



# Experimental Review of the IKK Query Recovery Attack: Assumptions, Recovery Rate and Improvements

Ruben Groot Roessink<sup>(✉)</sup>, Andreas Peter, and Florian Hahn

University of Twente, Enschede, The Netherlands  
r.grootroessink@alumnus.utwente.nl, {a.peter,f.w.hahn}@utwente.nl

**Abstract.** In light of more data than ever being stored using cloud services and the request by the public for secure, privacy-enhanced, and easy-to-use systems, Searchable Encryption schemes were introduced. These schemes enable privacy-enhanced search among encrypted documents yet disclose (encrypted) queries and responses. The first query recovery attack, the IKK attack, uses the disclosed information to (partly) recover what plaintext words the client searched for. This can also leak information on the plaintext contents of the encrypted documents. Under specific assumptions, the IKK attack has been shown to potentially cause serious harm to the security of Searchable Encryption schemes.

We empirically review the IKK query recovery attack to improve the understanding of its feasibility and potential security damage. In order to do so, we vary the assumed query distribution, which is shown to have a severe (negative) impact on the accuracy of the attack, and the input parameters of the IKK attack to find a correlation between these parameters and the accuracy of the IKK attack. Furthermore, we show that the recovery rate of the attack can be increased up to 10% points, while decreasing the variance of the recovery rate up to 78% points by combining the results of multiple attack runs. We also show that the including deterministic components in the probabilistic IKK attack can increase the recovery rate up to 21% points and decrease its variance up to 57% points.

**Keywords:** Searchable Encryption · IKK · Query recovery

## 1 Introduction

The use of currently available encryption schemes allows users to securely upload and retrieve documents anywhere in the world using cloud services. A user encrypts a set of documents and sends these encrypted documents to a server for storage. The server can return documents upon request by the user, which the user can decrypt to obtain the original documents, while the server is not capable of reading the contents of the documents. A downside of using encryption schemes is that they, in general, limit the functionalities of the cloud service. One such functionality is the possibility to search for word occurrences among documents. To overcome this loss in functionality, while also taking into regard data

confidentiality, Song et al. [15] introduced the notion of Searchable Encryption (SE).

In general, SE schemes provide a client with a way to search for the occurrence of a certain (plaintext) word, a *keyword*, among a set of encrypted documents, while neither the client nor the server has to decrypt all documents which the client wants to search among. A client generates a search token, a *query*, which it sends to a server hosting a set of (encrypted) documents. The server uses the query to find a subset of the encrypted documents corresponding with the search token and returns this subset to the client.

Nearly all of these SE schemes leak at least some information, usually in the form of *data access patterns* [8, 9, 15]. This means that an adversary can observe the issued queries from the client and the document identifiers of documents corresponding to said queries in the response by the server. This allows the client to make a connection between the queries and the corresponding documents. Some schemes were proposed [8, 10] that hide these access patterns. However, these schemes are quite inefficient as they require an extensive number of computations after each query. Other schemes propose to obfuscate the access patterns which can both lead to inconsistencies in the search results (false positives or false negatives) and an increase in communication and storage overhead [7].

Islam et al. [12] elaborate on the implications of the leakage of access patterns by proposing the first *query recovery attack*, dubbed IKK attack in subsequent research after the first initials of the authors of the paper. Their attack is a statistical attack which tries to map queries to their corresponding real-world keywords. This mapping process is dubbed *query recovery*. A correctly ‘recovered’ query tells an adversary what the client searched for and possibly even tells something about the contents of (encrypted) documents stored on the server. In their attack Islam et al. use the relative co-occurrence counts of queries, which denotes the number of documents a certain number of queries occur in together, relative to the total number of documents. These counts can be calculated from leaked access patterns. The IKK attack also assumes the adversary has access to (a close approximation of) the co-occurrence counts of the plaintext (key)words in these documents, dubbed *background knowledge*. Islam et al. show that a large percentage of queries is recoverable, expressed as the (*query*) *recovery rate*, if the adversary has perfect background knowledge, meaning that the co-occurrence counts of keywords exactly match the co-occurrence counts of their corresponding queries. They also briefly show that the recovery rate drops significantly in simulations with non-perfect background knowledge.

We revisit the IKK attack and empirically evaluate assumptions Islam et al. make in their proposal of the IKK attack. Additionally, we research correlations between certain parameters and the accuracy/recovery rate of the IKK attack and propose improvements to the attack that increase the recovery rate of the attack.

## Our Contributions

- We show that assumptions on the (Zipfian) distribution of queries/natural search behavior Islam et al. made positively influences the accuracy of the IKK attack.
- We show that there is a correlation between input parameters of the IKK attack and the accuracy of attack runs, independent of the (email) dataset used in the targeted Searchable Encryption scheme, potentially allowing an adversary to reuse the same values of parameters across different datasets.
- We show that the accuracy of the IKK attack can be increased significantly when combining multiple runs using a majority voting scheme as the median recovery rate is increased up to 10% points, whereas the variance of the recovery rate is decreased up to 78% points.
- We show that a more deterministic approach to select new states in the IKK attack, inspired by the Count attack [6], increases the accuracy of the attack, while decreasing the average number of states visited. The median recovery rate is increased up to 21% points and the variance of recovery rates is decreased up to 57% points, while the average visited number of states is decreased by 28%.

## 2 The IKK Query Recovery Attack

### 2.1 Searchable Encryption (SE)

The first SE scheme was proposed by Song et al. [15] to provide a client with a way to search for the occurrence of a plaintext word among a set of encrypted documents stored on a server without an adversary being able to learn the (plaintext) contents of these documents.

The server stores a set of encrypted documents, for example email files. The client wants to retrieve emails that contain information on an upcoming merger and thus requests all emails that contain the (plaintext) word *merger*. It does so by generating a so-called query using a keyed trapdoor function,  $Trapdoor_{sk}(merger)$ , for example, a keyed hash function. Only users with key  $sk$  can generate valid queries. More formally, a user knowing key  $sk$  can generate query token  $q_i$  for a keyword  $k_i$ . We assume, just like Islam et al. [12], that queries are deterministic, i.e.  $Trapdoor_{sk}(k_i) = Trapdoor_{sk}(k_j)$  if  $k_i = k_j$ .

In order to retrieve the corresponding documents, the client sends the query to the server, which on its turn, performs a matching algorithm. Most proposed SE schemes either encrypt every single word in a document (*In-place SE* [6]) to encrypt a document or encrypt every document using a traditional encryption scheme, such as AES, while also generating an encrypted inverted index of documents and trapdoors (*Encrypted-index SE* [6]) to allow the server to perform the search query. No matter the SE scheme, the server returns all the documents that match the search query, which can be decrypted by the client.

## 2.2 Access Pattern Disclosure

Just like Cash et al. [6], we deem the server the most likely adversary as it has access to the most information. Nonetheless, any adversary with access to the communication channels is able to connect a query to the identifiers of the documents that were returned and thus is able to see which documents were *accessed* upon the query of the client. This has been dubbed (*data*) *access pattern disclosure* in the literature [12]. Almost all SE schemes, except schemes that re-encrypt the documents or the encrypted index stored on the server after each query [8, 10], disclose access patterns of particular queries, and thus each query gives an adversary more information on which queries are connected to which documents. Some schemes propose to obfuscate access patterns which can lead to inconsistencies in the search results (false positives or false negatives) or an increase in communication and storage overhead [7]. In our research, we assume no such inconsistencies were added to the search results.

Although we note that different SE schemes leak different levels of information, we only research the disclosure of access patterns and assume that the adversary is able to get  $\langle \text{query}, \text{response} \rangle$  pairs, where *response* is a list of documents that matched the issued query. The leaked  $\langle \text{query}, \text{response} \rangle$  pairs allow the adversary to construct an inverted index from queries/trapdoors and documents. Each cell in the matrix contains a 1 if the document matched the query, i.e. the plaintext keyword occurs at least once in said document, or a 0 if not. An example of an (observed) query inverted index is shown in Table 1.

