only realizes the detectability of ciphertext data, but also enriches the retrieval form, retrieval structure and user management of ciphertext data to meet the requirements of more secure, accurate and efficient retrieval.

In 2000, Song [11] first studied the problem of searchable encryption and proposed the linear ciphertext scanning scheme of SWP. Although this scheme could basically realize word search, it required linear scanning of all documents, resulting in a linear relationship between the cost and file size, so its retrieval efficiency was low. In view of the disadvantages such as low retrieval efficiency of SE scheme proposed in [11], the following encrypted searchable schemes are usually to build a secure searchable index, generate a trapdoor matching index through the key, and use the trapdoor to match the index content hidden in the cloud to obtain the retrieval result of ciphertext. References [12-14] constructed a secure index structure based on inverted index. In 2003, Palupi [15] proposed a Z-IDX based on security index to quickly realize the search of massive ciphertext data. Qian [16] used Bloom Filter to construct the index of each document. Although the scheme had the advantage of efficient retrieval, the scheme was misidentified due to the Collision of Hash functions. In response to this problem, Mu [17] proposed an alternative scheme for security definition and structure, which made up for the defect of bit error rate. Wang [18] standardized Symmetric Searchable Encryption (SSE) and its security target, and proposed an encrypted searchable scheme based on inverted index. In this scheme, the list of documents corresponding to each keyword was encrypted and blurred into an array. But after keyword retrieval, the location and contents of the corresponding inverted list were exposed to the cloud server. Therefore, a keyword could only be retrieved once before the index was rebuilt. Then, the concept of dynamic searchable encryption was proposed, and a reverse index of encryption was constructed to support dynamic operations, such as document update. Xia [19] proposed a multi-keyword retrieval scheme that could dynamically add, delete, modify and check documents.

Earlier SE mechanisms can only support the single keyword retrieval, so these schemes have the same limitation that joint multi-keyword retrieval is not supported. In order to extend the query mode to achieve more accurate query, Ballard [20] proposed two efficient SE mechanisms to realize the connection keywords retrieval based on symmetric cryptography and public key cryptography respectively, but both of them needed to ensure that there were no duplicate keywords in the retrieval request. In the SE mechanism based on symmetric cryptography, a linear relationship between the size of trapdoor and the number of files was required. The SE mechanism based on public key cryptography solved this problem by using bilinear mapping to fix the size of the trapdoor. For multi-keyword retrieval, in addition to correlation rank query, fuzzy query proposed

was also an important part of the research on searchable encryption [21-22]. In order to enrich the application scenarios of searchable encryption schemes, relevant schemes for ciphertext retrieval of multi-user shared scenarios were proposed in [23-24].

All the above schemes are based on the SE mechanism with the single-server model architecture. With the increase of index file size, the memory and computing resources of the single-server architecture model can no longer satisfy the current huge amount of data retrieval and management, leading to the decrease of retrieval efficiency. In reference [25], a study on parallel ciphertext inversion index was proposed, and a distributed framework was used to construct the ciphertext inversion index on the server side in parallel. Although the efficiency of index construction and retrieval was improved, the plaintext data and key were exposed to the server during the construction of ciphertext index, which was lack of security. And the retrieval mode was single user single keyword, the application scenario had limitations. At the same time for the multi-keyword retrieval scheme, it had serious effect on the retrieval efficiency. In view of this problem, parallel searchable encryption schemes with multiple keywords have become a hot research area. The red-black tree structure was introduced as the index structure, so that the dynamic SSE could support multi-processor parallel retrieval [26].

Aiming at the problem of multi-keyword retrieval efficiency based on inverted index structure, this paper proposes a distributed parallel retrieval scheme. Based on the limitation of computer resources and the security of data retrieval, the limited single-machine computing and memory resources are fully utilized to expand the scale of data retrieval and optimize the retrieval efficiency.

# 3. Design of robot distributed multi-keyword parallel ciphertext retrieval scheme

The rapid growth of data scale in the cloud seriously affects the retrieval efficiency of ciphertext data. However, the searchable encryption scheme based on the single-server model cannot be applied to ciphertext retrieval in the big data environment due to the limitation of resources. On this basis, this paper uses the inverted index segmentation and the distributed model architecture to realize the multi-keyword parallel ciphertext retrieval scheme in the distributed environment. The distributed platform is used to extend the searchable encryption model architecture to improve the efficiency of multi-keyword retrieval. The structure of the distributed multi-keyword parallel ciphertext retrieval system is shown in figure 1. The scheme is based on the retrieval structure of inverted index, which is composed of client and server. The server side contains a master node and multiple distributed slave nodes. The client

