

# DARD: Deceptive Approaches for Robust Defense Against IP Theft

Alberto Maria Mongardini<sup>1</sup>, Massimo La Morgia<sup>2</sup>, Sushil Jajodia<sup>3</sup>, *Life Fellow, IEEE*,  
Luigi Vincenzo Mancini<sup>1</sup>, and Alessandro Mei, *Member, IEEE*

**Abstract**—With the rise of smart working and recent global events, the risk of cyberattacks is increasing steadily. Sometimes adversaries focus on stealing valuable data, such as intellectual property (IP): they exfiltrate a large volume of IP documents from a target company. They then identify those of their interest by leveraging automated methods. This work proposes the DARD (Deceptive Approaches for Robust Defense against IP theft) system, a framework designed to deceive adversaries who rely on automatic approaches to classify exfiltrated documents. Starting from an original repository of documents, DARD automatically generates a new deceptive repository that misleads popular automatic approaches, resulting in clusters of documents that are significantly different from the actual ones. By utilizing this approach, DARD aims to hinder the accurate clustering and the identification of the topic of documents by adversaries relying on automated techniques. The paper presents four deceptive operations (Basic Shuffle, Shuffle increment, Shuffle reduction, and Change topic) that DARD leverages to create a deceptive repository. We evaluate the efficacy of our approach by considering three different types of adversaries, each possessing varying levels of knowledge and expertise. Through extensive experiments, we show that the DARD system can deceive both automatic topic modeling and document clustering techniques, including widely-used commercial tools such as Amazon Comprehend. Hence, our solution provides a robust defense mechanism against Intellectual Property (IP) theft.

**Index Terms**—Deceptive repository, clustering, topic modeling, adversarial setting.

## I. INTRODUCTION

ACCORDING to recent cybersecurity reports from sources such as Deloitte [1] and Interpol [2], there is a noticeable rise in businesses suffering cyber-attacks. This upward trend

can be attributed to several factors, including the growing number of employees working remotely, which has become increasingly prevalent since the onset of the COVID-19 pandemic. Additionally, the outbreak of war in Ukraine has led to increased threats of cyberattacks against Western businesses, with reported attacks against European companies in particular [3]. In 2022, multiple exfiltration attacks occurred, where unauthorized individuals extracted data from targeted systems and publicly released it via platforms like TOR [4] or Telegram [5].

The Cybersecurity and Infrastructure Security Agency (CISA) registered an exfiltration attack within the Defense Industrial Base organization [6]. The adversaries infiltrated the organization's information system, compromised its network, and illicitly accessed and stole the organization's sensitive data. The press also reports significant thefts of Intellectual Property (IP) almost daily. The U.S. based cloud solution provider Blackbaud suffered a data breach that lasted from February to mid-May 2020, during which cyber criminals allegedly were able to exfiltrate a huge amount of data. In October 2020, cyber-criminals stole about 1TB of employee information and company documents from the German tech firm Software AG [7]. The Australian Toll Group in 2020 was hit by cyber criminals twice in three months, with an alleged data loss of over 200GB of corporate data [8]. In some cases, months might pass by before a successful compromise of an enterprise network is discovered. According to the 2021 Verizon's report [9], 20% of data breaches that occurred in 2020 were discovered several months after the attack, such as the SolarWinds cyber attack [10] that remained undetected for 9 months. Adversaries interested in a company's information could exploit the interval after intrusion and before detection to exfiltrate large amounts of IP documents from the company.

Given the vast amount of exfiltrated data, adversaries often employ a strategy of analyzing their contents to identify specific documents related to their interests. Using human domain experts is one option but it is a time-consuming activity for adversaries. Consequently, as a first step, they typically select documents related to certain topics of interest through an automated approach to focus their in-depth analysis only on a few documents. In the final phase, human domain experts come into play to assess the value of the few selected documents in terms of IP and the presence of innovative content. Our adversary model encompasses a broad range of adversaries, including Wikileaks users who possess the capability to employ topic modeling and clustering techniques. These techniques allow them to identify specific documents of interest within the vast collection published by Wikileaks, such as those containing highly sensitive information.

Manuscript received 26 May 2023; revised 19 January 2024 and 11 April 2024; accepted 10 May 2024. Date of publication 17 May 2024; date of current version 24 May 2024. This work was supported in part by the projects: Ministero dell'Università e della Ricerca (MUR) National Recovery and Resilience Plan, SEcurity and RIghts In the CyberSpace (SERICS), under Grant PE00000014; in part by the Horizon 2020 Framework Programme, Mobilization of Olive GenRes through pre-breeding activities to face the future challenges and development of an intelligent interface to ensure a friendly information availability for end users (GEN4OLIVE), under Grant 101000427; in part by AutoAD: USING ACTIVE DEFENSE TO DEFEAT CYBER ADVERSARIES, Sapienza University of Rome 2022, under Grant RG1221816C839BF9; and in part by the Office of Naval Research under Grant N00014-18-1-2670 and Grant N00014-23-1-2132. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Grigorios Loukides. (Corresponding author: Alberto Maria Mongardini.)

Alberto Maria Mongardini, Massimo La Morgia, Luigi Vincenzo Mancini, and Alessandro Mei are with the Department of Computer Science, Sapienza University of Rome, 00185 Rome, Italy (e-mail: mongardini@di.uniroma1.it; lamorgia@di.uniroma1.it; mancini@di.uniroma1.it; mei@di.uniroma1.it).

Sushil Jajodia is with the Center for Secure Information Systems, George Mason University, Fairfax, VA 22030 USA (e-mail: jajodia@gmu.edu).

Digital Object Identifier 10.1109/TIFS.2024.3402433

This paper aims to hinder the first phase of the attack using deceptive strategies [11]. To achieve this, it proposes the DARD (Deceptive Approaches for Robust Defense against IP theft) system. This system is designed to ensure that adversaries fail in their attempts to analyze a large repository of exfiltrated documents containing mainly text using automated tools. While the exfiltrated documents may include various data types beyond text, the DARD system exclusively focuses on applying deceptive operations to the text portion. The exploration of the feasibility of adapting the proposed operations to different data types, such as databases and images, is deferred to future research work.

By employing the DARD system, adversaries will be left with the only expensive option of using human domain experts to examine the entire repository thoroughly. Starting from an original repository  $\mathcal{R}$  of documents, DARD automatically generates a deceptive repository  $\mathcal{R}'$ . When an automatic clustering technique analyzes  $\mathcal{R}'$ , it produces a set of document clusters that is far away, in terms of the number of clusters and of individual documents grouped in each cluster, from what is actually present in  $\mathcal{R}$ . Specifically, to generate such a deceptive repository, this work presents four deceptive operations. A defender can use these deceptive operations to implement a defense strategy against IP thefts, hindering the use of both topic modeling and document clustering techniques. Regarding topic modeling, defenders can build a deceptive repository  $\mathcal{R}'$  that, when automatically parsed, presents topics of no interest to adversaries. This contrasts with the original repository,  $\mathcal{R}$ , which may contain topics of potential significance for the target organization. In this case, the adversaries have the following options: (1) trust the result found by the automated process; (2) attempt to reverse the deceptive operations, obtaining poor results, as shown in this paper; or (3) use human experts to identify documents of interest within  $\mathcal{R}'$ . In the case of document clustering, deceptive operations can build a new deceptive repository  $\mathcal{R}'$  in such a way that clustering techniques return clusters in which organization-relevant documents are distributed. Thus, adversaries interested in retrieving these documents must consider all clusters in their analysis. If they disregard specific clusters in their subsequent analysis, they risk neglecting pertinent documents within the excluded clusters. Retrieving the organization-relevant documents from the remaining clusters would still require human effort. A defender can combine the two strategies to deceive topic modeling and document clustering approaches, achieving higher levels of defense against IP thefts and effectively slowing down the adversaries and requiring increased effort on their part.

Legitimate users can transparently access deceptive repositories using a secure enclave-based architecture. When a legitimate user requests access, the secure enclave facilitates the restoration of the original document through the mapping of deceptive-original keywords. Due to the sensitive nature of this keyword mapping, it cannot be stored in the main file systems, as there is a potential risk of exfiltration along with other documents. To address this concern, Sec. VI introduces a Secure Enclave solution [12]. In this solution, the keyword mapping is stored in the dedicated Secure Memory, and the original document is restored using the Secure CPU, ensuring complete isolation.

The contributions of this work include:

- **Deceptive operations:** We designed and implemented four deceptive operations (Basic Shuffle, Shuffle

increment, Shuffle reduction, and Change topic) that select and replace some keywords present in the documents of the repository  $\mathcal{R}$  with deceptive keywords. These operations can be used to create a deceptive repository  $\mathcal{R}'$  that, when automatically parsed, results in a different number of clusters than those in the repository  $\mathcal{R}$  and produces new clusters containing documents initially belonging to different topics of  $\mathcal{R}$ .

- **Extensive experimentation:** The deceptive operations have been applied to a repository made of real papers collected through the Arxiv APIs. We evaluate the performance of three kinds of adversaries on an experimental repository and show that the adversaries cluster the documents as planned by the defender.
- **Topic modeling and commercial tool evaluation:** We evaluate the possibility of deceiving the adversaries on the actual topics covered within a deceptive repository. We find that the first 10 keywords by relevance retrieved by topic modeling algorithms in the deceptive repository are all deceptive keywords. This finding indicates that defenders can manipulate the topics retrieved by adversaries, presenting them with believable yet fake topics. Furthermore, we test the effectiveness of DARD against adversaries using commercial tools like Amazon Comprehend [13] and find that the adversaries were only able to retrieve deceptive keywords. This result underlines the effectiveness of the DARD system, even in the face of adversaries using commercial tools.

## II. ANALYSIS OF AN EXFILTRATED REPOSITORY

Assuming that adversaries have managed to exfiltrate a company's original repository, the purpose of this section is to show an example of how such adversaries could automatically infer the topics covered by each document in the exfiltrated repository and then select only the documents they are interested in. For simplicity, here we assume that the victim company has not adopted deceptive techniques in document production, and the adversaries are not aware of any of the topics covered by the documents of the repository. More powerful attack models will be defined in Sec. IV and evaluated in the experiments in Sec. V. This section assumes that the adversaries will follow the methodology defined in Sec. II-B since it represents the classical approach for document clustering and topic modeling tasks. Indeed, this pipeline is also used as a benchmark by other important proposals for new document clustering and topic modeling techniques described in the literature [14], [15].

