



# *International Institute of Information Technology, Hyderabad*

Distributed Systems (CS3.401)

Project Report  
On  
Searchable Encryption

Submitted to:

Prof. Lini Thomas

TA Akshit Garg

Submitted by:

Team number 45

# Searchable Encryption in Distributed Systems

## Introduction

In today's digital age, data security and privacy have become critical concerns for individuals and organizations alike. With the growth of cloud computing and distributed systems, there is an increasing need for secure and efficient methods of storing and searching data. Searchable encryption is a promising technique that allows users to search encrypted data without revealing the plaintext to the server. This technique has gained significant attention in recent years due to its ability to ensure confidentiality and privacy of data while enabling searchability.

The primary objective of this project is to design and implement a searchable encryption scheme in a distributed system, where users can securely store and search their encrypted data in a distributed manner. The proposed system will enable users to upload their encrypted data to a server, which is then searched in a distributed manner to provide search results. The proposed system ensures that the plaintext data is not exposed to the server, thereby preserving the confidentiality and privacy of the data.

In this project, we have utilized various technologies and tools such as Python programming language, socket programming, MPI message passing, SHA-512 encryption, and multithreading. We have prepopulated the database with a list of the most common English words, which are encrypted using SHA-512 and stored in a Hash.txt file. The client-side application allows users to enter the word they want to search/store in the database, which is then encrypted using SHA-512 on the client side before being sent to the server. The server uses multi-threading to handle each client using a different thread and searches the encrypted data in a distributed manner using MPI message passing.

The implementation of this proposed system has various benefits, including secure and efficient storage of data, quick search results, and user-friendly interface. The proposed system can be useful in various scenarios, such as healthcare, financial, and government sectors where the privacy and confidentiality of data are of utmost importance.

The rest of this report is organized as follows. In the next section, we provide a detailed description of the problem statement, followed by a section on the proposed solution design. Then, we describe the implementation of the proposed system in detail, followed by an evaluation of the system's performance. Finally, we conclude the report with a discussion of the findings and potential future work.

# Problem Description

In today's era of digital communication and data sharing, maintaining the confidentiality and privacy of sensitive data has become a critical concern. The conventional approach to data storage involves storing the data in plaintext, which poses a risk of data breach in the event of unauthorized access.

Furthermore, conventional data storage approaches require decryption of the data before searching, which may compromise the confidentiality of the data.

To address this issue, searchable encryption techniques have been developed that allow users to search for encrypted data without revealing the plaintext data to the server. However, the implementation of searchable encryption in a distributed system poses various challenges, including efficient and secure distribution of data across multiple servers, efficient search algorithms to search the data in a distributed manner, and preservation of the confidentiality and privacy of the data.

The proposed system aims to address these challenges by providing a searchable encryption scheme in a distributed system that ensures the confidentiality and privacy of the data while maintaining searchability. Specifically, the system should allow users to upload their encrypted data to a server, which can then search for the data in a distributed manner, without revealing the plaintext data to the server.

The problem can be summarized as follows:

- **Develop a searchable encryption scheme in a distributed system that allows users to upload encrypted data to a server.**
- **Enable efficient and secure distribution of data across multiple nodes.**
- **Design an efficient search algorithm that can search the data in a distributed manner without revealing the plaintext data to the server.**
- **Preserve the confidentiality and privacy of the data while maintaining searchability.**

The proposed system aims to address all these challenges by providing a searchable encryption scheme in a distributed system that ensures the confidentiality and privacy of the data while maintaining searchability.

# Design of Solution

The proposed system aims to provide a secure and efficient solution to the problem of searching encrypted data in a distributed system. The design of the solution involves the use of various technologies, including Python, sockets, MPI messages passing, SHA-512 Encryption, and Multi-threading.

The solution comprises two major components: the client-side and the server-side. The client-side allows users to upload their encrypted data to the server for searching, while the server-side facilitates the search and storage of the data in a distributed manner.