builds ciphertext index and submits retrieval request. The server side realizes the index distribution management and retrieval. Where, the main node of the server divides the index and retrieves the request. It distributes the sub-index and sub-request to the slave nodes of the distributed platform. Each slave node hands the retrieval result to the master node to process the result set, and the master node returns the processing result to the client.
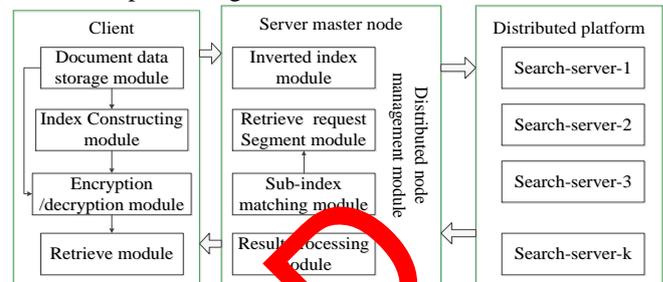


**Figure 1.** Distributed robot multi-keyword parallel ciphertext retrieval system structure

For multi-keyword retrieval request, it is divided into several independent and non-overlapping retrieval requests. It distributes the retrieval request to the corresponding distributed node for parallel retrieval. In order to realize the distributed management of inverted index, four steps need to be completed. 1) The inversion index and data encryption of ciphertext are constructed by the client based on the symmetric key. 2) it submits the ciphertext inverted index and ciphertext data to the server master node. 3) The server-side master node cuts the ciphertext inverted index into several complete and independent sub-indexes. 4) Distribution subindexes are managed by the distribution of individual compute nodes.

## 3.1. Inverted index construction and encryption

In this paper, inverted index is used to achieve efficient retrieval. In order to ensure the safe transmission and retrieval of data, two operations need to be completed before uploading data and index files. First, the user builds an inverted index for the uploaded document data on the client side. The other is to use the local key to match the number of documents before uploading the data.

Inverted indexes are composed of lexicons and Inverted Lists. The dictionary holds all the keywords and the logical Pointers to the inverted list and other information. The inverted list consists of all document identifiers containing keywords. The target document can be obtained by searching the dictionary, finding the corresponding logical pointer to retrieve the keyword, and traversing the relevant

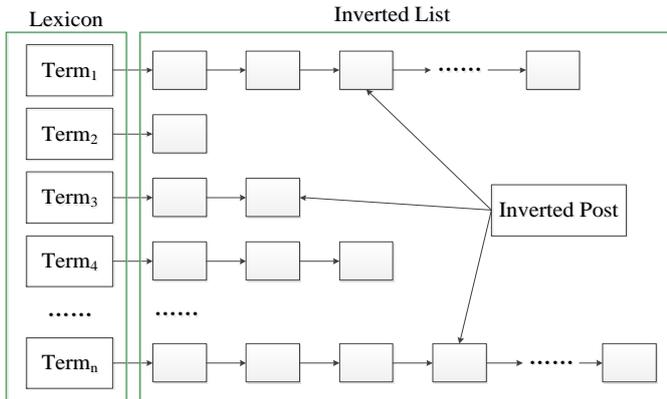inverted list. The inverted index structure is shown as figure 2.



**Figure 2.** Inverted index structure

From the structure of inverted index, the attacker is likely to reproduce the contents of the whole document by index information such as keywords of the dictionary file and document identification of the inverted list. In order to ensure the security of data retrieval, the inverted index file must be encrypted and uploaded.

Against the security threat of plaintext inverted index, it is necessary to encrypt and upload the inverted index to ensure the security of ciphertext data retrieval. Based on some data structures such as array, linked list and lookup table. In the encryption of index files, symmetric encryption should be used to encrypt the index to form a secure inverted index to encrypt related information of documents without changing the index structure and hide the information in the index to realize the retrieval of ciphertext data. The structure of the ciphertext inverted index is shown as figure 3.
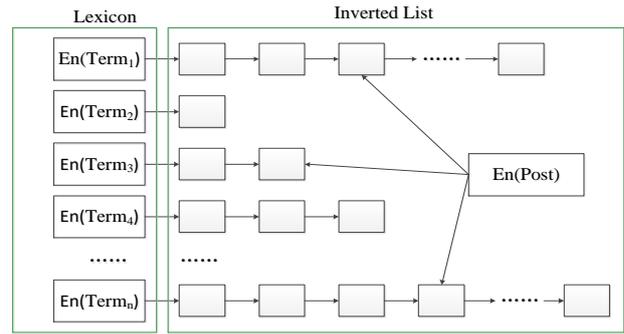


**Figure 3.** Ciphertext inverted index

En(term) encrypts the keywords in the dictionary, but does not encrypt the keywords in the index and the associated pointer to the permute table. En(post) encrypt the information of each inverted item in the inverted list, but does not encrypt the pointer information in the inverted list. Although the inverted ciphertext index files ensure the security of the data in the retrieval process, the increase of the data volume and the encryption of the index will lead to the increase of the index file size, which will affect the memory and computing resources and limit the retrieval efficiency. Therefore, for large-scale ciphertext data retrieval, it is necessary to realize the distributed management and retrieval of index files.