### A. The Repository

The exfiltrated repository presented in this subsection is also used in the experiments in Sec. III, Sec. IV, and Sec. V. This repository is a collection of 450 scientific papers, evenly divided into three different topics of computer science: Artificial Intelligence (AI), Database (DB), and Cryptography and Security (CR). The repository contains papers retrieved from ArXiv [16], an open-access archive for scholarly articles. ArXiv provides APIs<sup>1</sup> that allow users to retrieve documents specifying the domain (Computer Science), and a domain-related field (namely: Artificial Intelligence, Database, and Cryptography and Security). The documents

<sup>1</sup>[http://export.arxiv.org/api/query?search\\_query=query](http://export.arxiv.org/api/query?search_query=query)

in the repository are in Portable Document Format (PDF) and contain an average of 7,826 words each. The smallest document has 1,227 words, while the largest one 57,169. In the following, we refer to this repository as  $\mathcal{R}_d$ .

### B. Clustering the Documents and Retrieving Topics

Since adversaries know neither the exact number of clusters nor the topics covered by the repository, they want to discover both automatically and then focus their in-depth analysis only on documents related to topics of their interest. In the first automatic phase, the adversaries can use document clustering and topic modeling techniques. Document clustering and topic modeling are two data mining techniques used to automatically organize and retrieve information from unorganized collections of text documents. The goal of document clustering [17] is to organize a repository of documents into groups of similar documents. Instead, topic modeling techniques [18], [19] aim to build a latent semantic representation of the documents, detecting keywords that describe the subject dealt with by the documents. In particular, the latent semantic representation of a set of documents is called Topic. In the following sections, we describe the steps the adversaries should perform on the exfiltrated repository to retrieve the documents of their interest.

1) *Text Pre-Processing and Feature Extraction*: Before starting the analysis, the text has to be normalized and cleaned of all the elements that do not provide information about the topic (e.g., numbers). To this end, the pre-processing phase is a key component of every text classification tool [20]. Hence, the adversaries perform on the documents standard pre-processing operations such as tokenization, stemming, normalization of the upper and lowercase, and deletion of number and symbol characters. Once the documents in  $\mathcal{R}_d$  have been normalized, the adversaries proceed with the feature extraction. In text analysis, a document and its content are usually represented as a vector, where each position of the vector represents a term (i.e., one or more consecutive words in the document) with an associated weight.

In the feature extraction step, the adversaries extract the terms that occurred within the documents and assign them a weight through TF-IDF. The TF-IDF (Term Frequency-Inverse Document Frequency) [21] is a function that assigns a weight to a term in relation to a document. The greater the weight, the greater the importance of the term for the document. The idea behind the TF-IDF is to give more importance to terms that occurred within a document but are generally not frequent within the document repository. Therefore, terms that are characteristic only of a group of documents are considered significant.

By calculating the TF-IDF for each term, the adversaries obtain a TF-IDF matrix as the one in Fig. 1. Each column represents a document with a Document Vector containing the weights of the terms for that document. Instead, each row indicates a term with a Word Vector containing the weights of that term for each document in the repository.

2) *Document Clustering*: At this point, the adversaries are ready to group the documents according to the features extracted in the previous step. First, they need to estimate the correct number of clusters in the repository, which is one of the major challenges in cluster analysis [22]. The most popular approach proposed in the literature is internal clustering [23]. This method typically involves three steps: (1) apply to the dataset several clustering algorithms using

	Doc 1	Doc 2	Doc 3	Doc 4	Doc 5	Doc 6
Term 1	0	2	5	3	2	4
Term 2	4	1	2	6	0	1
Term 3	1	3	1	8	7	3
Term 4	2	3	4	9	1	1
Term 5	0	1	2	5	1	3
Term 6	6	7	0	1	3	2
Term 7	2	4	2	7	6	0
Term 8	0	2	1	5	4	4

↑  
Document vector

← Word vector

Fig. 1. Matrix representation of documents: each column identifies a document, while each row represents a term. A Document Vector is the column associated with a document and contains the weights of the terms for that document. A Word Vector is a row related to a term and contains its weights for each document in the repository.

different combinations of parameters, (2) compute the corresponding internal validation score for each obtained partition, and (3) detect the optimal number of clusters by choosing the partition with the best internal validation score. There are over thirty typologies of internal clustering evaluation [23] that can be used in steps 2 and 3 described above. The most commonly used metrics for assessing the quality of document clustering are: Silhouette Coefficient [24], the Calinski-Harabasz Index [25], and the Davies-Bouldin Index [26].

The **Silhouette Coefficient** is a measure that considers both cohesion (how close the objects are within the same cluster) and separation (how well-separated a cluster is from the nearest cluster) in a given clustering. It utilizes a distance metric (e.g., Manhattan distance, Euclidean Distance) to quantify these aspects. The measure range from  $-1$  to  $+1$ , where a higher value indicates that the data belonging to the same clusters are close to each other, and well separated from the point of the other clusters.

The **Calinski-Harabasz Index**, also known as Variance Ration Criteria (VARAC), is the ratio of the sum of the squared distances of the centroids (between-clusters dispersion) to the sum of the squared distances of the points from their centroid (inter-cluster dispersion). For the Calinski-Harabasz index, the higher the score, the more well-separated the clusters are.

The **Davies-Bouldin Index** provides a quantitative measure of how well-separated and internally cohesive the clusters are in a clustering solution. It is calculated by taking the average similarity ratio of each cluster with its most similar cluster, where similarity is defined as the ratio of intra-cluster and inter-cluster distances. The intra-cluster distance is the sum of the distances of the points belonging to that cluster from the centroid of the cluster. Conversely, the inter-cluster distance is the sum of the distances of the centroids. Differently from the previous cases, here, a lower score indicates better separations of clusters.

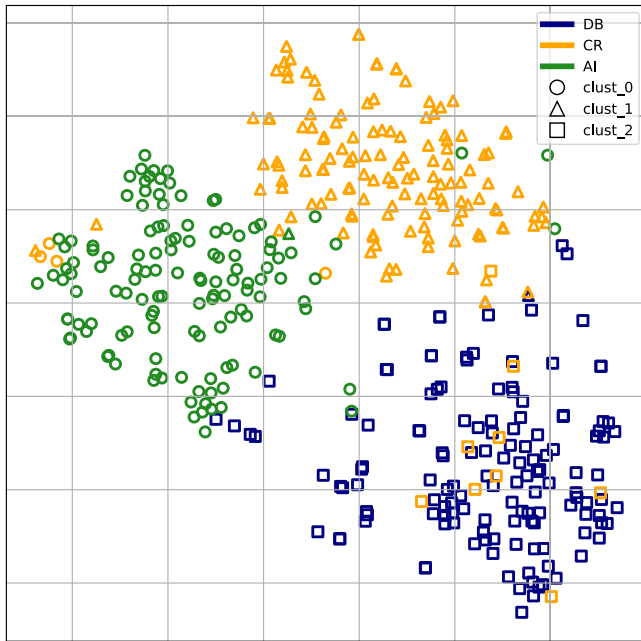


Fig. 2. Projection of the clusters obtained through K-means from the original repository composed of documents related to Artificial Intelligence (AI), Database (DB), and Cryptography and Security (CR).

Since these metrics have been extensively adopted in the literature, we assume adversaries rely only on these three internal validation scores to infer the number of clusters  $K$  in the repository.

Once detected the number of  $K$  clusters (three in this case), the adversaries rely on the TF-IDF weighting scheme and K-means to cluster the documents. K-means [27] is a popular clustering algorithm that takes as input a number  $k$  of expected clusters and finds a  $k$ -partition such that the squared error between the empirical mean of a cluster (centroid) and the points in the cluster is minimized.

Applying K-means on the document vectors of the TF-IDF weights, the adversaries obtain clusters that correctly group the exfiltrated documents according to their topics, as shown by their projection in Fig. 2. To visualize the clusters obtained by the adversaries in this section and the next ones, we performed a dimensionality reduction on the TF-IDF weights by applying the t-distributed stochastic neighbor embedding (t-SNE) [28], an algorithm that allows visualizing high dimensional data in a low dimensional space. The color of each item in the figure represents the original topic of the document, whereas the shape of the item represents the cluster the document belongs to. As we can see, the clustering result is remarkably similar to the real one. Therefore, the adversaries are able to correctly distinguish documents belonging to the three different topics.

3) *Topic Modeling*: In the previous section, adversaries cluster the documents according to the TF-IDF weights. Here, they want to infer the topic covered by each of the  $K$  clusters by retrieving the terms that describe each specific cluster. To this end, adversaries leverage Latent Dirichlet Allocation (LDA) [19], one of the most used topic modeling algorithms, that provides as output for each cluster a list of terms ordered by their relevance to the topic. We define as  $M$  keywords, the first  $M$  terms of the output list representing the topic of a given cluster. Tab. I shows the top 10 keywords extracted

TABLE I  
TOP 10 KEYWORDS EXTRACTED FROM EACH CLUSTER USING LDA

AI	DB	CR
learn	query	scheme
plan	node	protocol
agent	object	security
decision	logic	message
policy	xml	signature
action	attribute	public
network	tree	attack
constraint	semantic	service
strategy	update	random
intelligence	predicate	bit

from each cluster in the repository  $\mathcal{R}_d$  by the adversaries. The latter might infer from these keywords the three topics covered in  $\mathcal{R}_d$ , namely: Artificial Intelligence, Database systems, and Cryptography and Security. After this step, adversaries focus their analysis only on documents addressing specific topics in which they are interested. The paper does not cover the second phase since our techniques aim to deceive the first phase, of which results also influence the second one. Due to the proposed deceptive operations, the documents covering the topic of interest for the adversaries will be scattered throughout all clusters. As a result, adversaries can not focus on just one cluster but on all of them.

### III. OPERATIONS

#### A. Replacements of Terms

This subsection describes the idea behind the deceptive operations, or the term-replacement operations, illustrating the relationship between the term-replacement operations and the resulting changes in the TF-IDF matrix calculated on the sets of documents in  $\mathcal{R}$ . This paper refers to *keyword*  $k$  as the term to be replaced and *deceptive term*  $dk$  as a new term, not contained in  $\mathcal{R}$ , that replaces one or more keywords  $k$ .

To explain the effect of a term-replacement operation, we rely on the concepts of centroid and distance between centroids. Let  $T$  be the set of terms contained in the documents of  $\mathcal{R}$ ,  $S$  be a set of documents in  $\mathcal{R}$ , and consider a sub-matrix of the TF-IDF weights of  $\mathcal{R}$  that contains only the columns representing the documents in  $S$ . We define as the centroid of  $S$  the vector that contains the element-by-element average of the rows in this submatrix. Let  $S_1$  and  $S_2$  be two sets of documents, the distance between  $S_1$  and  $S_2$ , denoted  $d(S_1, S_2)$ , is the Euclidean distance between their centroids.

Given a repository of documents  $\mathcal{R}$ , this paper considers the following four strategies to replace keywords at the level of the documents set, where all occurrences of a certain keyword are replaced inside all the documents contained in a specific set.

(i) **1-to-1 replacement**: Consider a repository  $\mathcal{R}$  partitioned in  $n > 1$  sets of documents, such that  $\mathcal{R} = \{S_1 \cup \dots \cup S_n\}$ . Let  $k$  be a keyword for the repository  $\mathcal{R}$ , and  $dk$  a deceptive term. The 1-to-1 replacement operation changes all the occurrences of  $k$  in all the documents of the repository  $\mathcal{R}$ , with the deceptive term  $dk$ .