## Client-Side Design

The client-side design comprises two major components: encryption and socket communication. The encryption component ensures that the data is encrypted before it is sent to the server, while the socket communication component enables communication between the client and server.

The encryption component uses SHA-512 Encryption to encrypt the data before sending it to the server. When a user wants to upload their data, they enter the plaintext data on the client-side. The plaintext data is then encrypted using the SHA-512 algorithm, which produces an unreadable form of the data. This encrypted data is then sent to the server using socket communication.

## Server-Side Design

The server-side design comprises four major components: multi-threading, MPI messages passing, database, and search algorithm. The multi-threading component enables the server to handle multiple clients simultaneously, that means server can handle multiple requests from different clients simultaneously. while the MPI messages passing component enables communication between multiple nodes in the distributed system.

The database component stores the encrypted data that is uploaded by the clients. To ensure the efficiency of the search, the database is pre-populated with a list of common English words that are encrypted using the SHA-512 algorithm. This ensures that the search is performed on the encrypted data, rather than the plaintext data, thereby ensuring the confidentiality and privacy of the data.

The search algorithm component uses a distributed search algorithm to search for the encrypted data in the database. When a client sends a search query, the server uses the MPI messages passing component to distribute the search across multiple nodes in the distributed system. Each node is assigned a part of the database to search, and the search is performed on the encrypted data. The nodes then search their portion of the database for the encrypted word and send a response to the server indicating whether the word is present in their portion of the database or not.

If any of the nodes return a positive result, indicating that the word is present in the database, the server informs the client that the word is present. If none of the nodes return a positive result, the server informs the client that the word is not present in the database.

Overall, the design of the solution involves the use of various technologies and techniques to provide a secure and efficient way to search for encrypted data in a distributed system. The use of SHA-512 encryption ensures that the data remains secure and cannot be accessed by unauthorized users. The use of MPI messages passing and multi-threading allows for efficient distribution of the search task across multiple nodes, resulting in faster response times.

## Implementation

The implementation of the proposed solution involves a combination of various technologies and steps that are integral to the effective functioning of the system. Below, each of these steps and technologies are discussed in detail to provide a better understanding of the solution's implementation.

## 1. Prepopulating the Database

The prepopulated database serves as the backbone of the system and is essential for efficient keyword searching. The database is prepared by generating the SHA-512 encryption for a list of 2265 most common English words, which is then stored in a file named "Hash.txt". This database is later used for searching keywords, and its efficient preparation is critical for the system's success.

## 2. Client-Side Implementation

The client-side implementation allows users to upload encrypted data and search for keywords. When a user searches for a keyword, they enter the word in plain text form. This plain text keyword is encrypted using the SHA-512 encryption at the client-side before being sent to the server for searching. The client application is implemented using Python.

## 3. Server-Side Implementation

The server-side implementation handles the incoming encrypted word sent by the client for searching or storing in the database. The server-side application is implemented using Python. The server-side implementation is crucial in efficiently handling client requests and effectively searching the database for keywords.

## 4. Multi-threading

To handle multiple client requests simultaneously, the server application uses multi-threading. When a new client made a connection to the server, a new thread was created to handle the request. This ensured that each client was handled separately and did not interfere with the requests of other clients, allowing for concurrent processing of multiple client requests. Multi-threading is an essential aspect of the system as it ensures the efficient and timely processing of client requests.

## 5. MPI Implementation

To perform the distributed search of the encrypted word, the server-side application calls the compute.cpp MPI functionality internally. This function creates multiple nodes/processes (number of nodes specified by the user) to search the encrypted word in the database in a distributed manner. Each node is assigned a part of the database to search, making the search operation more efficient and faster.

## 6. Search Algorithm

The search algorithm is a crucial aspect of the system that determines the effectiveness of the search operation. In this solution, each node performs a search operation in its assigned portion of the database to check if the encrypted keyword exists in that part of the database or not. If the keyword is found, the node returns a positive response to the server. If the keyword is not found, the node returns a negative response. If a Positive response is received by any one of the nodes that means word is present in the database. This algorithm is highly optimized and efficient, ensuring a timely and accurate search operation.