Luence is a full-text search engine architecture that provides a complete query engine, index engine and partial text analysis engine. Lucene provides a simple but powerful application programming interface for full-text indexing and retrieval. In this paper, Luence is used to construct and retrieve the inverted index of text data, and then encryption is used to construct the inverted index of ciphertext.

Definition 1. $T = \{t_1, t_2, \cdots, t_n\}$ is the indexed dictionary file, $t_i (1 \le i \le n)$ is a key word in the dictionary, and $i$ is the number of key word.

**Algorithm 1**. Ciphertext inverted index construction algorithm.

Input. User *key*, document path *path*.

Output. Ciphertext inverted index $En_{key}(\Omega)$.

1) Initialization. Initializing the AES encryption module. Configuring the index *IndexWriter*, including the analyzer, the index store path *directory*.

2) Constructing the ciphertext inverted index.

  a. Reading the document name and encrypting it. Storing it in the document name index field *name_field*.

  b. Using the Chinese word segmentation tool *IKAnalyzer* to segment the document content and form

the keyword set $T = \{t_1, t_2, \cdots, t_n\}$. Ciphertext keyword dictionary *En(T)* is generated by encrypting each keyword and stored in the *content_field* of the document content index.

c. Adding the *name_field* and *content_field* to the index document. Writing the index document to *IndexWriter* and saving the ciphertext index $\Omega$ to disk.

## 3.2. Segmentation and distribution of inverted indexes

In view of the limitation of single machine resources in large-scale index files, the inverted index is divided into several sub-indexes to realize the distributed management and retrieval of index files. Each computation node is responsible for a portion of the index file. On the basis of limited single computer resources, the scale of data retrieval can be extended and the load capacity of the system can be improved to fully improve the efficient utilization of single computer memory and computing resources. To implement index file distribution, two operations are required. First, the whole ciphertext inverted index is divided into several complete and independent sub-indexes based on keywords. Second, it distributes the sub-index to each distributed node.

Definition 2. For any inverted index file $\Omega$, according to equations (1)-(2), it is divided into one independent and equal sub-index file, then they are distributed to each distributed node.

$$SLoad(S_m) \cong \frac{1}{k} \sum_{i=1}^{n} |P_i| \quad (1)$$

Where $L = \{P_1, P_2, \cdots, P_n\}$, P is the indexed inverted list file. $P_i = \{p_{i1}, p_{i2}, \cdots, p_{ij}\}$ is the inverted list mapped with the keyword $t_i$. $p_{ij}$ is the j-th inverted item in the inverted list $P_i$, which contains the document identification information. $k$ is the number of sub-indexes in the shard, as well as the number of distributed nodes.

Each distributed node $S_m$ in distributed node set $S = \{S_1, S_2, \cdots, S_k\}$ manages an approximately equal number of inverted item set $SLoad(S_m)$.

$$\sum_{i=1}^{m_1} t_i <df> \approx \sum_{i=m_1+1}^{m_2} t_i <df> \approx \cdots \approx \sum_{i=m_{k-1}+1}^{n} t_i <df> \quad (2)$$

Where $t_i <df>$ is the keyword $t_i$ corresponding to the number of inverted items in the inverted list. $m_k$ represents the maximum keyword number in the k-th keyword set.

---

**Algorithm 2**. Ciphertext index segmentation algorithm.

Input. Ciphertext inverted index $En_{key}(\Omega)$, distributed node number $k$.

Output. Sub-index set $\{En_{key}(\Omega_1), En_{key}(\Omega_2), \cdots, En_{key}(\Omega_k)\}$.

1) Calculating the segmentation criteria. Calculate the total number of inverted entries p contained in all inverted lists $P_i$ in the inverted index file. According to formula (1), the number of inverted items contained in each sub-index is calculated to determine the segmentation criteria.

2) Dictionary segmentation.

*SLoad(S)* is used as the segmentation standard. According to the keyword order number and the number of inverted items in the corresponding inverted list with formula (2), the order of keywords in the dictionary is divided into $k$ sets to form $k$ keyword dictionaries.