After the 1-to-1 replacement, the deceptive term  $dk$  appears in the same documents and with the same frequencies of  $k$  (Fig. 4(a)). Thus, there is a new row in the TF-IDF matrix of  $\mathcal{R}$  for  $dk$ , which has precisely the same weights as  $k$ . Moreover, since the keyword  $k$  no longer appears in the documents of  $\mathcal{R}$ ,

the row in the TF-IDF matrix associated with  $k$  disappears as well. Hence, the 1-to-1 replacement does not alter the relative position among the centroids of all the sets of documents  $S_i$ . In addition, since  $k$  was a keyword for  $\mathcal{R}$ , also  $dk$  will be a keyword for the deceptive repository  $\mathcal{R}'$ .

(ii) **1-to-N replacement:** Consider a repository  $\mathcal{R}$  partitioned in  $n > 1$  sets of documents, such that  $\mathcal{R} = \{S_1 \cup \dots \cup S_n\}$ . Let  $k$  be a keyword for  $\mathcal{R}$ , and  $\{dk_1, \dots, dk_n\}$  be a set of deceptive terms. The 1-to-N operation replaces all the occurrences of  $k$  with a different deceptive term  $dk_i$  in each document of  $S_i$ . Thus, for every  $i$ , after the 1-to-N replacement, the term  $dk_i$  appears in the documents of  $S_i$  instead of the keyword  $k$  and  $dk_i$  does not appear in the documents of  $\mathcal{R} \setminus S_i$  (Fig. 4(b)).

Note that after the 1-to-N replacement, the keyword  $k$  no longer appears in the documents of  $\mathcal{R}$  and, consequently, in its TF-IDF matrix. At the same time, after the replacement, in the TF-IDF matrix  $n$  new rows appear, one for each deceptive keyword  $dk_i$ . Finally, since the deceptive term  $dk_i$  appears only in the documents of the set  $S_i$ , its weight will be greater than zero in the documents that belong to  $S_i$  and zero for the others. Hence, the centroid of each set  $S_i$  tends to move away from the centroids of the sets  $S_l$  for each  $i, l$  with  $i, l \in \{1, \dots, n\}$  and  $i \neq l$ . In particular, the higher the rank of keyword  $k$  is, the more the centroids tend to move away from each other.

(iii) **N-to-1 replacement:** Consider a repository  $\mathcal{R}$  partitioned in  $n > 1$  sets of documents, such that  $\mathcal{R} = \{S_1 \cup \dots \cup S_n\}$ . Let the keywords  $\{k_1, \dots, k_n\}$  be a set of terms, such that  $k_i$  is a keyword for the set of documents  $S_i$ , while  $k_i$  is not a keyword for  $S_l$ , with  $i, l \in \{1, \dots, n\}$  and  $i \neq l$ . The N-to-1 operation replaces in every set of documents  $S_i$  all the occurrences of the keyword  $k_i$  with the deceptive keyword  $dk$ . Thus, after the N-to-1 replacement, the deceptive term  $dk$  appears in the documents of  $S_i$  instead of  $k_i$ , for every  $i$  (Fig. 4(c)). In the N-to-1 replacement, the goal is to bring closer the centroids of the set of documents  $S_i$ , replacing  $N$  different keywords with the same deceptive keyword  $dk$ . Differently, the 1-to-N replacement aims to move away the centroids of the set of documents  $S_i$ , replacing a unique keyword with  $N$  different deceptive keywords. Following the N-to-1 replacement, the TF-IDF weights of all the keywords  $k_i$  drop to zero for the documents in  $S_i$ , whereas a new row associated to the deceptive keyword  $dk$  appears in the TF-IDF matrix of  $\mathcal{R}$ . The TF-IDF weight of  $dk$  is greater than zero for all the documents in  $S_i$  that previously contained the keyword  $k_i$ . Hence, the centroid of each set  $S_i$  tends to get closer to the centroid of  $S_l$ , for every  $i, l$ . In particular, the higher the rank of the keywords  $k_i$  is, the more the centroids tend to get closer to each other.

(iv) **N-to-N replacement:** Consider a repository  $\mathcal{R}$  partitioned in  $n > 1$  sets of documents, such that  $\mathcal{R} = \{S_1 \cup \dots \cup S_n\}$ . Let  $\{k_1, \dots, k_n\}$  be a set of keywords, such that each  $k_i$  is a keyword for one set of documents  $S_i$ , while  $k_i$  is not a keyword for  $S_l$ , with  $i, l \in \{1, \dots, n\}$  and  $i \neq l$ , and  $dk_1, \dots, dk_n$  be a set of deceptive terms. The N-to-N operation replaces in every set of documents  $S_i$  all the occurrences of the keyword  $k_i$  with the deceptive keyword  $dk_i$ . Thus, after the N-to-N replacement, the deceptive term  $dk_i$  appears in the documents of  $S_i$  instead of  $k_i$ , for every  $i$  (Fig. 4(d)). The N-to-N replacement is similar to a 1-to-1 replacement applied to a single set of documents  $S_i$ , instead of to all the repository  $\mathcal{R}$ . After an N-to-N replacement, the deceptive term  $dk_i$  will

have in the TF-IDF matrix of  $\mathcal{R}$ , the same TF-IDF weights of the keyword  $k_i$ . Hence, since  $k_i$  was a keyword for the set  $S_i$ , also  $dk_i$  will be a keyword after the N-to-N replacement.

The previous paragraph describes the behavior of the four different replacement operations when executed once. However, applying the same operation several times on the repository is possible. Thus, in the following, we will refer to **m-multiple** replacement the use of the same operation  $m$  times on the repository. Applying  $m$  times a 1-to-1 or a 1-to-N replacement means that  $m$  different keywords will be replaced with  $m$  or  $m \times N$  deceptive keywords, respectively. Whereas applying  $m$  times the N-to-1 or N-to-N replacement means that  $m \times N$  different keywords will be replaced with  $m$  or  $m \times N$  deceptive keywords, respectively.

## B. Shuffle Clusters Operations

Starting from a repository  $\mathcal{R}$  where each cluster is made of documents that belong to a single topic, the Shuffle operation builds a deceptive repository  $\mathcal{R}'$  in which some or all the clusters contain documents of different topics. Thus, the adversaries cannot precisely cluster the documents of  $\mathcal{R}'$  according to the original topics of  $\mathcal{R}$ . The Shuffle operation, given a set  $L$  made of  $l$  different clusters, partitions each cluster in  $L$  into  $p$  subsets of documents and mixes these subsets among themselves, building a new set of clusters  $L'$ . In particular, each new cluster in  $L'$  is made of  $l$  subsets of documents each of which belongs to a different original cluster of the set  $L$ . Since each new cluster contains exactly one subset of each original cluster, the relationship between  $p$  and  $l$  determines the number of resulting new clusters in  $L'$ .

*Definition 1:* Shuffle  $(\mathcal{R}, C_1, \dots, C_l) \Rightarrow \mathcal{R}'$ .

Given a Repository  $\mathcal{R}$  composed of  $n > 1$  clusters, let  $\mathcal{CR}$  be the set of clusters in  $\mathcal{R}$ . Consider the clusters  $L = \{C_1, \dots, C_l\}$  in  $\mathcal{CR}$ , with  $1 < l \leq n$ . The Shuffle operation partitions each cluster  $C_i$  of  $L$  into  $p$  subsets of documents, with  $l - 1 \leq p \leq l + 1$  and  $l - 1 \geq 2$ , such that  $C_i = s_{i,1} \cup \dots \cup s_{i,p}$ , where  $s_{i,j}$  represents the subset  $j$  of the cluster  $C_i$ . The Shuffle operation replaces some keywords in  $L$  in such a way as to form  $p$  new clusters  $C'_i$ . Each new cluster  $C'_j$  is composed of  $l$  subsets of documents  $s_{i,j}$  in such a way that  $C'_j = s_{1,j} \cup \dots \cup s_{l,j}$  with  $j \in \{1, \dots, p\}$ . Depending on the relationship between  $p$  and  $l$ , the effect of the Shuffle operation on the repository  $\mathcal{R}'$  is different. In particular, there are three possible variants: The **Basic Shuffle** in which the number of partitions  $p$  is equal to  $l$  and thus the number of clusters in the repository  $\mathcal{R}'$  is  $n$ , the **Shuffle Increment** where the number of partition  $p$  is equal to  $l + 1$ ; in this case  $\mathcal{R}'$  contains  $n + 1$  clusters, and the **Shuffle Reduction** in which the number of partitions  $p$  is equal to  $l - 1$ , and the repository  $\mathcal{R}'$  contains  $n - 1$  clusters. After one of the three Shuffle operations, the adversaries that search for  $n - l + p$  clusters in the repository  $\mathcal{R}'$  will find the following set of clusters:  $\mathcal{CR}' = (\mathcal{CR} \setminus L) \cup L'$ , where  $L' = \{C'_1, \dots, C'_p\}$ .

In our implementation, the Shuffle operations first compute the keywords of all the clusters  $C_i$  in  $L$ . Then, it partitions the documents of each  $C_i$  into  $p$  subsets of documents such that  $C_i = s_{i,1} \cup \dots \cup s_{i,p}$ , with  $s_{i,j}$  that represents the subset  $j^{th}$  of the cluster  $C_i$ . The Shuffle operations select subsets of documents to compose the new clusters  $\{C'_1, \dots, C'_p\}$ , such that each new cluster  $C'_h$ , with  $h \in \{1, \dots, p\}$ , is made of  $l$  subsets of documents, one subset from each  $C_i$ . For the sake of simplicity, we assume that the Shuffle operations select the