**Table 1.** Example of a query inverted index

Query	Documents		
	<i>Doc</i> <sub>1</sub>	<i>Doc</i> <sub>2</sub>	<i>Doc</i> <sub>3</sub>
$q_1 = \text{Trapdoor}_{sk}(\text{merger})$	1	1	0
$q_2 = \text{Trapdoor}_{sk}(\text{corporate})$	1	1	0
$q_3 = \text{Trapdoor}_{sk}(\text{report})$	0	1	1

## 2.3 Statistical Processing

The IKK attack is grounded on the assumption that some words are more likely to occur together in any piece of natural language than others. Islam et al. [12] give the example of the words *New*, *York* and *Yankees*, where the words *New* and *York* are more likely to occur together than *New* and *Yankees* or *York* and *Yankees* because they are also used to refer to the city and the state and not only the baseball team. Islam et al. propose a model where the co-occurrence counts of 2 queries are used to recover which plaintext words correspond to which queries. The authors use a so-called co-occurrence matrix to express all co-occurrence counts of the queries in an attack run, as an co-occurrence matrix

lists the (relative) co-occurrence count for each of the queries with all other queries and with itself. The probability that two words appear together in a given document is expressed using the following formula by Islam et al.:

$$\beta = \frac{R_{Q_s} \cdot R_{Q_t}}{n} \quad (1)$$

In this formula,  $Q_s$  and  $Q_t$  are two queries, and  $R_{Q_x}$  denotes a vector with ones and zeros indicating whether the word corresponding to the query occurs at least once in the document corresponding with the place in the vector (1) or not (0). The co-occurrence count is calculated by taking the *dot* product of  $R_{Q_s}$  and  $R_{Q_t}$ . To get the relative co-occurrence count this value is simply divided by  $n$ , which denotes the total number of documents in the dataset. A query co-occurrence matrix simply lists all the co-occurrences of queries, is symmetric by nature and is easily generated using an inverted index as every row in the index, corresponding to query  $Q_x$ , is already represented as vector  $R_{Q_x}$ . The relative co-occurrence count of a query with itself is the total number of documents said query occurs in divided by  $n$ . The example inverted index in Table 1 gives us the co-occurrence matrix in Table 2.

**Table 2.** Example of a query co-occurrence matrix

Queries	$q_1$	$q_2$	$q_3$
$q_1$	0.67	0.67	0.33
$q_2$	0.67	0.67	0.33
$q_3$	0.33	0.33	0.67

**Table 3.** Example of a (perfect) background knowledge co-occurrence matrix

Keywords	<i>Corporate</i>	<i>Merger</i>	<i>Report</i>
<i>corporate</i>	0.67	0.67	0.33
<i>merger</i>	0.67	0.67	0.33
<i>report</i>	0.33	0.33	0.67

## 2.4 Background Knowledge Assumptions

In most of their simulations Islam et al. [12] assume the adversary has perfect background knowledge of the co-occurrence counts of plaintext words in the documents stored encrypted on the server. They mention that it is difficult, if not impossible, to obtain perfect background knowledge and briefly experiment on the accuracy of their attack in simulations with non-perfect background knowledge by adding various degrees of Gaussian noise to a co-occurrence matrix corresponding to perfect background knowledge.

Cash et al. [6] further research the effect of non-perfect background knowledge, but instead of adding various degrees of Gaussian noise to a perfect representation of the background knowledge (of the adversary) the authors assume the adversary (server) has access to a *fraction* of the plaintext documents and thus the adversary is capable of calculating both inverted indices and co-occurrence matrices from the documents it knows. The authors report that the IKK attack performs quite poorly if the background knowledge is made up of less than 99% of the documents. An example of a background knowledge co-occurrence matrix is shown in Table 3.

## 2.5 Simulated Annealing

Islam et al. [12] use two algorithms to recover queries from a query co-occurrence matrix and a background knowledge co-occurrence matrix: Their *Optimizer* (Algorithms 1 and 2) algorithm assigns a random 1-to-1 mapping for each query to a random keyword in the background co-occurrence matrix as the *initial state* variable. The mapping corresponds to a mapping between the query co-occurrence matrix and a subset of the background knowledge co-occurrence matrix which is equal in its dimensions to the query co-occurrence matrix. Each cell in the query co-occurrence matrix is therefore mapped to a single cell in the background knowledge co-occurrence matrix. The *initial state* is given as input to their ANNEAL (Algorithms 3 and 4) algorithm. This algorithm is a *Simulated Annealing* algorithm [13], which first copies the *initial state* to a *current state* variable and enters a while loop. Each iteration the algorithm randomly selects both a mapping (of a single query to a single keyword, i.e.  $q_1 \mapsto k_1$ ) and a keyword ( $k_2$ ) from the list of potential keywords. If the selected keyword is already mapped to another query, i.e.  $q_2 \mapsto k_2$  is in *current state*, the mappings are simply interchanged, i.e.  $q_1 \mapsto k_2$  and  $q_2 \mapsto k_1$ , otherwise the selected mapping is changed to  $q_1 \mapsto k_2$  to obtain the *next state*.

The algorithm determines whether it should accept or reject *next state* in favor of or against *current state* respectively. The algorithm calculates the sum of the squared Euclidean distance of the co-occurrence counts of each of the mappings with other mappings for both *current state* and *next state*. Depending on the calculated squared Euclidean distance either *next state* is accepted (and becomes *current state* in the iteration of the while loop) or is rejected. If the Euclidean distance of *next state* is lower than that of *current state* *next state* is accepted. A *next state* with a higher Euclidean distance is not necessarily rejected, but might be accepted with a small probability, depending how close the Euclidean distance is to 0. This is included to decrease the possibility of the algorithm finishing its run in a local optimum state as opposed to finding the global optimum state.

The ANNEAL algorithm takes three input parameters, next to the co-occurrence matrices. These input parameters *initial temperature*, *cool down rate* and *rejection rate* are used to ensure the algorithm has a finite run time. The *initial temperature* initializes the internal *current temperature* variable of the algorithm. Each iteration in the while loop *current temperature* is decreased by multiplying it with *cool down rate* (a value between 0 and 1, close to 1). The algorithm returns *current state* as the final mapping if the system *freezes*, i.e. *current temperature* becomes 0. *initial temperature* and *current temperature* therefore together determine the maximum number of loops the algorithm goes through. The algorithm can also finish before it freezes if no *next state* has been accepted for a certain (consecutive) number of iterations, which is determined by the value of *rejection rate*. The *current temperature* variable is also used also in deciding whether to accept a worse *next state* with a small probability.

## 2.6 Simulations

**Datasets Used.** Both Islam et al. [12] and Cash et al. [6] use the *\_sent\_mail* data folder of the ENRON dataset [2] (containing 30109 emails) as the dataset to run simulations on. Additionally, Cash et al. use the Apache Lucene project’s *java-user* mailing list [1] (containing 50116 emails) in their simulations.

**Tokenization/Stemming Algorithm.** Both papers tokenize all emails in the dataset to specific words before they are able to stem individual words, but neither elaborate on the tokenization algorithm used in their simulations. Stemming is done using Porter’s stemming algorithm [14] to get the stem of each word, meaning that words like ‘has’ and ‘have’, which in principle have the same meaning, are stemmed to the same word.

**Keyword Generation.** The stemmed keywords are sorted in decreasing order of overall occurrence. The 200 most occurring (stemmed) words in a dataset, that are likely to occur in every file (for example ‘a’ and ‘the’), are removed as they are not deemed useful in a Searchable Encryption scheme. The next  $x$  words are regarded as the keyword set.

**Query Generation.** Both Islam et al. and Cash et al. simulate a certain number of queries by using the Zipfian distribution on the keyword set. Due to the nature of the Zipfian distribution words with a higher occurrence count are more likely to be simulated as a query.

**Reported Results.** Islam et al. report recovery rates ranging from 60%–100% depending on the number of keywords, number of queries and the % of queries ‘known’ before the attack run. With different levels of Gaussian noise added to the background knowledge, the accuracy of the attack ranges between 40% and 85%. Cash et al. report recovery rates of the IKK attack ranging between 0% and 100% and show an exponentially decreasing correlation between the size of the input matrices (query and background knowledge co-occurrence matrices) and the recovery rate. Cash et al. also report recovery rates ranging between 0% and 60% for different percentages of documents ‘known’ to the adversary.

## 3 Revisiting the IKK Attack

Islam et al. [12] introduced the study on query recovery attacks by proposing the IKK query recovery attack. The authors report high query recovery rates that would allow an adversary to determine what a user searched for. More importantly, as Cash et al. [6] note in their paper, correctly recovered queries are inherently a part of the plaintext of encrypted documents and thus disclose part of the plaintext of the document stored on the server. We therefore stress that it is important to get a more broad understanding of query recovery attacks. In this research we revisit the following facets of the IKK attack:

- We evaluate the assumption on query distribution following the Zipfian distribution made by Islam et al. while simulating runs of the IKK attack (Sect. 4).
- We look at the correlation between the *initial temperature*, *cool down rate* and the *rejection rate* input parameters and the accuracy of the IKK attack (Sect. 5).

Furthermore, we propose and research the following improvements to the IKK attack:

- We propose to use a majority voting scheme to increase the accuracy of the IKK attack by combining the results of multiple runs (Sect. 6).
- We propose to (more) deterministically choose the next state of the ANNEAL algorithm to increase the accuracy of and decrease the number of visited states by the IKK attack. This method incorporates the (relative) word occurrence method, as proposed by Cash et al. in their Count attack (Sect. 6).