3) Inverted list segmentation. According to the segmentation scheme of dictionary keywords and the mapping relationship between keywords and inverted list $|P_i|$, the range of inverted list files is segmented.

---

Based on the segmentation of dictionary keywords and inverted list, $k$ complete and independent sub-index sets are obtained, namely $\{En_{key}(\Omega_1), \cdots, En_{key}(\Omega_k)\}$. The index segment diagram is shown in figure 4. The $t_i$ in each row represents a keyword, $d_j$ in each column represents a

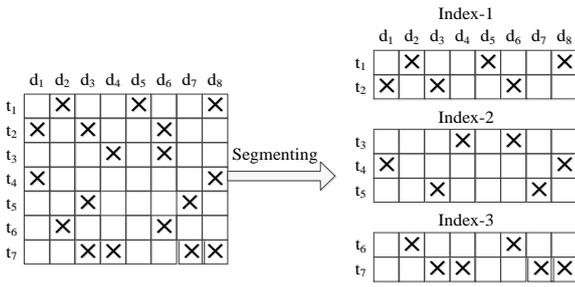document. The document set containing the keyword $t_i$ constitutes an inverted list.



**Figure 4.** Inverted index partitioning based on keywords

The obtained sub-indexes are then distributed to the distributed nodes. Each node is managed independently. Thus it realizes the distributed management of ciphertext inverted index.

## 3.3. Data retrieval

The segmentation and distribution of ciphertext inverted index realizes the distributed management of index files. In order to realize the complete multi-keyword distributed retrieval, three operations need to be completed. First, the multi-keyword search requests submitted by user are segmented and encapsulated into multiple sub-search requests corresponding to the sub-index. Second, the sub-retrieval request is distributed to the node that manages the corresponding sub-index to realize the distributed retrieval

of data, and the result set is handed over to the main node for relevancy processing and then returned to the user. Third, after the user receiving the retrieval result, it decrypts the result set and obtains the plaintext data. The user retrieval process is shown in figure 5.
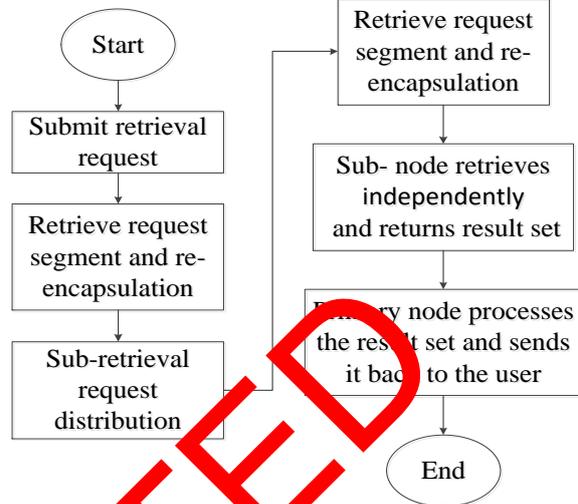


**Figure 5.** User retrieval process

Definition 3. User retrieval request $Q = \{td_1, td_2, \cdots, td_i\}$, where $td_i$ is the i-th ciphertext retrieval keyword in the retrieval request, and $td_i = En_{key}(q_i)$, $q_i$ is the retrieval keyword.

---

**Algorithm 3**. Multi-keyword search request partitioning algorithm.

Input. The user retrieving request Q.

Output. Sub-retrieving request set $\{Q_1, \cdots, Q_k\}$.

1) Segment retrieval request. The retrieval request submitted by the user is decomposed into several independent ciphertext keywords $td$ according to the order of the search words.

2) Keyword matching. According to the keywords segmentation scheme in the dictionary, the sub-index file $\Omega_j (1 \le j \le k)$ is calculated and determined. The search phrases belonging to the same sub-index are combined into the search word set.

3) Encapsulating the subset of retrieval keywords. After matching the sub-indexes for each retrieval keyword, each retrieval word set is then encapsulated into sub-retrieval request $Q_i = \{td_{i1}, td_{i2}, \cdots, td_{im}\}$ to obtain sub-retrieval request set $\{Q_1, \cdots, Q_k\}$. Where $Q_i$ is the sub-retrieval request corresponding to the sub-index $\Omega_j$. $td_{im}$ represents the m-th ciphertext

retrieval term in sub-retrieval request $Q_i$ .

4) Distributing retrieval requests. After the retrieval request is segmented, the primary node distributes each sub-retrieval request independently to the distributed node of the corresponding sub-index.

Definition 4. $R_i$ is the retrieval result set of the i-th node, and $d_{im}$ is the m-th retrieval document in the result set of the i-th node.

**Algorithm 4**. Distributed retrieval algorithm.

Input. Sub-retrieval request $Q_i$ .