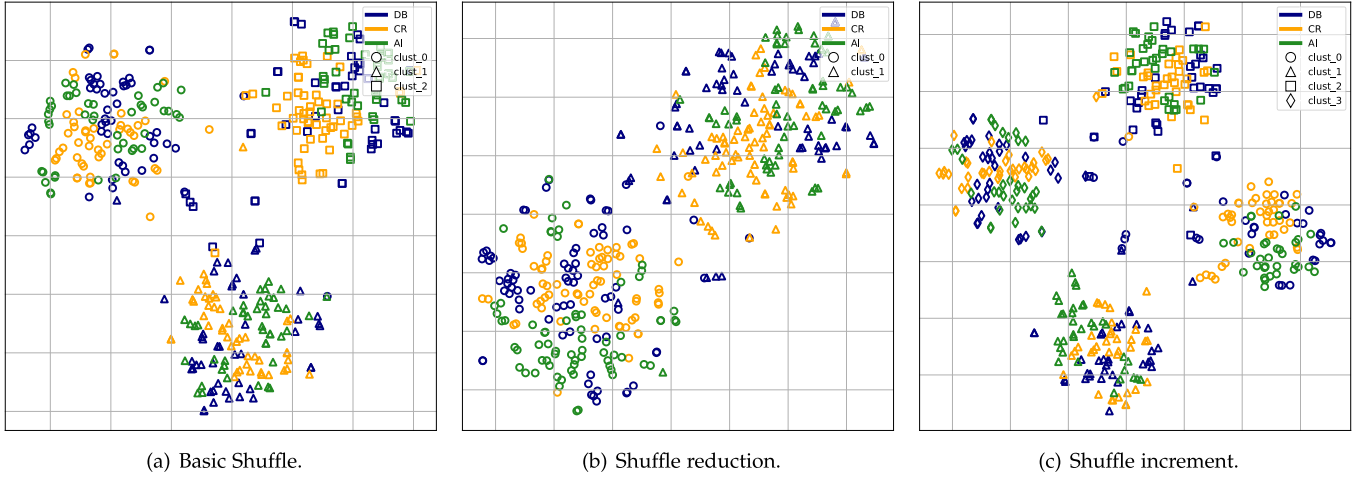


Fig. 3. Fig. 3(a) shows the projection of the repository (composed by AI, DB, and CR) modified by the Basic Shuffle operation. Fig. 3(b) and Fig. 3(c) show, respectively, the projections of the repository obtained by applying Shuffle reduction and Shuffle increment.

subsets of documents that compose the new cluster  $C'_h$  picking from each cluster  $C_i$  the  $h^{th}$  subset, thus  $C'_h = s_{1,h} \cup \dots \cup s_{l,h}$ . Finally, for each cluster  $C'_h$ , the Shuffle operations perform an  $l$ -to-1 replacement, which overall sums up to  $p$  times a  $l$ -to-1 replacement operations. An  $l$ -to-1 operation replaces  $l$  different keywords, each one computed on each  $C_i$ ,  $i \in \{1, \dots, l\}$  with the deceptive term  $dk_h$  (see Sec. III-D for details about the keywords selection). In particular, let  $k_i$  be a keyword of  $C_i$ . The  $l$ -to-1 operation replaces all the occurrences of  $k_i$  in the subset  $s_{i,h}$  with the deceptive term  $dk_h$ , with  $h \in \{1, \dots, p\}$  for every  $i \in \{1, \dots, l\}$ . Note that each of the  $h$  execution of the  $l$ -to-1 operation uses a different deceptive term  $dk_h$ .

A single  $l$ -to-1 replacement could not be enough to bring the centroids of all the subset  $\{s_{1,h}, \dots, s_{l,h}\}$  sufficiently closer to form the new cluster  $C'_h$ , for every  $h \in 1, \dots, p$ . Therefore, the Shuffle operation has to perform  $m$ -multiple  $l$ -to-1 replacements such that the following equation is satisfied:

$$\begin{aligned} d(s_{i,h}, (C'_h \setminus s_{i,h})) &< d(s_{i,h}, (C_i \setminus s_{i,h})) \\ \forall i \in \{1, \dots, l\} \wedge \forall h \in \{1, \dots, p\} \end{aligned} \quad (1)$$

The formula verifies that after each iteration of  $l$ -to-1 replacement, the centroid of  $s_{i,h}$  is closer to the centroid of  $C'_h \setminus s_{i,h}$  (left term of the formula) than to the one of  $C_i \setminus s_{i,h}$  (right term). In this way, each pair  $(C'_h \setminus s_{i,h}, s_{i,h})$  is close enough to build the new cluster  $C'_h$ , and the subsets of  $C_i$  will not cluster together.

Fig. 3(a) 3(b) 3(c) show the deceptive repository  $\mathcal{R}'_d$  after we applied on the repository  $\mathcal{R}_d$  the Basic Shuffle (Fig. 3(a)), the Shuffle Increment (Fig. 3(b)) and the Shuffle Reduction (Fig. 3(c)). To perform the three operations, we insert into the set  $L$  all the clusters of the repository  $\mathcal{R}_d$  and set the value of  $p$  as 2, 3, 4 respectively for the Shuffle Reduction, the Basic Shuffle, and the Shuffle Increment. After the operations, each cluster (denoted in the figures by the circle, square, diamond, and triangle markers) is made of a mixture of topics (green, blue, and yellow markers).

### C. Change Topic Operation

The Change Topic operation aims to change the original topic of a cluster of documents  $C_t$  in  $\mathcal{R}'$ . The Change Topic

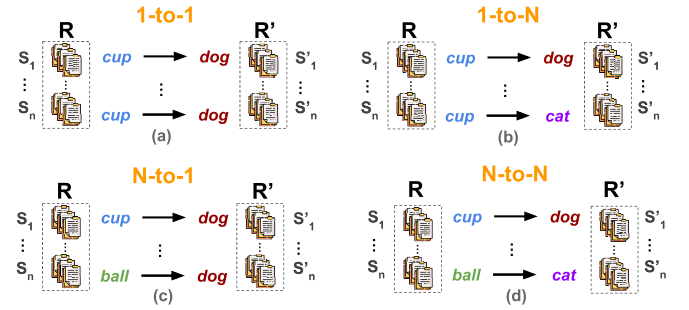


Fig. 4. Representation of the 1-to-1 replacement (a), 1-to-N replacement (b), N-to-1 replacement (c), and N-N replacement (d).

operation builds a repository  $\mathcal{R}'$  replacing several keywords of  $C_t$  with a set of deceptive terms  $\{dk_1, \dots, dk_l\}$ , in such a way that topic modeling performed on  $\mathcal{R}'$  returns a topic that depends on the deceptive terms, and such a topic can be different from the original one.

**Definition 2:**  $\text{Change Topic}(\mathcal{R}, C_t, \{dk_1, \dots, dk_l\}) \Rightarrow \mathcal{R}'$ . Given a repository  $\mathcal{R}$  that contains  $n > 1$  clusters, a target cluster  $C_t$  in  $\mathcal{R}$  and a set of  $l$  deceptive terms  $\{dk_1, \dots, dk_l\}$ , the Change Topic operation replaces  $l$  keywords with  $l$  deceptive terms, using one different deceptive term for each different keyword. At the end of the operation, the adversaries that perform topic modeling on the repository  $\mathcal{R}'$  will find for the cluster  $C_t$  the following keywords  $\{dk_1, \dots, dk_l\}$ .

In our implementation, the Change Topic operation computes the keywords of  $C_t$  and ranks them by their TF-IDF weight. Let  $\{k_1, \dots, k_l\}$  be the first  $l$  keywords of  $C_t$  in the rank. The Change Topic operation performs a 1-to-1 replacement on the documents of  $C_t$ , replacing all the occurrences of the keyword  $k_i$  with the deceptive term  $dk_i$ , for every  $i \in \{1, \dots, l\}$ . Overall, the Change Topic operation performs  $l$  times a 1-to-1 replacement on each document of  $C_t$ .

After the operation, the deceptive term  $dk_i$  is a keyword for the cluster  $C_t$  in  $\mathcal{R}'$ . Indeed,  $dk_i$  has the same TF-IDF weight in  $\mathcal{R}'$  as  $k_i$  has in  $\mathcal{R}$ . Thus, since  $k_i$  is a keyword for  $C_t$  in  $\mathcal{R}$ ,  $dk_i$  is a keyword for  $C_t$  in  $\mathcal{R}'$  as well. Moreover, since  $k_i$  and  $dk_i$  have the same weight in the TF-IDF of  $\mathcal{R}$  and  $\mathcal{R}'$ , respectively, the centroid of  $C_t$  is the same both in  $\mathcal{R}$  and in  $\mathcal{R}'$ .

TABLE II  
NUMBER OF TERMS TO INVOLVE FOR DECEPTIVE OPERATIONS  
SELECTING SUBSETS OF DOCUMENTS RANDOMLY OR BY  
GROUPING SIMILAR ONES AND REPLACING  
RANDOM TERMS

	B. Shuffle	Shuffled In.	Shuffle Red.
Random selection	8,6	10,2	7,4
Constrained K-means	8	10	6
K-means + random terms	71,6	86,6	41,9

#### D. Observations on the Keywords and the Selection of the Documents

This section describes some of the possible approaches to partition a cluster  $C$  of documents of  $\mathcal{R}$  into subsets suitable to be used by the deceptive operations described in the previous sections. In addition, we present the criteria we used to select the keywords to be replaced with the deceptive terms. When partitioning a cluster  $C$  to apply one of the deceptive operations, there are two main aspects to face: the number of documents each partition should contain and which documents of  $C$  should be grouped in the same partition.

The number of documents each partition is made is a crucial parameter to decrease the purity [29] of the resulting repository  $\mathcal{R}'$ . Purity is an external evaluation metric that assesses the quality of given clusters by indicating the percentage of the total number of correctly classified objects (documents). For instance, in the Shuffle operation, creating partitions with roughly equal numbers of documents leads to creating new clusters in  $\mathcal{R}'$  with a purity roughly equal to zero, which guarantees the greatest possible deception. In the following experiments, all the clusters in our test repository  $\mathcal{R}_t$  have the same number of documents. Hence, given the observations above, the best strategy in our case is to create the partitions used in each deceptive operation with the same number of documents. The second aspect to face is which documents of  $C$  should be placed in the same partitions. A trivial approach is to partition the documents of a cluster  $C$  in subsets of documents  $\{s_1, \dots, s_l\}$  through a random selection of the documents in  $C$ . This approach likely leads to group into the same subset documents uniformly spread among the cluster  $C$ , with the centroid of each subset  $s_i$  near the centroid of  $C$  and thus close to each other. However, the closer the centroids of the subsets are to each other, the more keywords the cluster operations need to replace in order to push the centroids away among them (See Tab. II). A better choice is to partition the documents so that the centroids of the subsets  $s_i$  result far away among them. An approach to generate such subsets  $s_i$  is to leverage a clustering algorithm, such as K-means. Since standard K-means may generate partitions with an unbalanced number of documents (*e.g.*, a partition with most of the documents and others with very few documents), we used the constrained version of K-means [30]. The constrained version of K-means extends the classic clustering algorithm by adding constraints on data point assignments. These constraints avoid local solutions with empty clusters or clusters having very few points. Moreover, they ensure that each partition has a roughly equal number of documents.

Concerning the selection of the keywords to be replaced with the deceptive terms, we select the keywords by their TF-IDF weights in descending order. This approach minimizes the number of deceptive keywords to be replaced to accom-

plish any of the cluster operations. Indeed, the effectiveness of the 1-to-N replacement and the N-to-1 replacement in pushing away or bringing close among them the centroids of partitions  $\{s_1, \dots, s_l\}$  is proportional to the TF-IDF weight of the replaced keywords (as discussed in Sec. III-A).