In order to run simulations of the IKK attack to address the points above we implemented the IKK attack as proposed by Islam et al. in Python3 and published it on Github [3]. The implementation allows the user to select:

- the distribution used to simulate queries (*Zipfian*, *reverse Zipfian*, *Uniform*)
- values for the parameters of the ANNEAL algorithm (*initial temperature*, *cool down rate*, *rejection rate*)
- sizes of the query and background knowledge co-occurrence matrices (resp. *number of queries*, *number of keywords*)
- datasets/email folders to use in the simulation (*ENRON/ \_sent\_mail*, *ENRON/inbox*, *ApacheLucene-java-user* (*Apache*))
- different methods to simulate non-perfect background knowledge (*Gaussian noise addition*, *using a fraction of the keywords*, *using a fraction of the documents*)
- the number of consecutive runs with exactly the same input parameters
- whether to more deterministically select new states using word occurrences as also proposed in the Count attack by Cash et al.

To give the reader an idea of the input parameters used in our simulations we mention the standard values for the different parameters of the IKK attack in Table 4.

We briefly capture our generalized method below. Simulation specific methodologies are elaborated upon in their correlated sections (Sects. 4, 5 and 6).

1. Tokenize and stem the words in all documents in a specific dataset. Tokenization is done by splitting the document on whitespaces. Stemming is done using Porter’s stemming algorithm [14].
2. Sort all unique (stemmed) words in decreasing order of occurrence (count) (the total number of times a word occurs in the dataset, not the number of matching documents).



**Table 4.** Standard parameter values in the IKK attack simulations

Variable	Value	Variable	Value	Variable	Value
<i>initial temperature</i>	1.0	<i>nr of keywords</i>	1500	<i>dataset /</i>	<i>ENRON /</i>
<i>cool down rate</i>	0.999	<i>nr of queries</i>	150	<i>email folder</i>	<i>_sent_mail</i>
<i>rejection rate</i>	50000				
		<i>nr of runs</i>	1	<i>keyword percentage</i>	1.0
<i>Gaussian</i>	0.0			<i>document percentage</i>	1.0
<i>noise scaling factor</i>		<i>distribution</i>	<i>Zipfian</i>		

- Disregard the first 200 most occurring words, just like Islam et al., and take the subsequent  $x$  words as *keyword set*.  $x$  is equal to the *number of keywords* input variable in our simulations.
- Simulate  $y$  queries from the  $x$  selected keywords using a specified query distribution as the *query set*.  $y$  is equal to the *number of queries* input variable in our simulations.
- Generate the query and background knowledge inverted indices from the selected queries and keywords, and the list of documents.
- Generate the query and background knowledge co-occurrence matrices from the inverted indices.
- Input the co-occurrences matrices and the input parameters *initial temperature*, *cool down rate* and *rejection rate* in the ANNEAL algorithm.
- Calculate the (query) recovery rate by dividing the number of correctly mapped queries, where  $query = keyword$ , by the total number of queries.

Islam et al. also use a *known queries* variable in their experiments, a method also adopted by Cash et al. This variable denotes  $\langle query, keyword \rangle$  pairs that the adversary knows to be mapped correctly before the attack run. We argue that the actual value of this variable is likely to be (close to) 0 and we therefore excluded this variable from our experiments.

## 4 Assumptions Evaluation

In their simulations, Islam et al. [12] make an assumption on distribution of queries in a real-world SE scheme in order to estimate real-world search behavior of users. They assume natural search behavior can be estimated by simulating queries using the Zipfian distribution as they argue that search behavior might follow a Zipfian distribution as the simulations are run on a natural language corpus. In their paper, the authors state that ‘according to Zipf’s law, in a corpus of natural language utterances, the frequency of appearance of an individual word is inversely proportional to its rank’ [17]. The Zipfian distribution is also used by Cash et al. [6] to simulate queries for their simulations.

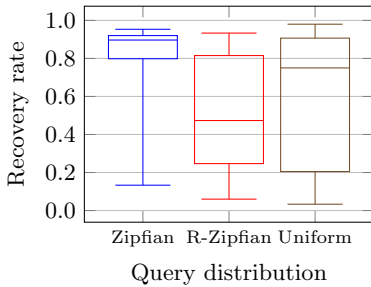
In order to simulate queries, from the simulated keyword set (of size  $x$ ), Islam et al. first sort the words in the keyword set in decreasing order of overall

occurrence. For the word in the  $j^{th}$  position in this list (rank  $j$ ) the following formulas are used to determine the probability the word is selected as a query:

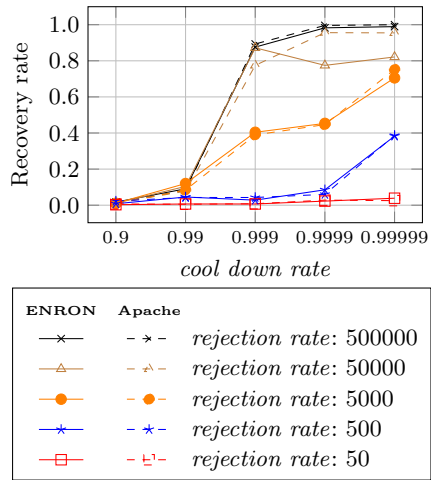
$$Pr_j = \frac{1}{N_x} = \frac{1}{j \times N_x}, \quad N_x = \sum_{i=1}^x \frac{1}{i} \tag{2}$$

A word with a higher total occurrence count is therefore more likely to be simulated as a query. Islam et al. also note that duplicate queries are removed. We argue that the assumption that search behavior follows a Zipfian distribution is counter-intuitive in the sense that users are more likely to search for a specific email in their mail archive and thus issue a (single word) query that is likely to return the sought after document while also not returning too much other emails (false positives). We therefore argue that search behavior might instead follow a reverse Zipfian distribution and thus a word that has a lower occurrence count has a higher chance of being selected as a query. The reverse Zipfian distribution can be calculated using the same formulas as the Zipfian distribution, but the list of word occurrences is sorted in ascending order of occurrence as opposed to descending order.

To compare the effect of the distribution used to simulate the queries we conducted three different simulations for the Zipfian distribution, reverse Zipfian Distribution and Uniform distribution respectively. The Uniform distribution denotes the setting where queries are simulated from the keyword set uniformly at random. The results of our simulations are shown in Fig. 1, where each box plot is the aggregation of 20 simulations.



**Fig. 1.** Correlation of different query distribution and recovery rate



**Fig. 2.** Correlation between cool down rate and recovery rate, with different values of rejection rate

It can be seen that the distribution chosen to simulate queries influences the results of the IKK attack quite a lot and that simulations where the queries were simulated using the Zipfian distribution in general have a much higher recovery rate than simulations where a different distribution was used. Unfortunately, we simply do not know what distribution real-world search behavior follows in a Searchable Encryption scheme as, to the best of our knowledge, there exists no dataset which contains query search behavior of real-world users in an SE setting. We can only conclude that the actual distribution determines the accuracy of the IKK attack and therefore whether using a Searchable Encryption scheme poses a risk for search privacy and potentially data confidentiality. These results are in line with a similar notion (high-selectivity keywords vs. low-selectivity keywords) as made in [4], which was published during our research.

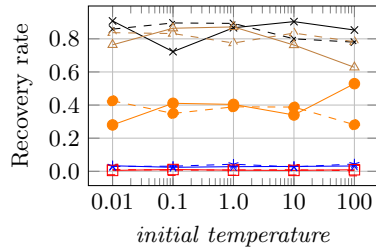
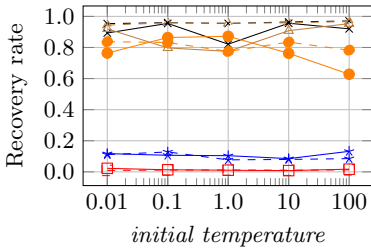
## 5 Recovery Rate Quantification

Islam et al. [12] show that their IKK attack allows an adversary to recover (most of the) queries in a simulated setting. Their ANNEAL algorithm (Algorithms 3, 4), which is part of their attack algorithm, takes three input parameters *initial temperature*, *cool down rate* and *rejection rate* to ensure the algorithm has a finite run time. These parameters are further explained in Sect. 2.5.

The values of these parameters have a significant influence on the number of visited states of the IKK attack as the *initial temperature* and *cool down rate* together determine the maximum number of states the algorithm visits, whereas the value of *rejection rate* determines whether the algorithm finishes before the *current temperature* reaches 0 or not. We argue that the accuracy of the IKK attack is therefore dependent on the values of these input parameters. This means that a proven correlation between these three input parameters, independent of the underlying dataset, and the recovery rate might allow an adversary to use simulations on another dataset to find the optimal input parameters for the IKK attack.