Output. Retrieval result set R.

1) When the primary node Broker finishes dispatching sub-retrieval requests, each distributed node independently retrieves the sub-index files managed by itself to obtain the retrieval result set.

$$R_i = Search(Q_i, En(\Omega_i)); 1 \le i \le k$$

2) After each distributed node gets the retrieval result, the retrieval result $R_i = \{d_{i1}, d_{i2}, \cdots, d_{im}\}$ is sent back to the master node. After the primary node receives the retrieval result set $\{R_1, \cdots, R_k\}$ of each node, it takes the intersection operation of the result set to obtain the final retrieval result set R of user's retrieval request Q.

3) Return the retrieval results. After the primary node obtains the final retrieval result set, the result set is sent back to the client.

After receiving the retrieved ciphertext retrieval result set data, the user adopts the key to decrypt the data and get the plaintext data.

## 4. Experiments and analysis

This paper firstly compares the retrieval efficiency of SSE-1 scheme, CP-BSE [27], SSUR [28], with proposed index encryption scheme in this paper. The ciphertext retrieval efficiency of single index retrieval scheme and segmented multi-index retrieval scheme under single keyword are also compared. At the same time, the changes of the size of each index file with the increase of data volume are compared under the two schemes. Then, a multi-keyword distributed search scheme is implemented in the distributed model architecture. The retrieval efficiency of different retrieval keywords is compared and the feasibility of the distributed scheme proposed in this paper is verified.

In this paper, ciphertext inversion index is constructed and retrieved, based on Luence4.10. Data and indexes are encrypted symmetrically by using AES. System programming is implemented in Java. The computer is configured with the Intel Celeron 2955U, 1.40ghz CPU, a

4GR hard drive with 500GB and Windows 7. The experimental data comes from the text data of the major story websites in the network, including various types of classics and novel stories. The data format is TXT text.

The distributed architecture model includes a primary Broker node and several search-server child nodes. The Broker master node is responsible for the monitoring and scheduling of retrieval tasks, the segmentation and management of inverted indexes. The sub-nodes are computing nodes, which are responsible for data retrieval and index file management. It is configured with the Intel Celeron2955U, 1.40ghz CPU, RAM 4GB, hard drive 500GB and Windows 7.

The retrieval efficiency of ciphertext data is a very important performance index in the searchable encryption scheme. First of all, this paper compares the three index schemes and the proposed index encryption scheme in the case of single index file and single keyword respectively for different number of document sets. The experimental data are shown in figure 6.
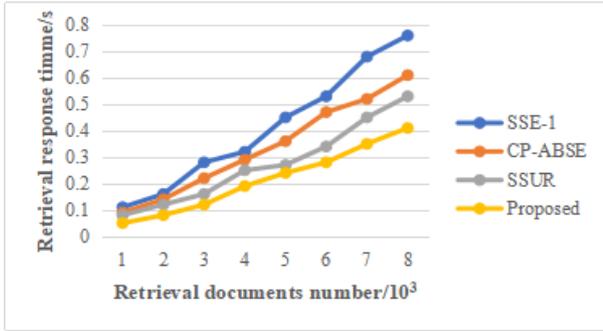
**Figure 6.** Comparison of retrieval efficiency between different index encryption methods

It can be seen from figure 6 that:

1) the retrieval time of encryption scheme in this paper is less than that of SSE-1 scheme. Because the inverted index pointer in SSE-1 scheme also needs to be encrypted. In the retrieval process, the previous node must be decrypted before accessing the content of the next node. However, this scheme only encrypts the key words and document information in the inverted index, and does not encrypt the pointer and structure of the whole index, which can achieve more efficient retrieval.

2) With the increase of documents number, the index encryption scheme in this paper has more obvious advantages in terms of retrieval efficiency. Because the increase in the number of documents makes more documents related to the search terms, and the inverted chain table in the inverted index is longer, which leads to the increase in the decryption calculation for Pointer in SSE-1 scheme, it reduces the retrieval efficiency.