To better understand how the documents and the keywords selection affect the number of keywords needed to perform a deceptive operation, we evaluated the deceptive operations in the following three settings: partitions created with a random selection of the documents and keywords selected by TF-IDF weight; partitions created leveraging the constrained K-means and keywords selected by TF-IDF weight; and partitions created leveraging the constrained K-means and keywords selected randomly. For each of the above settings (except the constrained k-means version), we repeat the experiment 10 times and compute the average number of keywords needed to perform the cluster operations. Tab. II shows the results of this experiment. The best combination to minimize the number of keywords replaced is the one based on the constrained k-means and the keyword selected by TF-IDF weight (constrained k-means in the table). This is in line with our previous observations in this section. Randomly selecting the keywords increases the number of keywords drastically to be replaced. For example, in the case of the Shuffle Increment operation, the number of keyword replacements increases from about 10 to more than 80. Building the partitions by randomly selecting the documents requires a few more replacements than partitioning the documents via k-means.

#### E. Deceiving the Number of Topics

An accurate clustering result requires the right estimation of the number of clusters in a repository of documents. By our assumption, the adversaries that exfiltrated the repository  $\mathcal{R}'$  do not know the number of clusters that the repository contains. Thus, they have to estimate the number of clusters in the repository  $\mathcal{R}'$  through internal cluster indices (see Section II-B.2). This section aims to illustrate our proposed technique to deceive adversaries in such estimation, making them believe that the repository  $\mathcal{R}'$  contains a given numerical value  $K_d$  for the number of clusters which is deceptive.

In the literature, there are several internal cluster indices that the adversaries can leverage to estimate the number of clusters of  $\mathcal{R}'$ . The main idea behind these indices is to evaluate the compactness (how close are the items of the same cluster), the separation (how distant are the clusters from each other), or a combination of them. Every index evaluates these criteria accordingly with the different evaluation methodologies they use (*e.g.*, average distance, minimum distance, the sum of square error).

However, it is not possible to know in advance the validation indices the adversaries will use. Therefore to deceive adversaries, the clusters in  $\mathcal{R}'$  have to be enough compact and separated so that for all the indices, or at least most of them, the resulting number of clusters is  $K_d$ . In terms of our cluster operations, we propose to redefine the stopping criteria for the number of term-replacement to be performed (recall that both the 1-to-N and the N-to-1 operation contribute to pushing away or bringing close clusters among them) so that their number is greater than or equal to those defined in Eq. 1.

To evaluate the minimum number of term-replacement to deceive adversaries on the estimation of the number of clusters

in the repository  $\mathcal{R}'$ , we introduce the following function:

$$f(\mathcal{R}, Op, K_d, K_{max}, T_{max}, \mathcal{S}_{ivi}) \quad (2)$$

The function  $f$ , given a repository  $\mathcal{R}$ , a cluster operation  $Op$ , and a Set of internal validation indices  $\mathcal{S}_{ivi}$ , computes the minimum number of term-replacement operations such that all the indices in  $\mathcal{S}_{ivi}$  evaluate  $K_d$  as the estimated number of clusters. Since it is impractical to evaluate all the possible numbers of clusters, we reduce the search space of the number of clusters from 2 up to  $K_{max}$ . Finally,  $T_{max}$  represents the maximum number of term-replacement operations we are willing to perform. It is important to set  $T_{max}$  because internal validation indices, depending on how they evaluate the compactness and the separation, could cause an unlimited number of term-replacement operations when evaluating particular data distribution (e.g., presence of outliers, skewed distribution) [31].

Computing  $f$  on the repository  $\mathcal{R}_d$  for the Basic Shuffle, the Shuffle Increment, and the Shuffle Reduction operations, we find that to deceive adversaries about the number of clusters contained in  $\mathcal{R}'_d$ , for the Basic Shuffle operation, we have to perform 41 term-replacement instead of 8, 38 replacements for the Shuffle Increment instead of 10, while 18 for the Shuffle Reduction. For the above-mentioned results, we compute the function  $f$  evaluating the following indices: the Silhouette Coefficient (SIL), the Calinsky-Harabasz index (CH), and the Davies-Bouldin Index (DB), and we set as 100 the maximum number of replacement operation  $T_{max}$ . For the Shuffle operation, we set  $K_{max}$  as 9 since we divided the repository into 9 partitions, whereas  $K_d$  as 3 because we aim to make the adversaries believe that  $\mathcal{R}'_d$  contains 3 clusters. For the Shuffle Increment operation, we set  $K_{max}$  as 12, and  $K_d$  as 4. Finally, for the Shuffle Reduction operation, we set for  $K_d$  and  $K_{max}$  respectively 2 and 6.

Tab. III shows the scores of the internal validation indices computed on  $\mathcal{R}'_d$  varying the number of clusters that the adversaries are looking for. As we can see, the number of clusters estimated by the adversaries after each operation coincides with the predetermined deceptive number of clusters  $K_d$ .

#### IV. POSSIBLE ADVERSARIES

##### A. The Attack Model

This subsection defines three models, each representing an adversary with different knowledge of both the content of the repository  $\mathcal{R}'$  and the deception techniques adopted.

- **Black Box adversaries:** They are the weakest kind of adversaries we consider in this work. They are not aware of the proposed deceptive techniques and believe that the exfiltrated repository  $\mathcal{R}'$  is the original one.
- **Gray Box adversaries:** These adversaries suspect that some deceptive operations may have been executed on the repository  $\mathcal{R}'$ . Nonetheless, even though they know the deceptive operations presented in this paper, Gray Box adversaries neither know how many and which specific deception operations were performed, nor how many and which deceptive keywords have been used to perform each operation.
- **Enhanced Gray Box adversaries:** They have the same knowledge as the Gray Box adversaries. However, these adversaries also leverage the **Oracle Function** to obtain an ordered list of terms in  $\mathcal{R}'$  that may have been replaced by the deceptive operations (details in Sec. IV-A.1). The

TABLE III

SCORES OF THE DAVIES-BOULDIN INDEX (DB), CALINSKI-HARABASZ INDEX (CH), AND SILHOUETTE COEFFICIENT (SIL) BASED ON THE NUMBER OF CLUSTERS SEARCHED FOR

		Estimated Number of clusters				
		2	3	4	5	6
Shuffle	DB	4.13	<b>3.76</b>	3.86	4.25	4.50
	CH	19.54	<b>20.45</b>	16.26	13.90	12.52
	SIL	0.044	<b>0.063</b>	0.063	0.046	0.038
Shuffle-Incr	DB	4.36	4.07	<b>3.55</b>	3.62	3.57
	CH	17.16	17.90	<b>19.50</b>	15.26	12.77
	SIL	0.040	0.056	<b>0.070</b>	0.062	0.059
Shuffle-Red	DB	<b>3.90</b>	4.68	4.59	4.43	4.76
	CH	<b>23.17</b>	16.01	12.97	11.31	10.15
	SIL	<b>0.052</b>	0.037	0.030	0.032	0.032

ability to invoke the Oracle Function makes Enhanced Gray Box adversaries the strongest. These adversaries represent the most challenging scenarios for evaluating the effectiveness of the proposed deceptive operations. It is crucial to emphasize that the Enhanced Gray Box adversaries include those with the ability to foresee potential deceptive keywords within the documents. Moreover, as there is no real **Oracle Function** available, these adversaries serve as the worst-case scenario, highlighting the strength of our deceptive operations when faced with adversaries capable of predicting some potential deceptive keywords.

We assume that Gray Box and Enhanced Gray Box adversaries can remove deceptive keywords from the repository  $\mathcal{R}'$ , as discussed in Sec. IV-B. In addition, we assume that all the adversaries described in this section: cannot access the mapping of the keywords replacements, use K-means as the clustering algorithm and the Silhouette Coefficient, the Calinsky-Harabasz Index, and the Davies-Bouldin Index to estimate the number of clusters contained in the exfiltrated repository. It is important to note that none of the three adversaries can self-estimate metrics such as the Purity or the ARI on their clustering since they do not know the ground truth of the repository  $\mathcal{R}$ .

1) *The Oracle Function:* Knowing which are the deceptive keywords in the repository  $\mathcal{R}'$  may be a great advantage for the adversaries. Indeed, leveraging this information, they can eliminate those terms to subvert the operations.

In this section, we define the Oracle Function to emulate adversaries that somehow gained access to the list of terms and thus are able to select the deceptive terms we used to build the deceptive repository.

*Definition 3:* Oracle ( $\mathcal{R}'$ )  $\Rightarrow L_k$ .

Consider the repository  $\mathcal{R}'$  made of  $M$  keywords, such that  $D$  keywords are all and only the deceptive keywords used to generate the repository  $\mathcal{R}'$ , and the remaining  $M - D$  terms are the original keywords that are both in  $\mathcal{R}$  and  $\mathcal{R}'$ . The Oracle Function takes as input a deceptive repository  $\mathcal{R}'$  and returns as output an ordered list of  $M$  keywords  $L_k = \{dk_1, \dots, dk_D, k_{D+1}, \dots, k_M\}$ , where the first  $D$  items of the list are deceptive keywords, while the remaining  $M - D$  items are unchanged keywords. In particular, the deceptive keyword  $dk_i$ , with  $i \in \{1, \dots, D\}$ , represents the  $i^{\text{th}}$  deceptive keyword used to generate the repository  $\mathcal{R}'$ . The remaining  $M - D$  keywords are ordered as they were the next



terms to be replaced by the operation used to generate the repository  $\mathcal{R}'$ .

Although Enhanced Gray Box adversaries, through the Oracle Function, can access in an ordered way all the deceptive keywords, they still do not know the number  $D$  of deceptive keywords contained in the repository. Thus, Enhanced Gray Box adversaries cannot be sure if the  $j^{\text{th}}$  keyword, with  $j \in \{1, \dots, M\}$ , provided by the Oracle Function, is a deceptive keyword or not.

### B. Countering the Operations

The Gray Box adversaries are aware of the deceptive operations described in this work. In this section, we explore a possible approach that this kind of adversary could carry on to counter the deceptive operations and build a new repository  $\mathcal{R}'$  that is more significant than  $\mathcal{R}'$ .

The adversaries, to smooth the effect of the deceptive operations, have to solve the following two problems: (1) estimate the right number of clusters contained in the repository, and (2) estimate how many and which deceptive keywords are in the repository  $\mathcal{R}'$ . Recall that if the adversaries evaluate the number of clusters in  $\mathcal{R}'$  leveraging standard techniques such as the Silhouette score or other internal validation measures as explained in Sec. II-B.2, they find out the deceptive clusters accordingly with Sec. III-E. Therefore, adversaries have to counter the deceptive operations to obtain meaningful information from the exfiltrated repository.

The adversaries know that, whatever the adopted policy to replace keywords with deceptive keywords, the subsets of documents that form the clusters in  $\mathcal{R}'$  are held together, or separated, among them by the keywords with higher TF-IDF weight. Thus, a possible approach to reduce the effect of the deceptive operations and restore the original clustering of documents is to remove those keywords that are likely deceptive keywords from the repository  $\mathcal{R}'$ .