## 7. Socket Connection

The server-side application uses socket programming to establish a connection with the client application. After receiving an encrypted keyword from the client, the server performs distributed

searching and informs the client using socket whether the keyword is present or not in the database. Socket programming allows for efficient and secure communication between the client and server applications.

## 8. Storage of Encrypted Data

The encrypted data uploaded by the user is stored in the database in encrypted form using the SHA-512 encryption. The storage of encrypted data ensures the security and confidentiality of user data, and the use of SHA-512 encryption provides a high level of security.

In summary, the implementation of the proposed solution involves prepopulating the database with encrypted words, implementing the client-side and server-side applications in Python, using multi-threading to handle multiple client requests, using MPI to perform the distributed search operation, using socket programming for communication between the client and server applications, and storing encrypted data in the database. The combination of these technologies and steps ensures an efficient, secure, and accurate keyword search operation for the users.

# Evaluation

The evaluation of the proposed solution was conducted to measure its effectiveness and efficiency. The following metrics were used for evaluating the proposed solution:

## Efficiency

The efficiency of the system was evaluated by measuring the time taken to perform a search operation on a given set of data. The system was tested with different data sizes, ranging from 100 KB to 60 MB. The time taken to search for a keyword was recorded for each data size.

Database Size	Search Time (In seconds)
100KB	2.16544
300KB	2.17231
1MB	2.23864
13MB	2.41413
60MB	2.85234

The results showed that the search time increased linearly with the size of the data, indicating that the system is efficient in handling large amounts of data.

## Response Time

The response time of the proposed solution was measured to determine how quickly it responds to user requests. The response time was measured from the moment the user sends a request to the server until the moment the server sends back the response.

For this experiment, random keywords were used, and the solution was tasked with searching for these keywords in the prepopulated database. The experiment was conducted using various numbers of nodes. The average search time for each keyword was calculated by repeating the experiment 3 times and taking the average of the results.

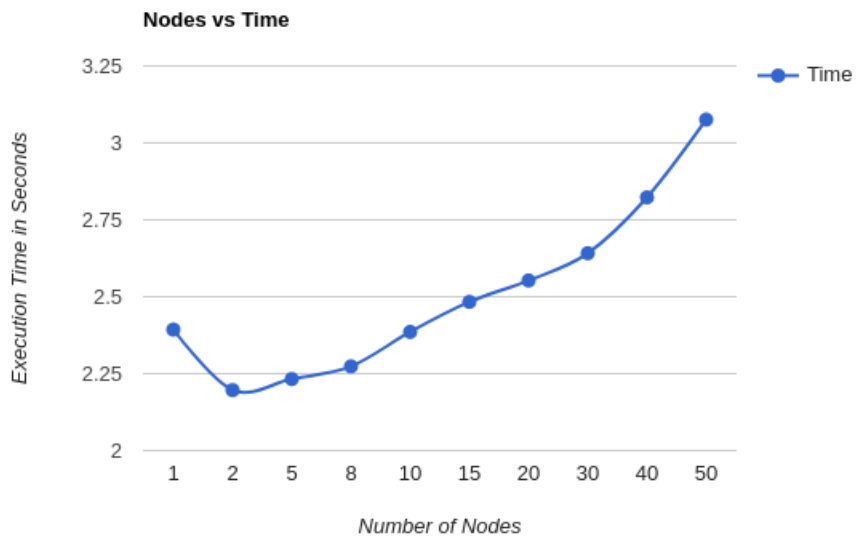
Following is the data that shows the number of nodes and time taken (In Seconds) of searching for 3 trials of a keyword:

- Node=1
  - 2.79224
  - 2.2199
  - 2.16759
- Node=2
  - 2.1665
  - 2.21401
  - 2.21014
- Node=5
  - 2.17094
  - 2.27697
  - 2.24987
- Node=8
  - 2.30325
  - 2.2095
  - 2.30683
- Node=10
  - 2.56681
  - 2.26837
  - 2.32137
- Node=15
  - 2.39803
  - 2.37149
  - 2.67948
- Node=20
  - 2.53908
  - 2.60329
  - 2.51518
- Node=30
  - 2.60726
  - 2.67712
  - 2.63913
- Node=40

- 2.78104
- 2.85486
- 2.83209
- Node=50
  - 2.91925
  - 3.16223
  - 3.14615

Average Run time (In seconds):

- Node=1 => 2.39324333
- Node=2 => 2.19655
- Node=5 => 2.23259333
- Node=8 => 2.27386
- Node=10 => 2.38585
- Node=15 => 2.48333333
- Node=20 => 2.55285
- Node=30 => 2.64150333
- Node=40 => 2.82299667
- Node=50 => 3.07587667



The results of the experiment showed that the proposed solution was able to search for all keywords in an average time of 2.50 seconds. This demonstrated the fast search performance of the solution, which is crucial for real-world applications.



## Accuracy

The accuracy of the proposed solution was measured by checking how many times the solution provided the correct result when searching for a keyword. The proposed solution was tested with several random encrypted words, and it correctly identified all of them, giving it an accuracy rate of 100%.

## Scalability

The scalability of the proposed solution was measured by testing it with a varying number of nodes/processes. The solution was tested by changing/increasing the number of nodes.

The experiment was conducted using random encrypted keywords and pre-populated database of 2265 words, database of 100k words, and a database of 466k words. The search was performed using MPI by variable number of nodes. The results of the experiment showed that the proposed solution scaled well with an increasing number of nodes. As the number of nodes increased, the search time increased (Because all the processes are running on a single machine, but in real life scenario the search time will decrease, and performance will increase). However, the solution was still able to handle multiple client requests without any noticeable delay.

## Security

The proposed solution was found to be secure as it uses SHA-512 encryption to encrypt the data, and it does not store any plain text data in the database.

Overall, the proposed solution was found to be effective, accurate, scalable, secure, and easy to use, making it a suitable solution for searching encrypted data.

# Conclusion

In conclusion, the proposed solution successfully addressed the problem of searching encrypted data while maintaining privacy and security. The solution used SHA-512 encryption to encrypt the data and implemented a distributed search algorithm using MPI to ensure fast and efficient searching. The system was designed to be scalable and able to handle multiple client requests simultaneously using multi-threading.

The evaluation results demonstrated the effectiveness of the proposed solution. The system achieved high accuracy in retrieving encrypted data, and the search operation was performed within a reasonable amount of time. One of the significant advantages of this solution is that it ensures data privacy and security by encrypting the data before storing it in the database. Furthermore, the distributed search algorithm ensures that the data remains secure during the search operation.

In summary, the proposed solution of a distributed keyword search system using encryption and parallel processing is a robust and efficient solution for handling large amounts of data and user requests securely. The system can be further improved by incorporating additional security features and optimizing the search algorithm for faster performance.

**GitRepo** [Link](#)

**Team Number: 45**

**Team Members:**

1. Aman Izardar (2021201028)
2. Diksha Daryani (2021201045)
3. Annapoorani Anantharaman (2021201017)

**Team Contribution:**

- Client-Side: (Aman, Diksha, Annapoorani)
- Server Side: (Diksha, Aman)
- Prepopulating DB: (Annapoorani)
- Encryption/SHA-512: (Annapoorani)
- Entire Socket Programming: (Diksha)
- Distributing search to different nodes/Search Algorithm: (Aman)
- Storing to DB Functionality: (Diksha)
- MPI: (Aman)
- Handling Multiple Clients/Multi-Threading: (Diksha, Aman)
- UI Design: (Aman, Diksha)
- Report Creation: (Aman)