Comparing the document retrieval efficiency of document sets with different data amount under the proposed index encryption scheme. The experimental results are shown in figure 7.
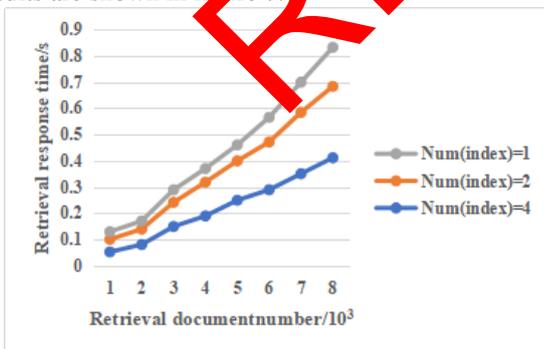


**Figure 7.** Comparison of single keyword retrieval efficiency.

Figure 7 shows that the retrieval time consumption increases with the increase of the number of documents. Because the increasing of documents number leads to the increase of index files size, resulting in the increase in the time consumption of index loading and retrieval. In addition, the multi-index retrieval scheme based on index segment is more efficient, and with the increase of the document data number and segment number, the advantages of retrieval are more obvious. Because in the retrieval process, the single index scheme needs to load the complete index file for calculation. However, the index segment-based retrieval scheme only needs to load the sub-indexes corresponding to the search words accurately in memory, which reduces the time consumption of index data loading and calculation, thus achieving high efficiency of retrieval.

By constructing ciphertext indexes for documents with different data volumes, the changing trend of the data size versus the size of index files during the retrieval operation is analyzed under different index schemes. The experimental results are shown in figure 8.
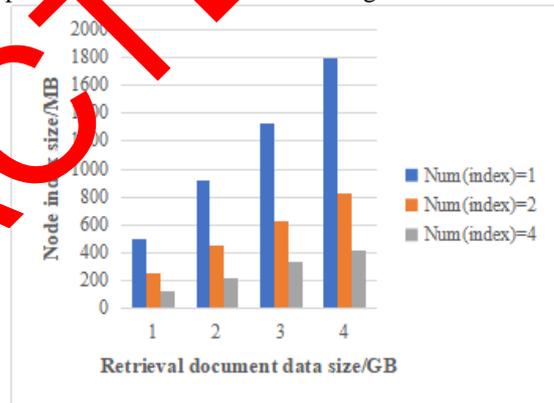


**Figure 8.** Index file size comparison during retrieval.

It can be found from figure 8 that in the single keyword retrieval and single index scheme, the increase of data volume leads to the segment increase of index file size. However, in the index-based segment scheme, with the increase of segment number, the size of each sub-index file increases slowly. Under the same computing resources, the sub-index retrieval scheme based on index segment can effectively improve the utilization rate of single machine resources and expand the scale of data retrieval in the retrieval process, so as to better adapt to the big data retrieval scenario.

Figure 9 depicts the retrieval time consumption of multiple keywords in a distributed retrieval system with different number of nodes.
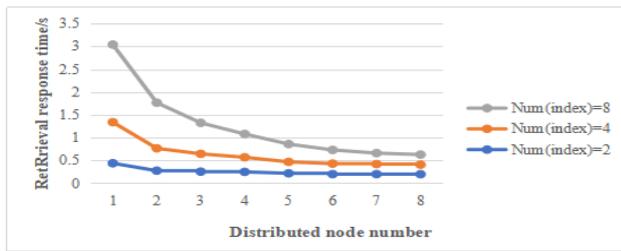
**Figure 9.** Comparison of distributed multi-keyword retrieval efficiency under different node number

Through the analysis of the comparison results in figure 9, it can be concluded that:

1) For ciphertext retrieval with multiple keywords, the distributed retrieval model architecture realized by index segmentation can effectively improve the retrieval efficiency. Because with the increase of compute nodes and index segment number, the retrieval keywords are more dispersed and fewer keywords are retrieved on each node, thus improving the retrieval efficiency in parallel retrieval mode.

2) with the increase of keywords in retrieval requests, the retrieval efficiency decreases gradually. This because the increase of keywords leads to the increase of the computation of each node retrieval. At the same time, the computation of the primary node Broker in segmenting the retrieval request and merging the retrieval result increases, which leads to the decrease of the whole retrieval process efficiency.

3) With the increase of segment nodes, the retrieval efficiency of different search terms tends to be stable. This because with the increase of nodes, the calculation amount of each node decreases. Under the condition that the retrieval time consumption of each node does not dominate, the time consumption is mainly concentrated in the main node Broker's index segmenting and sub-index matching processing stage, so the overall retrieval time tends to be stable.