The Gray Box adversaries do not know how many keywords they have to remove from  $\mathcal{R}'$ . To estimate the number of keywords to remove and the real number of clusters in  $\mathcal{R}'$ , they may perform the following iterative approach: using an internal validation index (e.g., Silhouette Coefficient), they infer the optimal number of clusters  $K_{init}$  for the repository  $\mathcal{R}'$ . For each cluster  $C'_i$  in  $\mathcal{R}'$ , with  $i \in \{1, \dots, K_{init}\}$ , they rank the keywords by TF-IDF weight and build the ordered list of keywords  $LK_i$  for  $C'_i$ . Let  $T$  be the maximum number of keywords the adversaries are willing to remove from each document. The choice of  $T$  is a trade-off for the adversaries: the more keywords the adversaries delete from the documents, the higher the probability of discarding both deceptive and original keywords. The adversaries perform  $T$  times the following procedure. For each cluster  $C'_i$  in  $\mathcal{R}'$ , they select the keyword  $k$  from  $LK_i$  with the highest TF-IDF weight. Then, the adversaries delete all occurrences of  $k$  from the documents in  $\mathcal{R}'$ , and remove  $k$  from the list  $LK_i$ . Let  $K_{estim}$  be the maximum number of clusters the adversaries suppose to be in  $\mathcal{R}$ . At the end of each step, the adversaries assess on the repository the internal validation score for  $K$  different number of clusters, with  $K \in \{2, \dots, K_{estim}\}$ . At the end of the  $T * K_{deceptive}$  steps, the adversaries evaluate all the internal validation scores they computed and select the configuration that achieved the best score accordingly with the internal validation index they used. If  $\mathcal{R}'$  is the repository that achieves the best internal validation score, the adversaries

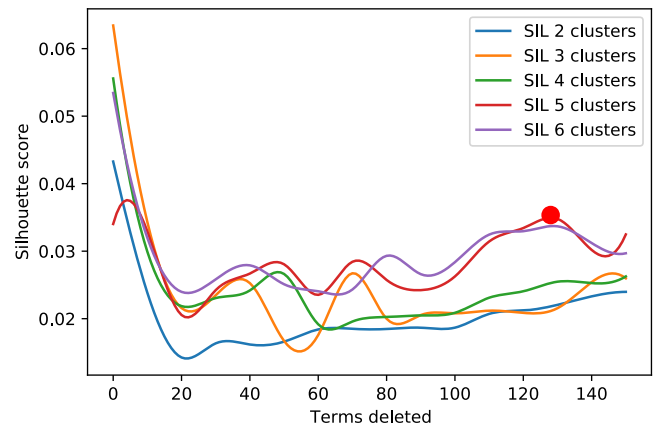


Fig. 5. Values of the Silhouette Coefficient obtained on the modified repository removing up to 300 terms from each identified cluster. The red dot indicates the number of terms to remove in order to achieve the best Silhouette score.

assume as  $\mathcal{R}'$  the restored repository of  $\mathcal{R}$ . Therefore, the adversaries consider the keywords deleted from  $\mathcal{R}'$  to achieve  $\mathcal{R}'$  as the deceptive keywords of  $\mathcal{R}'$  and the number of clusters of  $\mathcal{R}'$  as the number of topics covered by  $\mathcal{R}$ .

Enhanced Gray Box adversaries follow the same approach as Gray Box, but they have a significant advantage in determining the deceptive keywords since they can rely on the Oracle Function. Indeed, Enhanced Gray Box adversaries leverage the Oracle Function to build lists of potential deceptive keywords to remove.

For instance, assume that Enhanced Gray Box adversaries exfiltrate the repository  $\mathcal{R}'_d$ , that has been generated starting from  $\mathcal{R}_d$  using the Basic Shuffle operation and 60 deceptive terms for each document (see Fig 3(a)). The adversaries apply the procedure outlined in this section to obtain a repository  $\mathcal{R}'$ . Enhanced Gray Box adversaries set  $T$  and  $K_{max}$  respectively to 300 terms and 6 clusters. Fig. 5 shows the Silhouette scores obtained by Enhanced Gray Box adversaries removing the terms from the repository  $\mathcal{R}'_d$ . Each colored line represents the Silhouette score for a different number of clusters. At the end of the assessment, the adversaries find out that the Silhouette score is maximized when removing 130 terms from  $\mathcal{R}'_d$  and searching for 5 clusters (red dot in the figure). Thus, the adversaries build the repository  $\mathcal{R}'$  accordingly with the configuration found.

## V. RESULTS

### A. The Deceptive Repositories

To evaluate the adversaries' performances, we built 6 different deceptive repositories starting from the repository  $\mathcal{R}_d$  described in Sec. II-A. To build the 6 repositories, the following deceptive operation has been applied on  $\mathcal{R}_d$ : Basic Shuffle, Shuffle Increment, and Shuffle Reduction. Each deceptive operation has been executed twice, once partitioning the documents through random selection and once through the constrained version of K-means. With both approaches, all the subsets have been made approximately of the same number of documents. The keywords have been replaced in descending order by their TF-IDF weight. Finally, for the sake of comparison, all the operations performed the same number of m-terms-replacements. We set the number of m-terms-replacements to 60, since, according to our experiments,

TABLE IV

ADJUSTED RAND INDEX (ARI) VALUES ACHIEVED BY THE ADVERSARIES ON THE REPOSITORY  $\mathcal{R}'_d$  AGAINST BASIC SHUFFLE, SHUFFLE INCREMENTATION (SHUFFLE INCR.), AND SHUFFLE REDUCTION (SHUFFLE RED.) APPLIED TO SUBSETS OF DOCUMENTS SELECTED RANDOMLY (RANDOM) AND BY GROUPING SIMILAR ONES (SIMILAR). NOTE THAT THE ARI SCORE OBTAINED BY THE ADVERSARIES ON THE ORIGINAL REPOSITORY  $\mathcal{R}$  IS 0.94

	m-replacement	Black Box		Gray Box			Enhanced Gray Box			Original Repository	
		#clust	ARI	#del	#clust	ARI	#del	#clust	ARI	#clust	ARI
<i>Basic - Shuffle<sub>random</sub></i>	60	3	-0.004	50	5	0.33	130	5	0.56	3	0.94
<i>Shuffle - increment<sub>random</sub></i>	60	4	-0.005	60	6	0.28	130	5	0.56	3	0.94
<i>Shuffle - reduction<sub>random</sub></i>	60	2	-0.003	50	5	0.23	130	5	0.56	3	0.94
<i>Basic - Shuffle<sub>similar</sub></i>	60	3	-0.004	60	6	0.27	110	6	0.51	3	0.94
<i>Shuffle - increment<sub>similar</sub></i>	60	4	-0.005	80	6	0.21	110	6	0.51	3	0.94
<i>Shuffle - reduction<sub>similar</sub></i>	60	2	-0.003	40	6	0.20	110	6	0.51	3	0.94

it is the minimum number of replacements such that the Silhouette Coefficient, the Calinski-Harabasz Index, and the Davies-Bouldin Index return the deceptive number of clusters for the repositories. For example, Fig. 3(a), 3(b), 3(c), respectively, show the deceptive repositories built with the Basic Shuffle, the Shuffle Reduction, and Shuffle Increment operations, using the constrained version of K-means to build the subset of documents.

### B. Attacking the Deceptive Repositories

To evaluate the adversaries' performance, we use the Adjusted Rand Index (ARI) [32]. Given a predicted clustering (the one obtained by the adversaries, in our case) and the clustering given by the true labels of the documents, the ARI measures the similarity between these two clusterings. The value of the ARI varies between  $-1$  and  $1$ , where a value of  $1$  indicates a perfect match between the two clusterings, a value close to  $0$  a random labeling of the predicted clustering, and a negative value a labeling worst than a random one.

The Black Box adversaries believe they have exfiltrated the unmodified repository. Therefore, Black Box adversaries analyze the exfiltrated repository as described in Sec. II. At the end of the analysis, Black Box adversaries will discover just the deceptive clusters. Since our operations build each deceptive cluster by grouping together with a uniform distribution of documents of different topics, such clustering of the Black Box adversaries achieves an ARI approximately equal to  $0$ , which corresponds to the same result that the adversary would have by grouping the documents randomly.

Enhanced Gray Box and Gray Box adversaries are aware that the exfiltrated repositories could be deceitful. However, they can not be sure of that. They have two options. They can consider the repositories not deceitful and analyze the repositories as the Black Box adversaries obtaining the same results. Alternatively, they can try to get rid of the deceptive keywords, for example, applying the algorithm described in Sec. IV-B, and building for each exfiltrated repository a recovered version  $\mathcal{R}'$ . We assume that the adversaries analyze the deceptive repositories by searching for 2 up to 6 possible clusters, attempting to remove up to 150 keywords, and evaluating the repositories using the Silhouette Coefficient, the Calinski-Harabasz Index, and the Davies-Bouldin Index. At the end of the analysis, the adversaries find out that the best configurations for each repository  $\mathcal{R}'$  are the ones reported in Table IV. The table reports only the results with higher ARI obtained by the Enhanced Gray Box and Black Box adversaries. All the reported results have been achieved

by leveraging the Silhouette Coefficient. Indeed, it was the internal evaluation index that allowed attackers to reach the best ARI scores in all the experiments. Enhanced Gray Box adversaries are those that achieve the highest ARI, between  $0.51$  and  $0.56$ , while the Gray Box obtains an ARI between  $0.20$  and  $0.33$ . The highest performances of the Enhanced Gray Box adversaries are due to the Oracle Function. Although the Enhanced Gray Box adversaries remove several original keywords from the documents, on average 60 original keywords from each document (over a total of 2,433 terms), they can completely clean up the document from the deceptive keywords. Conversely, the Gray Box adversaries remove a few original keywords from the documents (about 10 original keywords). However, in the Gray Box scenario, each document still contains more than 10 deceptive keywords on average that are sufficient to keep the document in the deceptive clusters created with the deceptive operations described in this work. It is worth noting that the repositories built using the constrained version of K-means (*similar* in Tab. IV and Tab. VII) are those that lead both White-box and Gray-box adversaries to obtain the worst results (*i.e.*, incorrectly clustering the documents). Hence, the *similar* approach appears to be the most robust against the counter operation. Enhanced Gray Box adversaries are able to achieve notable results. Indeed, an ARI of 56% can intuitively be interpreted as the 56% of documents are correctly clustered. However, some considerations have to be taken into account. Even assuming the Enhanced Gray Box adversaries can self-estimate the ARI achieved by the repositories, they still can not infer which documents are correctly labeled and which are not (*i.e.*, the adversaries do not know the original topics of the documents). An ARI of 56% means that if the adversaries pick one document in a cluster  $C'$ , there is roughly 50% probability that  $C'$  contains the majority of documents with the same original topic.