To answer the question whether there is a correlation between the three input parameters and the recovery rate, independent of the dataset, we used the same datasets as used by Islam et al. and Cash et al. as these datasets are most common in literature. The first dataset is the ENRON dataset [2], specifically its *\_sent\_mail* data folder which contains 30109 emails. Cash et al. also experiment on the *java-user* mailing list of the Apache Lucene project (henceforth Apache) [1] (reportedly containing about 38.000 emails). However, the exact dataset they used was unavailable and thus we crawled the archive site of the java-user mailing list and retrieved 50116 emails. The crawled Apache dataset is included in our Python3 implementation of the IKK attack on Github [3]. In order to test our hypothesis we conducted three different experiments. In all of the experiments we kept one of the input parameters (*initial temperature*, *cool down rate*, *rejection rate*) constant while varying the other two. The experiments were repeated for both the Apache and ENRON dataset, with both the query and background

knowledge co-occurrence matrix from the same dataset and with perfect background knowledge. Each point in Figs. 2, 3 and 4 is the average of 5 simulations of the IKK attack.



ENRON	Apache	
—x—	-*-	cool down rate: 0.99999
-△-	-△-	cool down rate: 0.9999
-●-	-●-	cool down rate: 0.999
-*-	-*-	cool down rate: 0.99
-□-	-□-	cool down rate: 0.9

ENRON	Apache	
—x—	-*-	rejection rate: 50000
-△-	-△-	rejection rate: 50000
-●-	-●-	rejection rate: 5000
-*-	-*-	rejection rate: 500
-□-	-□-	rejection rate: 50

**Fig. 3.** Correlation between *initial temperature* and Recovery rate, with different values of *cool down rate*

**Fig. 4.** Correlation between *initial temperature* and Recovery rate, with different values of *rejection rate*

Figure 2 shows the aggregation of simulation results with a constant *initial temperature*. The results of simulations on the ENRON dataset and the Apache dataset are roughly the same. The only exceptions are the simulations with a *rejection rate* of 50000 and a *cool down rate* of 0.9999 respectively 0.99999, which we attribute to the relatively low number of simulations (5) aggregated in each data point. With more simulations these results might become more similar. Furthermore, the recovery rate increases with both the *cool down rate* and the *rejection rate*. This makes sense as the maximum number of loops is increased with a *cool down rate* closer to 1 and a higher *rejection rate* increases the likelihood of finding the best mapping as the algorithm does not halt prematurely.

Figure 3 shows the aggregation of simulations with a constant *rejection rate*. We see that the value of recovery rate is only dependent on the value of *cool down rate* as the correlation between *initial temperature* and recovery rate is relatively constant. The recovery rate is also not dependent on the underlying dataset used as the results for both the ENRON and Apache dataset are roughly the same.

Figure 4 shows the aggregation of simulations with a constant *cool down rate*. We see that the value of recovery rate is dependent on the value of *rejection rate* and not on the value of *initial temperature* as we again see a constant correlation between *initial temperature* and recovery rate. We can also see that the value of recovery rate is not dependent on the underlying dataset used as the results are quite similar for both the ENRON and Apache datasets.

We conclude that the values of *rejection rate* and *cool down rate* significantly influence the recovery rate of the IKK attack. Furthermore, we conclude that it is possible for an adversary to find the optimal values for *cool down rate* and *recovery rate* using simulations on a different dataset as the recovery rate is independent of the underlying dataset used. This means that it is possible to use simulations on the ENRON dataset to select the optimal input parameter values for runs on the Apache dataset and vice versa. We argue that email datasets are quite similar due to the nature of the files they contain as emails are structured in a certain way, are limited in length and are used for specific purposes and thus might contain similar data. More research should be conducted to find out whether our findings hold true for completely different datasets as well.

## 6 Improvements

### 6.1 Combining Multiple Runs

The IKK attack returns a 1-to-1 mapping between queries and keywords. An adversary cannot, from the mapping alone, determine which queries were recovered correctly and which were not, as even with perfect background knowledge the IKK attack shows a lot of variance. For example, Fig. 5 shows that the recovery rates of 20 simulations of the IKK attack with equal input parameters and perfect background knowledge return recovery rates ranging between 0.1 and 0.98, which almost spans the entire range of possible recovery rates.

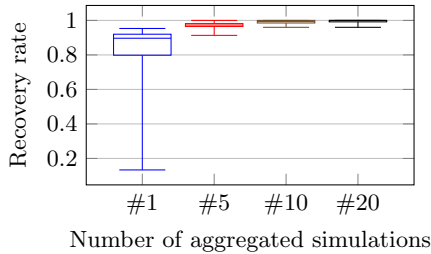
The IKK attack is a probabilistic algorithm in the sense that the algorithm uniformly at random selects a new mapping to change in the current state to determine the next state to explore. We argue that the IKK attack shows a large variance in the recovery rate as the algorithm merely approximates the optimal state yet does not necessarily always return it. A deterministic algorithm that simply visits every possible mapping is less likely to show a large variance, but the single attack runs will have to evaluate far more different states in order to be successful, making such an algorithm quite inefficient as the total number of potential states is given using the following equation:

$$\text{nr. of states} = \frac{(\text{nr. of keywords})!}{(\text{nr. of keywords} - \text{nr. of queries})!} \quad (3)$$

For 50 observed queries and 500 keywords in the background knowledge this would mean that there are already  $7.039 * 10^{133}$  potential states to explore.

As every single run of the IKK attack still approximates the optimal mapping, we argue that it is possible to combine the results of different attack runs using a simple majority voting scheme to better approximate the optimal solution. We conducted 20 simulations each consisting of 20 attack runs on the same query co-occurrence matrix and background knowledge co-occurrence matrix (representing perfect knowledge) per simulation. In each of the simulations, we combined a certain number of runs by selecting the most prevalent keyword mapped to each of the queries. If no prevalent keyword could be found (two or

more keywords are most prevalent) the majority voting scheme did not assign a most prevalent keyword to a query and instead assigned the *None* value.



**Fig. 5.** Aggregation of a number of different runs of the same simulation using a majority voting scheme

Figure 5 shows the results of combining multiple runs on the same query and background co-occurrence matrices. It can be seen that the accuracy of the attack significantly decreases the variance that is observed with single runs of the IKK attack. When combining 5 runs per simulation (#5) the results are already very promising, which is even more the case when the aggregation contains either 10 or 20 runs per simulation. Our proposed aggregation method also has the advantage that the single attack runs can be executed in parallel and then aggregated, ensuring the execution time overhead is limited. The median recovery rate between 1 run per simulation (#1) and 20 runs per simulation (#20) is increased with more than 10% points, whereas the variance is decreased with 78% points.

## 6.2 Deterministic IKK Attack

As the IKK attack is a probabilistic algorithm, it does not necessarily return the optimal query-to-keyword mapping. We argue a more deterministic approach to finding the right mapping might increase the recovery rates of the IKK attack.

The Count attack, as proposed by Cash et al. [6], takes a more deterministic approach to map queries to keywords, by eliminating candidate mapping keywords using the relative document occurrence count of keywords. Cash et al. assume the adversary has access to not only the co-occurrence counts of queries and keywords, but is also in possession of the (relative) document occurrence counts from queries and keywords, i.e. the number of documents a query or keyword occurs in (relative to the total number of documents in the dataset). The theory behind this is that, while assigning a keyword to a query, a lot of potential keywords can already be eliminated as their relative document occurrence count is not within a certain range of the relative document occurrence count of the query. These keywords therefore are not likely to be the right keyword corresponding to the query and thus can be disregarded.

The Count attack incorporates eliminating candidate keywords using their document occurrence count ‘and brute-forces all possible mappings for a small number of queries and returns the mapping which maximizes the number of disambiguated queries’. Cash et al. report much higher recovery rates from their deterministic Count attack as opposed to the probabilistic IKK attack. We propose to incorporate the candidate keyword elimination method of the Count attack while selecting new mappings in the IKK attack to both decrease the number of potential states to visit as well as increase its accuracy. We also argue that the accuracy of the attack will increase as the algorithm is likely to visit better states on average as the worst potential states are eliminated. Our method still differs from the method as used by Cash et al. as they propose a deterministic algorithm, whereas our algorithm still makes use of the probabilistic nature of the IKK attack.