Figure 10 depicts the comparison of the multiple keywords retrieval efficiency based on the literature [27] scheme and the proposed scheme under the condition that the number of search words and the number of distributed nodes are four under the set of different document quantities.
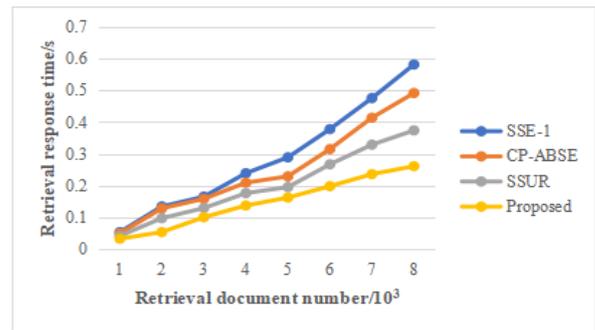


**Figure 10.** Comparison of retrieval efficiency between two distributed multi-keyword retrieval schemes

Through the comparison result in figure 10, it can be seen that the multi-keyword retrieval scheme in this paper has more advantage than the other three schemes, and the advantage gradually become obvious with the increase of the documents number. Because several search terms in the retrieval scheme are independent of each other and the generation of the retrieval trap does not depend on the index file. While the retrieval scheme in CP-ABSE needs to couple multiple search terms together to generate polynomials in the process of creating the retrieval trap, and realize the multiple keywords retrieval according to the polynomial matrix. Both index construction and retrieval process require a lot of polynomial matrix calculation, which seriously affects the retrieval efficiency. In addition, the calculation of retrieval trap polynomial matrix relies on the dictionary file of inverted index. The increase of the documents number will lead to the increase of the dictionary file size in the index, which will make the calculation of retrieval trap more complicated and reduce the retrieval efficiency. Therefore, it is not suitable for the application scenarios with rapidly increasing data volume.

The experimental results show that the distributed retrieval scheme based on index segmentation can effectively solve the problem of single machine resource limitation and improve the efficiency of ciphertext data retrieval with multiple keywords. With the continuous growth of data scale, the distributed model architecture has good scalability and can better adapt to the application scenarios with the rapid growth of data volume.

## 5. Conclusion

Searchable encryption is an important ways to realize ciphertext data retrieval. In this paper, based on the inversion index segmentation of keywords, a segmentation-based ciphertext inversion index is

proposed. By index segmentation and accurate loading of sub-index files, the limitation of stand-alone resources in the single-server model architecture is solved to improve the efficient utilization of stand-alone memory resources and computing. In view of the multi-keyword retrieval mode, a multi-keyword parallel ciphertext retrieval scheme is proposed in the distributed environment, which realizes the distributed management of index files and the parallel retrieval of keywords. It improves the retrieval efficiency of multi-keyword, and makes the ciphertext retrieval model architecture have the characteristics of scalability and flexibility. The experimental results show that the distributed ciphertext retrieval system can realize efficient retrieval for ciphertext data, and the system is effective and feasible.

## Acknowledgements.

## References

[1] Shoulin Yin, Hang Li and Jie Liu. A New Provable Secure Certificateless Aggregate Signcryption Scheme [J]. Journal of Information Hiding and Multimedia Signal Processing, Vol. 7, No. 6, pp. 1274-1281, November 2016.

[2] Shoulin Yin and Jie Liu. A K-means Approach for Map-Reduce Model and Social Network Privacy Protection[J]. Journal of Information Hiding and Multimedia Signal Processing, Vol. 7, No. 6, pp. 1215-1221, November 2016.

[3] Li P, Chen Z, Yang L T, et al. An Incremental Deep Convolutional Computation Model for Feature Learning on Industrial Big Data[J]. IEEE Transactions on Industrial Informatics, vol.15, no.3, pp. 1341-1349, 2019.

[4] Qingchen Zhang, Changchuan Bai, Laurence T. Yang, Zhikui Chen, Peng Li and Hang Yu, A unified smart Chinese medicine framework for healthcare and medical services[J]. IEEE/ACM Transactions on Computational Biology and Bioinformatics, 2019, DOI: 10.1109/TCBB.2019.2914447.

[5] Ling Teng, Hang Li and Shoulin Yin. IM-MobiShare: An Improved Privacy Preserving Scheme Based on Asymmetric Encryption and Bloom Filter for Users Location Sharing in Social Network [J]. Vol. 30, No. 3, pp. 59-71, Journal of Computers (Taiwan). 2019.

[6] Fu Z, Ren K, Shu J, et al. Enabling Personalized Search over Encrypted Outsourced Data with Efficiency Improvement[J]. IEEE Transactions on Parallel and Distributed Systems, 2015:1-1.

[7] Qiang Z, Qin L, Guojun W. PRMS: A Personalized Mobile Search Over Encrypted Outsourced Data[J]. IEEE Access, 2018, 6:31541-31552.

[8] Lina Zou, Xueying Wang, Shoulin Yin. A Data Sorting and Searching Scheme Based on Distributed Asymmetric Searchable Encryption[J]. International Journal of Network Security, Vol. 20, No. 3, pp. 502-508, 2018.