### C. Increasing the Number of True Topics

In this subsection, we explore how deceptive operations perform by scaling up the number of true topics in the original repository. For this analysis, we built 6 different deceptive repositories, each of them containing a different number of true topics, from 2 up to 7. All the repositories have been built leveraging the Basic Shuffle operation (see Sec. III-B) and partitioning the documents randomly. In particular, the 7 true topics we used to build the repositories are the three described in Sec. II ( Artificial Intelligence (AI), Database (DB), Cryptography and Security (CR)), and four new topics (Robotics (RO), Computers and Society (CS),

TABLE V  
SCORES OF ARI ACHIEVED BY THE ADVERSARIES IN THE ORIGINAL REPOSITORY AND AFTER THE BASIC SHUFFLE OPERATION WHEN INCREASING THE NUMBER OF TOPICS

	Number of topics involved					
	2	3	4	5	6	7
Original	0.96	0.94	0.92	0.92	0.89	0.85
Black Box	-0.0002	-0.004	-0.005	-0.005	0.006	0.0003
Gray Box	0.51	0.33	0.29	0.31	0.14	0.12
Enh. Gr. Box	0.68	0.56	0.48	0.47	0.43	0.28

Logic in Computer Science (LO), Computational Complexity (CC)) that we collected in the same fashion always from ArXiv. To assess the performance of the Shuffle Operation, we compute for each repository the ARI before applying the deceptive operation (*i.e.*, the ARI the adversaries would have gotten exfiltrating the plain repository) and the ARI obtained by the three kinds of adversaries after the operation has been applied.

Table V sums up the results of these experiments. As for the previous experiments, the Black Box adversaries achieve an ARI of about 0 for all the configurations. This result is straightforward since the Black Box adversaries attempt to cluster the deceptive repositories without applying any countermeasures. Instead, in the other cases (Original repository, Enhanced Gray Box, and Gray Box adversaries), the ARI drops as we increase the number of topics into the repository. However, while analyzing the original repositories, the ARI falls of few points ( $-0.11$ ), from 0.96 on the repository containing 2 topics to 0.85 on the repository containing 7 topics, for the Gray Box and the Enhanced Gray Box the ARI fall down drastically. The Gray Box adversaries achieve an ARI of 0.51 in the case of a deceptive repository made of 2 topics, while an ARI of 0.12 for the deceptive repository made of 7 topics. The Enhanced Gray Box adversaries go from 0.68 to 0.28. These results show that the deceptive operations become much more effective as the number of clusters in the deceptive repository increases.

#### D. CORD-19 Dataset

This section assesses the effectiveness of our deceptive operations when attackers exfiltrate a repository made of a large amount of documents protected with DARD. To this end, we employ the CORD-19 dataset [33]. This dataset consists of over 140,000 scientific articles on viruses related to the coronavirus family. To create a representative sample, we select 10,000 documents from the CORD-19 dataset using the same methodology outlined by Eren et al. [34]. Subsequently, we preprocess the sampled documents using the method described in Section II-B.1.

To determine the optimal number of clusters of the repository based on the CORD-19 dataset, we use the approach detailed in Section II-B.2. In particular, we evaluate the Silhouette Coefficient on the repository for a number of clusters ranging from 10 up to 50. At the end of the process, we find that the Silhouette Coefficient achieves its maximum value with 20 clusters, which aligns with the results obtained in [34]. Since the CORD-19 dataset is not labeled, we consider the outcome of the clustering process as the ground truth for the subsequent experiments. Similar to the approach described in Section V-A, we build six new deceptive repositories of 10,000 documents each, as described below. We apply each

deceptive operation twice: once by partitioning the documents through random selection and once by using the constrained version of K-means. Starting from the original repository using the Basic Shuffle, we generate a new deceptive repository of 20 clusters. Starting from the original repository made of 20 clusters, we generate a new deceptive repository made of 25 clusters employing the Shuffle Increment. Finally, with the Shuffle Reduction, we reduce the number of clusters in the original repository, producing a deceptive repository consisting of 15 clusters. We set the number of m-terms-replacement to 70, such that the three internal-validation indexes return the deceptive numbers of clusters for all the repositories. Also, in this scenario, we simulated attackers willing to remove up to 150 keywords and search for 2 up to 50 clusters by evaluating the repositories through the Silhouette Coefficient. Table VI reports the ARI values achieved by the different typologies of attackers on the six different deceptive repositories. In this scenario, our operation also completely deceives the Black Box adversaries. Similar to previous experiments, the Gray Box and the Enhanced Gray Box tend to overestimate the number of clusters in the original repository. Additionally, a significant decrease is observed when comparing the ARI scores for both Gray Box and Enhanced Gray Box adversaries with the experiment in Sec. V-A. Specifically, the scores decrease from the range 33-20% to 25-16% for the Gray Box and from 56-51% to 30%, for the Enhanced Gray Box adversaries. This appears to be evidence that the quality of clustering achieved by these adversaries deteriorates as the size of the repository increases.

#### E. Investigating the Runtime Execution of DARD

For assessing the execution time of DARD defense, we use a desktop computer with an Intel Core i7-12700 CPU operating at 2.10GHz and 32 GB of RAM running ManjaroLinux 22.1.0 as the operating system. The deceptive operation has been coded in Python 3.10.10. The scalability of the proposed deceptive operations, as the size of the dataset grows, depends on the computation time of the Equations 2. More specifically, this equation necessitates performing numerous TF-IDF and document clustering operations. For instance, on the Arxiv papers dataset comprising 450 documents, it took approximately 13 seconds to compute the TF-IDF once and 0.83 seconds for a single clustering operation. However, the computation time increased when applied to a larger dataset, such as CORD-19, consisting of 10,000 documents. Calculating the TF-IDF and clustering one time on CORD-19 took 124 seconds and 29 seconds, respectively. In Equation 2, we employ a binary search strategy to determine the minimum number of keywords to be replaced. This binary search is executed within the list of keywords that holds the maximum number of keywords the defender intends to substitute, which, in this experiment, is set at 150. Implementing the binary search strategy results in almost a tenfold reduction in the frequency of TF-IDF and clustering computations. In summary, to calculate Equation 2, DARD requires approximately 81 seconds on the Arxiv papers dataset (450 documents) and 7,122 seconds on the larger CORD-19 dataset (10,000 documents). Results can be enhanced by implementing DARD in a more efficient programming language. From the perspective of adversaries, it is worth noting that when analyzing a deceptive repository, as outlined in the attacking strategy detailed in Section V-B, we estimated that attackers need

TABLE VI

ADJUSTED RAND INDEX (ARI) VALUES ACHIEVED BY THE ADVERSARIES ON THE REPOSITORY  $\mathcal{R}'_d$  AGAINST BASIC SHUFFLE, SHUFFLE INCREMENTATION (SHUFFLE INCR.), AND SHUFFLE REDUCTION (SHUFFLE RED.) APPLIED TO SUBSETS OF DOCUMENTS OF THE CORD-19 DATASET SELECTED RANDOMLY (RANDOM) AND BY GROUPING SIMILAR ONES (SIMILAR). THE NUMBER OF CLUSTERS IN THE ORIGINAL REPOSITORY IS 20

	m-replacement	Black Box		Gray Box		Enhanced Gray Box			
		#clust	ARI	#del	#clust	ARI	#del	#clust	ARI
<i>Basic - Shuffle<sub>random</sub></i>	70	20	-0.001	50	50	0.25	100	50	0.30
<i>Shuffle - increment<sub>random</sub></i>	70	25	-0.004	60	50	0.27	100	50	0.30
<i>Shuffle - reduction<sub>random</sub></i>	70	15	-0.001	40	48	0.17	100	50	0.30
<i>Basic - Shuffle<sub>similar</sub></i>	70	20	-0.004	40	49	0.17	100	50	0.30
<i>Shuffle - increment<sub>similar</sub></i>	70	25	-0.005	40	49	0.18	100	50	0.30
<i>Shuffle - reduction<sub>similar</sub></i>	70	15	-0.003	40	49	0.16	100	50	0.30

TABLE VII

PERCENTAGES OF DECEPTIVE KEYWORDS RETURNED BY LDA ACCORDING TO THE NUMBER OF KEYWORDS RETRIEVED FROM EACH TOPIC AND THE DECEPTIVE OPERATION APPLIED ON THE REPOSITORY

	Number of topic keywords				
	10	20	30	40	50
<i>Basic - Shuffle<sub>random</sub></i>	100%	98%	98%	98%	95%
<i>Shuffle - increment<sub>random</sub></i>	100%	93%	92%	86%	80%
<i>Shuffle - reduction<sub>random</sub></i>	100%	95%	88%	85%	75%
<i>Basic - Shuffle<sub>similar</sub></i>	100%	98%	96%	96%	89%
<i>Shuffle - increment<sub>similar</sub></i>	100%	90%	86%	83%	80%
<i>Shuffle - reduction<sub>similar</sub></i>	100%	100%	95%	85%	74%

much longer execution time. Approximately 360 seconds and 36,000 seconds on the Arxiv papers and CORD-19 datasets, respectively. This high cost is because the attackers must compute the TF-IDF and the clustering operation for each deceptive keyword they suspect. In our experiment, we assume that the adversary aims to remove 150 keywords, which is consistent with the number employed in the aforementioned experiment conducted by DARD.

### F. Topic Modeling

The previous subsections analyzed the effect of deceptive operations on document clustering. In particular, this subsection shows how deceptive operations affect topic modeling. We consider adversaries that try to infer the underlying topics of the exfiltrated repository leveraging LDA, one of the most popular topic modeling algorithms. Given a repository of documents and a number of topics  $T$ , LDA computes for each term in the documents the probability that the term belongs to one of the  $T$  topics. The set of  $n$  terms that have the highest probability according to the LDA algorithm is defined as the *topic keywords* of each topic. Note that the most used number of topic keywords is 10 [19].

Assume that a Black Box adversary wants to infer the topic of each cluster of documents into the exfiltrated repository  $R_d$ . Thus, the number of topics  $T$  that such adversaries seek is equal to the number of deceptive clusters.

To assess if our deceptive operations are able to deceive the Black Box adversaries in inferring the topics, we use the same deceptive repositories described in Sec. V-A, and the LDA version of Scikit-learn library [35]. In particular, we measure the presence of deceptive keywords that the adversaries retrieve with LDA for each topic with  $n$  varying from 10 to 50.

Analyzing the topic keywords computed by LDA on the repositories, we make the following observations:

(i) With  $n = 10$ , for each topic, all the keywords retrieved by LDA are deceptive keywords. Instead, with  $n = 50$ , the percentage of deceptive keywords varies between 74% and 95%, depending on the deceptive operations used to build the repository. Thus, even considering a significant amount of topic keywords ( $n = 50$ ), very few real topic keywords are retrieved using LDA. Moreover, as we can see from Tab. VII, the real topic keywords are mainly at the lower position of the rank, which means they have a low probability of being significant for the specific topic. Indeed, looking at the real topic keywords returned by LDA, we find that they are generic keywords of computer science, such as: *lemma*, *root*, *induct*, *resource*, *program*, *rate*, *label*.