In order to eliminate candidate keywords Cash et al. construct a confidence interval for the document occurrence count of each of the keywords using Hoeffding’s inequality [11]. The lowerbound ( $LB$ ) and upperbound ( $UB$ ) of the confidence interval per keyword  $k$  are calculated using the following formula(s):

$$LB_k, UB_k = \frac{c_k^s}{p_{pk}} \mp \sqrt{0.5 n \log 40} \quad (4)$$

In this formula  $c_k$  computes the document occurrence count of keyword  $k$  in the background knowledge dataset and  $p_{pk}$  denotes the size relativity between the query and background knowledge dataset.  $\frac{c_k}{p_{pk}}$  therefore denotes the expected document occurrence count of  $k$  in the query dataset.  $\epsilon = \sqrt{0.5 n \log 40}$  is used by Cash et al. to ensure the confidence interval has a confidence level of 95%.  $n$  denotes the number of documents in the query dataset. After calculating a confidence interval for each of the keywords the candidate keywords for a query can be calculated as  $S_q = \{k' \in K | LB_{k'} \leq c_q \leq UB_{k'}\}$ .  $S_q$  denotes the candidate keyword set,  $K$  the keyword set and  $c_q$  denotes the document occurrence count of query  $q$ .

The Original IKK attack maps queries to keywords in two places in the algorithm, namely when selecting the *initial state* (Algorithms 1 and 2) and while selecting a *next state* (Algorithms 3, 4 and 5). We therefore incorporated the method of Cash et al. in two places in our Deterministic IKK attack:

While selecting the initial state we first assign a *None* value to queries of which  $S_q$  is an empty set, meaning that no keywords are in range. These queries are left unchanged throughout the entire algorithm run and thus are assigned *None* in the final mapping as well. This also allows an adversary to determine which queries were not mapped to a keyword. Then all queries with a non-empty candidate set  $S_q$  are ordered in ascending order of the size of  $S_q$  and each of the queries, starting at the query with the lowest size of  $S_q$ , is assigned a random keyword in  $S_q$  that was not yet assigned to another query. As we enforce the 1-to-1 mapping property of the IKK attack this potentially creates the edge case where all keywords in  $S_q$  of a query are already assigned to other queries. The algorithm tries, with a depth of one, whether it is possible to re-assign one of the other queries to ‘free up’ a keyword in  $S_q$ . If it succeeds the ‘freed’

keyword is assigned to the query, otherwise the query is assigned *None* and is thus disregarded during the rest of the algorithm run.

While selecting a new state we choose a random query, keyword mapping, e.g.  $q_1 \mapsto k_1$ , from the queries in the *current state* that were not assigned *None* and we select a random keyword  $k_2$  from  $S_{q_1}$  as opposed to the full keyword set, while ensuring  $k_1 \neq k_2$ . Then, just like in the Original IKK attack there are two possibilities:

If  $k_2$  was mapped to a query  $q_2$  we try whether keyword  $k_1$  is in range of query  $q_2$  and interchange the mapping if so. If not, we keep (uniformly at random) selecting a new keyword  $k_2$  and checking whether the new  $k_2$  adheres to the right properties. If we cannot find a satisfactory candidate  $k_2$  for a certain number of loops (2 times the size of the keyword set in our simulations) the algorithm returns the current state as the next state, which is rejected as the Euclidean distance is not better than the old current state (as they are the same). If  $k_2$  was not mapped to any query in the current state we change the next state so that  $q_1 \mapsto k_2$ .

In order to compare our deterministic version of the IKK attack to the Original IKK attack, especially in cases where the adversary does not have full background knowledge, we defined a metric that captures the correlation between the similarity between the co-occurrence matrices and the recovery rate as we argue that it is important to research the effect of our improvements on simulations with different levels of background knowledge to get a broad understanding of the effects of our improvements. We have used four different methods to simulate non-perfect background knowledge (omitting a percentage of the selected keywords, omitting a percentage of the documents in a dataset, adding Gaussian noise and using a different dataset as background knowledge). These methods, the defined metric and the results with the original IKK attack are included in Appendix A.

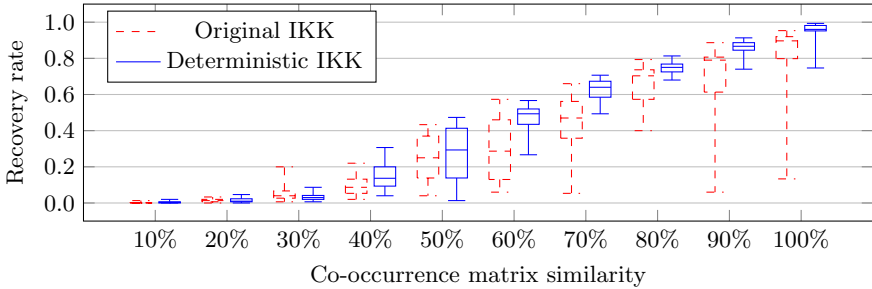
**Table 5.** Nr. of states visited by the IKK and Deterministic IKK attack

Parameters	IKK version	Min.	Max.	Avg.
# <i>Total loops</i>	Original	531,733	737,741	733,213
	Deterministic	196,790	737,741	526,038
# <i>Accepted loops</i>	Original	7783	9535	8533
	Deterministic	7287	101,145	14,252
$\frac{\# \text{ Accepted loops}}{\# \text{ Total loops}}$	Original	0.0105	0.0159	0.0117
	Deterministic	0.0099	0.1371	0.0263

Table 5, the aggregation of 500 simulations of both algorithms, shows that the Deterministic IKK algorithm visits much less total states on average than the Original IKK attack. Additionally, the average number of iterations where the next state is accepted is much higher and the ratio between the number of

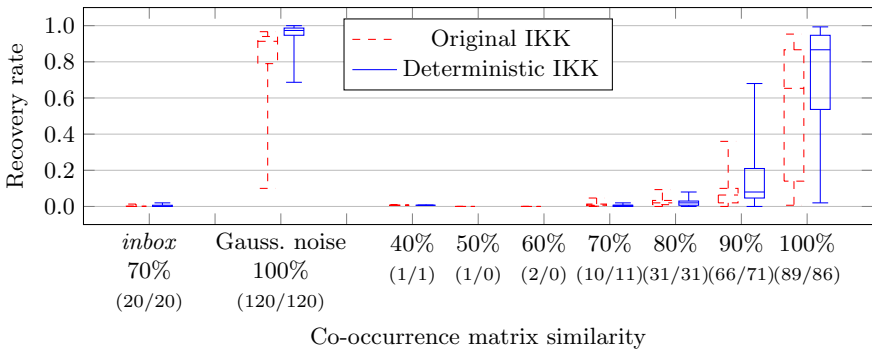


*accepted* loops and total loops is more than twice as high for the Deterministic IKK attack. It is useful to note that both attacks at most visit 737,741 different states and then return their current state as the final mapping. This is due to the chosen values of the input parameters of the ANNEAL algorithm and explains the values in the **Max.** column of the *# Total loops* row.



**Fig. 6.** Original/Deterministic IKK recovery rates and recovery rate, simulating non-perfect background knowledge by regarding a *percentage* of keywords

In Fig. 6 we compare recovery rate of the Original and Deterministic IKK attack when only a fraction of the actual keywords simulates background knowledge. The recovery rates of the Original IKK attack and methods to generate non-perfect background knowledge are the same as expressed in Fig. 8 and each bucket in the figure is the aggregation of 20 simulations. We see the same linear correlation for the Deterministic IKK attack as we saw before for the Original IKK attack, however, recovery rates of the Deterministic IKK attack show much less variance as well as a higher median value.



**Fig. 7.** Original/Deterministic IKK recovery rates and recovery rate, simulating non-perfect background knowledge using other methods

Figure 7 shows the comparison of the Original and Deterministic IKK attack when non-perfect background knowledge is simulated using other methods than using a percentage of all keywords. Figure 7 therefore also contains the same information as Figs. 9 and 10. The (non-percentage) numbers between brackets denote the number of simulations aggregated in that box plot. Due to the way in which we generate non-perfect background knowledge these are not the round number of 20 simulations per box plot as is the case in Fig. 6.

In simulations where we took a different, but similar dataset as background knowledge (*ENRON/inbox*) we see that both attacks have recovery rates close to 0 and in simulations where we added Gaussian noise to the background knowledge we see that the Deterministic IKK attack again shows less variance and higher recovery rates.

In simulations where a percentage of user folders in a dataset was used to simulate background knowledge we see the same exponential correlation between co-occurrence similarity and the recovery rate as we see for simulations using the Original IKK attack. Additionally, we see that the Deterministic IKK attack achieves higher recovery rates on average, but we do not see the drop in variance that we saw in simulations using the other methods to simulate non-perfect background knowledge.