[9] Bing Wang, Wei Song, Wenjing Lou, et al. Inverted index based multi-keyword public-key searchable encryption with strong privacy guarantee[C]// IEEE INFOCOM 2015 - IEEE Conference on Computer Communications. IEEE, 2015.

[10] Zhang R, Xue R, Yu T, et al. Dynamic and Efficient Private Keyword Search over Inverted Index--Based Encrypted Data[J]. ACM Transactions on Internet Technology, 2016, 16(3):1-20.

[11] SONG D X, WAGNER D, PERRIG A. Practical techniques for searches on encrypted data[C]// Proceedings of the 2000 IEEE Symposium on Security and Privacy. Piscataway: IEEE, 2000: 44

[12] Curtmola R, Garay J, Kamara S, et al. Searchable symmetric encryption: Improved definitions and efficient constructions[J]. Journal of Computer Security, 2011, 19(5):895-934.

[13] Naveed M, Prabhakaran M, Gunter C A. Dynamic Searchable Encryption via Blind Storage[C]// 2014 IEEE Symposium on Security and Privacy (SP). IEEE, 2014.

[14] Hongwei Li, Yi Yang, Yuanshun Dai, et al. Achieving Secure and Efficient Dynamic Searchable Symmetric Encryption over Medical Cloud Data[J]. IEEE Transactions on Cloud Computing, 2017, PP(99):1-1.

[15] Dwi Satya Palupi, Eduardus Tandelilin, Arief Hermanto, et al. Fluctuation of LQ45 index and BCA stock price at Indonesian Stock Echange IDX[J]. International Journal of Engineering Research and Applications, 2017.

[16] Qian Jiangbo, Huang Zhipeng, Zhu Qiang, et al. Hamming Metric Multi-Granularity Locality-Sensitive Bloom Filter[J]. IEEE/ACM Transactions on Networking:1-14, 2018.

[17] Yi Mu, Yupu Hu, Leyou Zhang, et al. Fully secure hierarchical inner product encryption for privacy preserving keyword searching in cloud[J]. International Journal of High Performance Computing & Networking, 2018, 11(1):45.

[18] Guofeng Wang, Chuanyi Liu, Yingfei Dong, et al. IDCrypt: A Multi-User Searchable Symmetric Encryption Scheme for Cloud Applications[J]. IEEE Access, 2018, 6:2908-2921.

[19] Xia Z, Wang X, Sun X, et al. A Secure and Dynamic Multi-keyword Ranked Search Scheme over Encrypted Cloud Data[J]. IEEE Transactions on Parallel and Distributed Systems, 2015:1-1.

[20] Ballard L, Kamara S, Monrose F. Achieving Efficient Conjunctive Keyword Searches over Encrypted Data[C]// Information and Communications Security, 7th International Conference, ICICS 2005, Beijing, China, December 10-13, 2005, Proceedings. Springer-Verlag, 2005.

[21] Fu Z, Wu X, Guan C, et al. Toward Efficient Multi-Keyword Fuzzy Search Over Encrypted Outsourced Data With Accuracy Improvement[J]. IEEE Transactions on Information Forensics and Security, 2016, 11(12):2706-2716.

[22] Guo Z, Zhang H, Sun C, et al. Secure Multi-Keyword Ranked Search over Encrypted Cloud Data for Multiple Data Owners[J]. Journal of Systems and Software, 2017, 137:380-395.

[23] Kim I T, Quan T H, Duc L V, et al. An Efficient Searchable Encryption Scheme in the Multi-user Environment[C]// International Conference on Green and Human Information Technology. Springer, Singapore, 2018.

[24] Teng Lin, Li Hang, Liu Jie, Yin Shoulin. An efficient and secure Cipher-Text retrieval scheme based on mixed homomorphic encryption and Multi-Attribute Sorting Method Under Cloud Environment[J]. International Journal of Network Security, v 20, n 5, p 872-878, September 1, 2018.

[25] Song N, Zou L, Liu Y. Research on the Inverted Index Based on Compression and Perception[C]// International Conference on Trustworthy Computing and Services. 2014.

[26] Tolson B, Matott L S, Gaffoor T A, et al. Parallel and Preemptable Dynamically Dimensioned Search Algorithms for Single and Multi-objective Optimization in Water Resources[C]// 2015.

[27] Hui Yin, Jixin Zhang, Yingjie Xiong, et al. CP-ABSE: A Ciphertext-Policy Attribute-based Searchable Encryption Scheme[J]. IEEE Access, 2019, 7(99):1-1.

[28] Fateh Boucenna, Omar Nouali, Samir Kechid, et al. Secure Inverted Index Based Search over Encrypted Cloud Data with User Access Rights Management[J]. Journal of Computer Science and Technology, 2019.