All in all, we conclude that using components of the Count attack by Cash et al. [6], that make the IKK attack more deterministic, is a promising method to both decrease the number of states visited in a single attack run (28% decrease) and increase the recovery rate, as the median recovery rate is increased up to 21% points (Fig. 7, 100% box plot) and the variance is decreased up to 57% points (Fig. 6, 100% box plot).

## 7 Related Work

The first Searchable Encryption scheme was introduced by Song et al. [15] to allow for (plaintext) search among a set of encrypted documents. Their paper introduces the first In-place SE scheme which uses a stream cipher to scan for the occurrence of a plaintext word as well as introduces the notion of the potentially more efficient Encrypted-Index SE schemes. Song et al. already note that these schemes leak access patterns and that statistical attacks might disclose information of encrypted documents, but do not research this further. The notion of Oblivious RAM (ORAM) [10], introduced before the first SE scheme, is frequently mentioned as a method to not disclose access patterns. Oblivious RAM, however, in a Searchable Encryption scheme is computationally quite expensive. A less expensive version specifically targeted for encrypted search, proposed by Curtmola et al. [8], still is computationally inefficient. Other papers propose to obfuscate access patterns by introducing inconsistencies in the search results by modifying the internal encrypted index of the SE scheme [7] or by using Bloom filters [5,9]. These schemes are reportedly computationally expensive as well.

The first statistical attack, the IKK attack, on Searchable Encryption schemes that leak access patterns was proposed by Islam et al. [12]. This

attack uses co-occurrence counts of observed queries to determine what plaintext word(s) the client searched for. Cash et al. [6] recognize that a recovered query inherently discloses part of the plaintext of encrypted documents and propose their Count attack as a response to the IKK attack. The Count attack uses the (relative) document occurrence counts next to the co-occurrence counts of queries to deliver better results faster as opposed to the IKK attack. Cash et al. also define different levels of leakage of SE schemes and coin the term *leakage-abuse attacks* to more broadly describe attacks that are intended to disclose information on the contents of encrypted documents in SE schemes as opposed to attacks that only disclose what the client searched for. Leakage-abuse attacks were further researched by Blackstone et al. [4].

Both the IKK attack and the Count attack are passive attacks, meaning that the adversary acts according to the protocol of the SE scheme, but tries to additionally obtain as much information and potentially runs calculations in parallel. Zhang et al. [16] show that an adversary capable of injecting files into a Searchable Encryption scheme that leaks access patterns ‘is devastating for query privacy’.

Blackstone et al. [4], which was published during our research, deserves a special mention. This paper focuses on new attacks, but also includes experiments using the IKK and Count attacks as these experiments are used for the comparison to newer attacks. Our paper, instead, studies the IKK attack in-depth to shed more light on its assumptions and practicality. The reported IKK recovery scores align with the results in this paper.

## 8 Conclusion

In this paper, we revisited the IKK query recovery attack on Searchable Encryption schemes as proposed by Islam et al. [12].

We show that the assumption that queries in a Searchable Encryption scheme follow a Zipfian (query) distribution, as Islam et al. made while simulating queries, positively influences the recovery rate of the IKK attack. Furthermore, we show a correlation between input parameters of the IKK attack, of which the values were left unexplained by Islam et al. and the recovery rate of the IKK attack, independent of the underlying dataset used in the SE scheme. This potentially allows the adversary to optimize the parameter values using a different dataset before executing the actual attack.

We also propose improvements to the IKK attack by showing that the accuracy of the attack can be improved significantly by combining multiple attack runs, as we show that median recovery rates can be increased up to 10% points, whereas the variance of recovery rates of simulation can be decreased up to 78% points. In addition, we show that the accuracy of the IKK attack can be increased, while the number of states visited can be decreased by incorporating deterministic components, based on notions made by Cash et al. [6] in their Count attack, to the IKK attack. The average number of states visited

is decreased by 28%, the median recovery rate is shown to be increased up to 21% points in different simulations, whereas its variance is decreased up to 57% points.

In recent literature, e.g. [7], obfuscation methods have been proposed to combat the effectiveness of the IKK attack. We leave research into the effects of these obfuscation methods on the effectiveness of our Deterministic IKK open for the reader.

## Appendices

### A Co-occurrence Matrix Correlation (Partial Background Knowledge)

Both Islam et al. [12] and Cash et al. [6] both briefly elaborate on the recovery rate of the IKK attack in the case where the adversary only has partial background knowledge. Islam et al. add various degrees of Gaussian noise to individual cells in the co-occurrence matrix representing perfect background knowledge to simulate this setting, whereas Cash et al. simulate non-perfect background knowledge co-occurrence matrix by taking a fraction of all documents in the dataset. Both papers show that the accuracy of the attack is greatly dependent on the level of background knowledge the adversary has. We therefore argue that it is important to get a better understanding of the correlation between the level of background knowledge the adversary has and the recovery rate of the IKK attack. We also argue that the level of background knowledge can be expressed as a similarity between the query and background knowledge co-occurrence matrices, i.e. the *co-occurrence matrix similarity*.

In order to express co-occurrence matrix similarity we propose a metric that returns a similarity score between 0 (no similarity) and 1 (equivalent matrices) between two co-occurrence matrices of the same dimensions. For two matrices  $M_1$  and  $M_2$  and arbitrary words  $a, b$  (corresponding to a row and column) the following formulas are used:

$$\Delta_{a,b}^2 = \begin{cases} (M_1[a, b] - M_2[a, b])^2, & \text{if } a, b \in M_1 \text{ and } a, b \in M_2 \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

$$\Delta_{total}^2 = \sum_{\forall a, b \in M_1} \Delta_{a,b}^2 \quad (6)$$

$$\epsilon_{a,b} = \begin{cases} 1, & \text{if } a, b \in M_1 \text{ and } a, b \in M_2 \\ 0, & \text{otherwise} \end{cases} \quad (7)$$

$$\epsilon_{total} = \sum_{a, b \in M_1} \epsilon_{a,b} \quad (8)$$

$$Co - ocsim. = \left( 1 - \frac{\Delta_{total}^2}{\epsilon_{total}} \right) * \left( \frac{K_{overlap}}{K_{total}} \right) \quad (9)$$

Equations 5 and 6 are used to calculate the total squared Euclidean distance of cells that occur both in  $M_1$  and  $M_2$ .

**Table 6.** Example co-occurrence matrices

$M_1$	a	b	c	$M_2$	a	b	c	$M_3$	a	b	c	$M_4$	a	b	d
a	1	1	1	a	1	1	1	a	0	0	0	a	1	1	1
b	1	1	1	b	1	1	1	b	0	0	0	b	1	1	1
c	1	1	1	c	1	1	1	c	0	0	0	d	1	1	1

Equations 7 and 8 are used to calculate the number of cells that occur in both  $M_1$  and  $M_2$ .

In Eq. 9 we calculate the average squared Euclidean distance of cells that occur in both matrices and multiply this by the ratio of *row identifiers that occur in both matrices* ( $K_{\text{overlap}}$ ) to the *total number of rows in both matrices* ( $K_{\text{total}}$ ).

In Table 6 matrices  $M_1$  and  $M_2$  are exactly the same.  $\Delta_{\text{total}}^2$  is 0 as the squared Euclidean distance between each of the cells is  $(1 - 1)^2 = 0$ . The average is therefore also 0. As all keywords in both matrices also occur in the other matrix,  $\frac{K_{\text{overlap}}}{K_{\text{total}}} = 3/3 = 1$ . The similarity between the matrices is calculated as *co-ocsim.* =  $(1 - 0) * 1 = 1$  meaning that the matrices are exactly the same. The calculation for the similarities between matrices  $M_1$  and  $M_3$ , and  $M_1$  and  $M_4$  gives the values 0 and 2/3 respectively. In our simulations we calculate the co-occurrence matrix similarity using the perfect background knowledge co-occurrence matrix  $M_F$  (which corresponds with an unqueried query co-occurrence matrix) and a non-perfect background knowledge co-occurrence matrix  $M_P$ . Both matrices have the same dimensions. In order to simulate partial background co-occurrence matrix  $M_P$  we use the following methods:

**Gaussian noise addition** - We use the method by Islam et al. to add Gaussian noise in various degrees to the cells in  $M_F$  to obtain  $M_P$ .

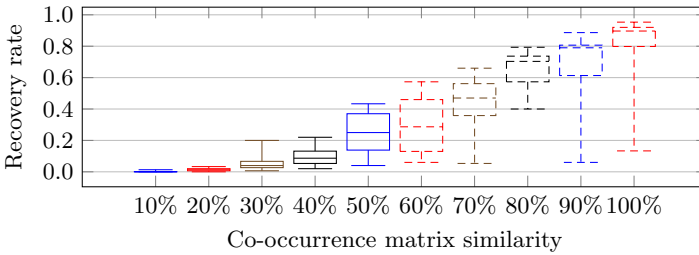
**Document percentage** - In this setting we use 10% to 100% of the user folders in a dataset to generate  $M_P$ . This method differs a bit from the method by Cash et al. as we argue that the adversary is more likely to obtain a percentage of the mail boxes of users (and all documents that are in these folders) than a percentage of all documents, selected uniformly at random, in a dataset. We believe that this choice might influence the results as different users are likely to use specific language in (all of) their emails.

**Keyword percentage** - In this setting we, uniformly at random, select 10% to 100% of the keywords in  $M_F$  to obtain  $M_P$ . To keep the dimensions of  $M_P$  consistent throughout all our simulations we supplement the selected keywords with words with a lower occurrence count in the dataset used, i.e. that were not in the keyword set.

**Different input folder** - In this setting we use a different, but similar dataset to generate  $M_P$ . In our simulations we use the *inbox* folder (containing 44859 emails) of the ENRON dataset.

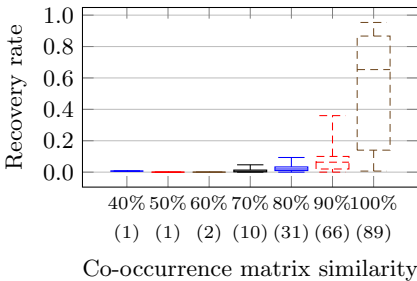
The results of the different methods are shown in Figs. 8, 9 and 10. In these figures we group the values into certain buckets to group similarity scores. If a

co-occurrence similarity score is between 0 and 0.1 it is put in the 10% similarity bucket, a value between 0.1 and 0.2 is put in the 20% bucket and so on.

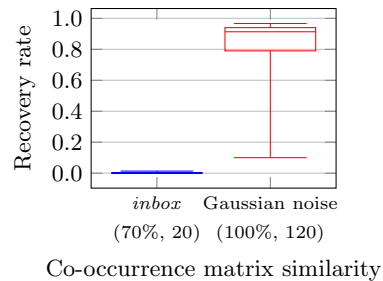


**Fig. 8.** Correlation between co-occurrence matrix similarity and recovery rate, simulating non-perfect background knowledge by regarding a *percentage* of keywords

Figure 8 shows the correlation between the co-occurrence similarity and the recovery rate in the setting where we use a certain percentage of the keywords in the keyword set to simulate partial background knowledge. Each bucket represents 20 simulations. Due to the nature of our similarity score using 90% of keywords from the keyword set will result in a similarity score of exactly 90%. The figure shows a clear linear correlation between the co-occurrence similarity score and the recovery rate. This makes sense as in this setting entire rows (and thus also columns) that are present in  $M_F$  are changed while simulating  $M_P$ . The highest possible percentage of recoverable queries therefore is linearly dependent on the percentage of keywords regarded. As the individual cell values are not changed while simulating  $M_P$  (as opposed to the other methods) the algorithm is likely to recover (most of the) queries of which the corresponding keywords were selected in  $M_P$  as these have the optimal Euclidean distance of 0.



**Fig. 9.** Correlation between co-occurrence matrix similarity and recovery rate, simulating non-perfect background knowledge by regarding a *percentage* of user folders in the dataset



**Fig. 10.** Correlation between co-occurrence matrix similarity and recovery rate, simulating non-perfect background knowledge by using the *ENRON/inbox* data or by adding Gaussian noise

Figure 9 shows the correlation between the co-occurrence matrix similarity and the recovery rate in simulations where non-perfect background knowledge is simulated by taking a percentage of user folders in a dataset to generate  $M_P$ . We ran 20 simulations for each percentage ranging from 10%, 20% to 100% each. The first thing that we notice is that the buckets do not contain the same number of simulations per bucket, which is shown in the figure as the number between brackets. This is due to the fact that we uniformly at random select a percentage of all user folders in a dataset and these user folders do not contain the same number of documents. Different writing styles can also be of influence to the overall co-occurrence similarity. The results in Fig. 9 show a different correlation than the results in Fig. 8. This makes sense as in these simulations both row/column identifiers as well as individual cell values are changed. The algorithm is less likely to correctly map queries to keywords that do occur in  $M_P$  as the changed individual cell values, in Fig. 9, make it less likely to find the optimal mapping. We note that the results in the 40%-60% bucket do not give much information, as each bucket consists of a single simulation. We conclude from the rest of the results that the co-occurrence matrix similarity and the recovery rate show an exponential correlation in Fig. 9.

The results in Fig. 10 show the correlation between the co-occurrence matrix similarity and the recovery rate of simulations where  $M_P$  was generated using a similar, but different dataset (*ENRON/inbox*) and when we add Gaussian noise to various degrees. First of all, if we generate  $M_P$  using the *ENRON/inbox* data folder we obtain a similarity score of approximately 70% to the *ENRON/\_sent\_mail* dataset. The recovery rate of almost 0 is consistent with our results in Fig. 9.

With the addition of various degrees of Gaussian noise (with  $C$  values 0.0, 0.2, ..., 1.0) the similarity of the co-occurrences matrices is always between 0.999 and 1.0. This can be explained as this method does not change the row/column identifiers, but only changes the individual cell values (co-occurrence counts). As only a little noise is added most of these cell values stay relatively the same. The recovery rate distribution among 120 simulations is relatively high as opposed to other methods to generate non-perfect background knowledge.

We conclude that it is not possible to use our matrix similarity metric to find a single correlation between the similarity of co-occurrence matrices and the recovery rate. The different methods change non-perfect background knowledge  $M_P$  in different manners and this influences the results of the IKK attack a lot. The IKK attack correctly recovers queries if the co-occurrence counts in  $M_P$  exactly match (or are close to) those in  $M_F$ , which is shown in Figs. 8 and 10 (*Gaussian noise addition*). If the co-occurrence counts in  $M_P$  are further away from those in  $M_F$ , which is the case in Figs. 9 and 10 (*Inbox folder*) the accuracy of the IKK attack decreases drastically.

We argue that the scenario where the background knowledge, represented as  $M_P$ , is generated by taking a percentage of the user folders in a dataset is the most realistic one in a real-world scenario. It is not unlikely that an adversary, somehow, gets access to a certain set of the plaintext contents of the email boxes

of specific users. The IKK attack proves to be a powerful attack which can break the privacy of queries as well as data confidentiality of documents stored encrypted of the server, yet it is only exploitable by a powerful adversary, which has access to a dataset which results in background knowledge that is at least 90% similar the actual dataset encrypted on the server, as can be seen in Fig. 9.

## B IKK Algorithms

In this section, we cite (part of) the Simulated Annealing (SA) algorithms as proposed by Islam et al. [12] as well as formalize our proposed algorithms for the Deterministic version of the IKK attacks. In short:

- Algorithm `Optimizer` (Algorithms 1 and 2) is used to select the initial state of the `ANNEAL` algorithm.
- Algorithm `ANNEAL` (Algorithms 3 and 4) is the heart of the IKK attack and is the actual algorithm that maps queries to keywords (apart from setting the initial state). The algorithms displayed are a simplified version of the `ANNEAL` algorithm as presented by Islam et al. and are mainly included to illustrate changes we made to more deterministically select the `nextState`.

---

### Algorithm 1: Optimizer

---

```

input
  V : variable list
  // List of all (non-mapped) queries
  D : domain list
  // List of all (non-mapped) keywords
  K : known assignments
  // Known query-keyword mappings
  Mc          // Query co-occurrence matrix
  Mp          // Background knowledge co-oc matrix
  Q = {q: Sq} // Queries and their candidate keywords
1 initState ← {} // Initial state
2 valList ← copy D
  // Copies values in D to variable valList
3 foreach var ∈ V do
4   val ← random.choice(valList)
  // Randomly selects a keyword from valList
5   add {var ↦ val} to initState
  // Adds mapping to initState
6   remove val from valList
  // As query to keyword mappings are 1-to-1
end
7 nonAssignableQueries ← {}
  // Var. containing None-assigned queries
8 sortedQ ← sort(Q, key=len(Sq), ascending=True)
  // Sorts Q on number of candidate keywords

```

---



**Algorithm 2:** Optimizer (cont.)

---

```

9  foreach var, Svar ∈ sortedQ do
10  candKeywords ← Svar
    // Gets candidate keywords for query v
11  if len(candKeywords) == 0 then
12  |   add v to nonAssignableQueries
    // Query v added to nonAssignable Queries
    else
13  |   assignedKWs = []
    // Every candidate keyword per query is
14  |   nonAssignedKWs = []
    // already assigned to another query or not
15  |   foreach cand ∈ candKeywords do
16  |   |   if cand ∈ valList then
17  |   |   |   add cand to nonAssignedKWs
    // cand not assigned to another query
    |   |   else
18  |   |   |   add cand to assignedKWs
    // cand assigned to another query
    |   |   end
    |   end
    |   if len(nonAssignedKWs) ≠ 0 then
19  |   |   val ← random.choice(nonAssignedKWs)
    // Selects keyword from nonAssignedKWs
20  |   |   add {var ↦ val} to initState
    // Adds mapping to initState
21  |   |   remove val from valList
    // Query/keyword mappings are 1-to-1
22  |   |   else
23  |   |   |   foreach k ∈ assignedKWs do
24  |   |   |   |   q ← initState.getByValue(k)
    // Get query q, mapped to keyword k
25  |   |   |   |   if k ∈ Sq then
26  |   |   |   |   |   remove { q ↦ k } from initState
    // Removes old mapping from initState
27  |   |   |   |   |   add { q ↦ val } to initState
    // Adds new mapping to initState
28  |   |   |   |   |   add { var ↦ k } to initState
    // Adds new mapping to initState
    |   |   |   |   |   break
    |   |   |   |   end
    |   |   |   end
29  |   |   |   if initState.get(var) == None then
    |   |   |   |   add var to nonAssignableQueries
    // If no suitable mapping could be found
    |   |   |   end
    |   end
    end
    end
    end
30  add K to initState
    // Adds known mappings to initState
31  return ANNEAL(initState, D, Mp, Mc)
    // Returns result of function ANNEAL()
32  return ANNEAL(initState, D, Mp, Mc, nonAssignableQueries, Q)

```

---

---

**Algorithm 3: ANNEAL**

---

```

input           // Simulated Annealing parameters
  initState
  D
  Mc, Mp
  initTemperature
  // initial temperature variable
  coolingRate // cool down rate variable
  rejectThreshold
  // rejection rate variable
  nonAssignableQueries
  // List of None assigned queries
  Q = {q; Sq} // Queries and their candidate keywords
1 currentState ← initState
  // Search continues until temp. reaches 0
2 succReject ← 0 // or the system is frozen (no new state
3 currT ← initTemperature
  // is accepted for large number of times)

```

---



---

**Algorithm 4: ANNEAL (cont.)**

---

```

1 while (currT ≠ 0 and succReject < rejectThreshold) do
2   | currentCost, nextCost ← 0, 0
3   | nextState ← findNextState(currentState, D)
  // Selects nextState using the method by Islam et al.
4   | nextState ← findNextStateDet(currentState, D, Q)
  // Selects nextState deterministically
5   | E ← costCalculation(nextState, currentState)
  // Calculates cost difference of two states,
6   | probability = exp(-E/currT)
  // using the method by Islam et al.
7   | acceptNewState = (E < 0) or (random.choice j probability)
  //nextState accepted if E < 0 or with prob. exp(-E/currT)
8   | if acceptNewState then
9     | | succReject, currentState ← 0, nextState
  else
10    | | succReject++
  end
11  | currT = coolingRate*currT
  // temperature is decremented each loop
  end
12 foreach query ∈ nonAssignableQueries do
13  | | add {query ↦ None} to currentState
  // Maps non-assignable queries to None
  end
14 return currentState

```

---

**Algorithm 5:** findNextStateDet

---

```

input
  currentState // 1-to-1 mapping of all queries to keywords
  Q = {q: Sq} // Queries and their candidate keywords
  D
1 nextState ← copy currentState
2 {x ↦ y} ← random.choice(nextState)
3 Sx ← Q.get(x) // Gets candidate keywords for query x
4 cand ← Sx.remove(y)
  // Keyword y cannot be selected again
5 y' ← None // Initializes y'
6 if len(cand) ≠ 0 then
  | y = random.choice(cand)
  | // Selects random keyword from candidates
else
  | return currentState
  | // No new mapping could be found
end
7 count ← 0 // Initializes count variable
8 while {z ↦ y'} ∈ currentState and y ∉ Sz do
9   | y' = random.choice(cand)
  | // Selects a new candidate keyword y'
  | if len(count) ≤ 2 * len(D) then
  | | return currentState
  | | // No new mapping could be found
  | end
10  | count += 1
end
11 remove {x ↦ y} from nextState
12 add {x ↦ y'} to nextState
13 if {z ↦ y'} ∈ currentState then
14   | remove {z ↦ y'} from nextState
  | // If y' is already mapped to query z
15   | add {z ↦ y} to nextState
  | // Map query z to y instead of y'
end
16 return nextState

```

---

- Algorithm `findNextStateDet` (Algorithm 5) is our proposed sub algorithm of the ANNEAL algorithm which selects a new state of the algorithm more deterministically.

In order to easily annotate differences between the Original IKK attack (as proposed by Islam et al.) and our proposed Deterministic IKK attack we use the following colors:

- **Black** annotates lines that are present in both the Original and Deterministic IKK attack.

- **Red** annotates lines that are present in the Original IKK attack, but not in the Deterministic IKK attack. Red lines are replaced by blue lines.
- **Blue** annotates lines that are not present in the Original IKK attack, but are in the Deterministic IKK attack. Blue lines replace red lines.

The Original IKK attack and the Deterministic IKK attack are elaborated upon in Sects. 2.6 and 6 respectively.

We note that the pseudo code in the algorithms as shown below does not fully match with our implementation of the Original/Deterministic IKK attack [3]. First of all, we used Python3 specific methods to easily implement both attacks and, in order to reduce the number of lines we re-used as much of the code of the Original IKK attack as possible to implement our Deterministic IKK attack.

## References

1. Apache Lucene java-user email dataset, September 2001-July 2011. [http://mail-archives.apache.org/mod\\_mbox/lucene-java-user/](http://mail-archives.apache.org/mod_mbox/lucene-java-user/). Accessed 19 May 2020
2. ENRON email dataset, version 7th May 2015. <https://www.cs.cmu.edu/~.enron/>. Accessed 19 May 2020
3. IKK query recovery attack implementation (Python). <https://github.com/rubengrootroessink/IKK-query-recovery-attack>. Accessed 27 June 2020
4. Blackstone, L., Kamara, S., Moataz, T.: Revisiting leakage abuse attacks. *IACR Cryptol. ePrint Arch.* **2019**, 1175 (2019)
5. Boneh, D., Kushilevitz, E., Ostrovsky, R., Skeith, W.E.: Public key encryption that allows PIR queries. In: Menezes, A. (ed.) *CRYPTO 2007*. LNCS, vol. 4622, pp. 50–67. Springer, Heidelberg (2007). [https://doi.org/10.1007/978-3-540-74143-5\\_4](https://doi.org/10.1007/978-3-540-74143-5_4)
6. Cash, D., Grubbs, P., Perry, J., Ristenpart, T.: Leakage-abuse attacks against searchable encryption. In: *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, pp. 668–679. ACM (2015)
7. Chen, G., Lai, T.H., Reiter, M.K., Zhang, Y.: Differentially private access patterns for searchable symmetric encryption. In: *IEEE INFOCOM 2018-IEEE Conference on Computer Communications*, pp. 810–818. IEEE (2018)
8. Curtmola, R., Garay, J., Kamara, S., Ostrovsky, R.: Searchable symmetric encryption: improved definitions and efficient constructions. *J. Comput. Secur.* **19**(5), 895–934 (2011)
9. Goh, E.J.: Secure indexes. *IACR Cryptology ePrint Archive*, 216 (2003)
10. Goldreich, O., Ostrovsky, R.: Software protection and simulation on oblivious rams. *J. ACM (JACM)* **43**(3), 431–473 (1996)
11. Hoeffding, W.: Probability inequalities for sums of bounded random variables. In: *The Collected Works of Wassily Hoeffding*, pp. 409–426. Springer (1994)
12. Islam, M.S., Kuzu, M., Kantarcioglu, M.: Access pattern disclosure on searchable encryption: ramification, attack and mitigation. In: *NDSS*. vol. 20, p. 12. Citeseer (2012)
13. Kirkpatrick, S., Gelatt, C.D., Vecchi, M.P.: Optimization by simulated annealing. *Science* **220**(4598), 671–680 (1983)
14. Porter, M.F.: *Snowball: A language for stemming algorithms* (2001)

15. Song, D.X., Wagner, D., Perrig, A.: Practical techniques for searches on encrypted data. In: Proceeding 2000 IEEE Symposium on Security and Privacy. S&P 2000, pp. 44–55. IEEE (2000)
16. Zhang, Y., Katz, J., Papamanthou, C.: All your queries are belong to us: the power of file-injection attacks on searchable encryption. In: 25th USENIX Security Symposium (USENIX Security 16), pp. 707–720 (2016)
17. Zipf, G.K.: Selected studies of the principle of relative frequency in language (1932)