(ii) Each topic retrieved by LDA describes with its top ten keywords a different deceptive cluster within the deceptive repository. Consequently, adversaries relying on LDA will infer the deceptive topics designed by the defenders.

(iii) Each deceptive cluster of the deceptive repository has one topic associated with it.

Combining these observations, we have that LDA finds as topic keywords for a certain deceptive cluster, the same deceptive keywords used to build the deceptive cluster itself. Thus, the defender has the ability to manipulate the topic that the adversaries will infer, choosing wisely the deceptive keywords to use with the deceptive operations. The defender has multiple strategies to pick the deceptive keywords in order to show a specific deceptive topic to the adversaries. Suppose the defender selects deceptive keywords that fit both the deceptive context of the sentence and the part of speech of the terms to be replaced. In that case, it will be more challenging also for a domain expert to recognize the documents modified by our operations. To automatically perform this task, the defender can leverage language modeling techniques such as [36] and [37]. This paper does not discuss the possible strategies of this extension, as it affects the second phase of the attack, while we focus on the first (see Sec. I).

Note also the following interesting consequence of the observations depicted above. Adversaries that first leverage LDA to feed a clustering algorithm at the end of the computation will group the documents of  $R_d$  in a way very close to the deceptive clusters constructed by the defender. Indeed, on our 6 deceptive repositories, the Black Box adversaries group on average the 97% of the documents in  $R_d$  accordingly to the deceptive clusters built by the deceptive operations.

As a further experiment, to assess the robustness of the deceptive operations against LDA, we also leverage a

commercial topic modeling tool, the **Amazon Comprehend-Topic Modeling** [13]. It is a topic modeling tool developed by Amazon Inc. that leverages the Amazon SageMaker Latent Dirichlet Allocation (LDA) algorithm, a custom version of LDA developed by Amazon Inc.. Given a repository of documents, Amazon’s tool provides as output two CSV files. The first file reports for each document of the repository the topic number that the document is assigned to, and the proportion of the document concerned with that specific topic. Instead, the second file contains the top 10 topic keywords for each topic.

As for the results obtained with SKlearn’s implementation of LDA, at the end of the experiments on the deceptive repositories  $R_d$ , all the topic keywords returned by Amazon’s tool are deceptive keywords. All the topic keywords retrieved belong to a single deceptive cluster, and each deceptive cluster has a topic associated with it. The result of the experimentation with Amazon’s tool shows that our proposed deceptive operations are robust also if adversaries employ an implementation of LDA not known to the defender, and with optimized parameter and data processing pipeline.

Note that leveraging the Amazon Comprehend tool makes sense only for the Black Box adversaries. Indeed, the Enhanced Gray Box and Gray Box adversaries are aware of the presence of deceptive keywords in the exfiltrated repository. Thus, they do not perform topic modeling but attempt to get rid of the deceptive keywords and cluster the documents achieving the poor results shown in the previous sections.

## VI. RESTORING THE ORIGINAL DOCUMENTS

This section describes a possible architectural design involved in restoring a deceptive document to provide the original document to a legitimate user. This ensures that the DARD deception remains completely transparent from the perspective of legitimate users. The keywords mapping, the most sensitive data of the DARD system, cannot be stored on the main file system, risking being exfiltrated together with the other repository documents. To address this concern, a secure application must be developed for file restoration, utilizing a Secure Enclave Solution [12]. By adopting this approach, the mapping can be securely stored within the Secure Memory of the Secure CPU. Performing the restore operation on the Secure CPU ensures that the keywords mapping remains isolated and safeguarded. In Fig. 6, we report a possible architecture to explain how the restore process can work. In particular, there is a deceptive repository stored in a storage server, and  $N$  authorized users can access it through their workstations. The mapping of the deceptive keywords with the relative original keywords is kept in the secure enclave of the workstation. Moreover, the users’ document reader/editor application has a plug-in installed that can interact with the secure enclave. When a user requests access to a document, if he has the privilege to read/write it, the storage server sends the deceptive document to the user. At this point, the plug-in interacts with the secure enclave to authenticate the user. If the secure enclave recognizes the user as legitimate, the modified document is restored and visualized to the end user. The original document is always maintained in the workstation’s RAM. When a user saves an update to a document, the plug-in interacts with the secure enclave to replace the existing keywords with the corresponding deceptive ones. After receiving the updated deceptive document from the secure enclave, the plug-in sends it to the storage server, where the revised version

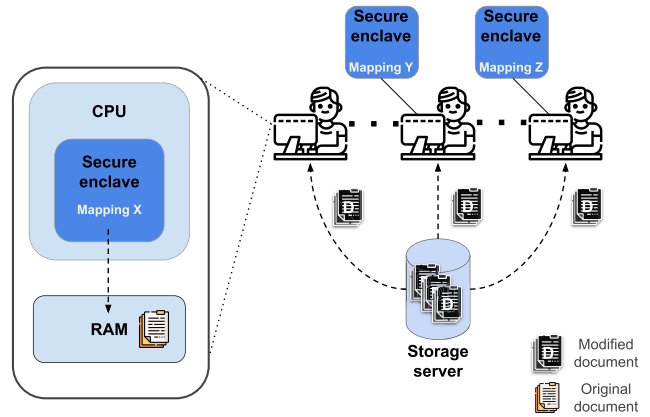


Fig. 6. A possible architecture design to restore deceptive documents.

of the document is then stored. Note that the original document is never stored on the workstation’s hard disk nor transmitted over the network.

## VII. RELATED WORK

### A. Adversarial Setting

Several studies propose attacks to clustering algorithms through the generation of adversarial settings [38], [39]. Generally, an adversary injects malicious examples into the data to impact the clustering results. There are two main typologies of these attacks: poisoning and obfuscation. The poisoning attack aims to worsen the clustering results as much as possible by corrupting the data. The strategy is to create new clusters or bridges between clusters by adding samples within the dataset. In the first case, the purpose is causing the misclassification of a single cluster into more clusters. As for the bridge case, the goal is to pretend that two different original clusters are one. Instead, an obfuscation attack aims to hide a specific data set. Typically, it consists in adding a set of samples to join the target cluster to hide with another one. As a result, clustering methods return a unique cluster that conceals the target cluster with another one. As in the work of [38] and [39], DARD can generate or reduce the number of clusters within a repository. However, unlike previous approaches, DARD does not require the introduction of new samples into the repository to execute deceptive operations. Obfuscation allows one to completely hide a cluster by merging it into another cluster. Nevertheless, there’s no assurance that the newly formed cluster will successfully conceal all the keywords from the hidden cluster. Thus, while obfuscation effectively misleads regarding the number of clusters, it doesn’t guarantee the concealment of topics or deception about the underlying topics within a repository. Similarly, poisoning worsens the performance of clustering algorithms, increasing the number of clusters within a repository. However, in this case, the topics remain retrievable using machine learning techniques.

Unlike previous studies that focused on fooling image classifiers [38] and malware classifier [39], our research addresses explicitly text document classifiers. Specifically, DARD introduces a novel approach by utilizing adversarial settings not for attacking models but as a defense mechanism against adversaries relying on automated techniques to classify exfiltrated repositories. This paradigm shift also alters the threat model. Previous studies assumed adversaries would carry out stealthy

attacks and target a specific system, typically knowing its parameters. However, in our work as defenders, we cannot assume in advance the number of clusters the attackers will seek or whether they will attempt to circumvent the DARD system.

### B. Adversarial Text

These techniques aim to reduce machine learning ability to automatically analyze text without providing a deceptive layer. For example, obfuscation techniques can hide the real author of a document from automatic Authorship Attribution [40], [41] by changing the writing style but keeping the content of the document understandable to the reader. However, these techniques are not able to cover the underlying topic of a document nor obfuscate the number of topics present in a repository. Indeed, these techniques are usually limited to inserting minor typos or replacing a word with a synonymous, attempting to not alter the meaning of the sentences contained in the text.

### C. Deceptive Repositories

Finally, there are solutions similar to ours that aim to create deceptive repositories. Chakraborty et al. propose Forge [42]. This system leverages ontologies to generate new fake documents, credible to unauthorized readers, from a set of original documents. The resulting deceptive repository will contain fake and original documents that are indistinguishable to human readers. Identifying the fake documents within the repository requires extensive reading and, thus, a significant investment of time to differentiate them from the original ones. WE-Forge [43], an extension of Forge, goes beyond ontologies and utilizes word embeddings to automatically generate fake documents that closely resemble authentic ones, enhancing the credibility of the forged documents. While Forge and WE-Forge focus on building a deceptive repository to mislead the human reader by creating fake new documents, we use the same documents of the original repository to generate a new deceptive repository to mislead automatic classification systems. Moreover, Forge and WE-Forge mainly aim to hide details of the original documents and not the underlying topics. The DARD solution and Forge/WE-Forge solution are orthogonal and can be jointly used to hinder both human analysis and automatic systems. For example, Forge could be applied to create a deceptive repository to resist a human analyst, followed by DARD to make the repository deceptive against automatic analysis.

## VIII. DISCUSSION ON DARD'S APPROACH

This section compares DARD with traditional encryption and current obfuscation practices, considering specific aspects such as focus on deception and deceptive operations.

### A. Focus on Deception

While encryption primarily focuses on data confidentiality by securing information through cryptographic algorithms, DARD emphasizes topic confidentiality and deceiving attackers over data secrecy. Instead of merely encrypting data, DARD creates a deceptive repository intentionally crafted to attract attackers, giving them the impression that they have discovered valuable and easily accessible information.

### B. Deceptive Operations

DARD's deception tactics go beyond encryption and obfuscating data. It involves carefully crafting the repository's structure, including the number of clusters and their topics, to make attackers believe they have discovered valuable data. This entices attackers to conduct further analysis, wasting their time and resources on false content. This proactive approach not only facilitates intrusion detection but also actively engages attackers by leading them away from genuine data.

### C. Ease of Deployment and Usability

In terms of ease of deployment and end-user usability, DARD is comparable to traditional encryption methods. Both approaches require secure storage for keeping secrets (e.g., encryption keys or keyword mappings) and keeping plaintext confined to RAM during data processing. Therefore, organizations can integrate DARD into their existing security infrastructure without significant additional overhead or complexity.

### D. Response Efficiency

Unlike encrypted data, where attackers immediately recognize encrypted content upon discovery, DARD's deceptive repositories delay attackers' realization that they have encountered engineered content. This delay provides security teams valuable time to detect and respond to threats before attackers move on to other servers or resources in search of unencrypted data.

In addition, there are cases where the DARD system could be used complementary to encryption to enhance data security and privacy. By combining both approaches strategically, organizations can achieve a more comprehensive and robust security posture. Here are some scenarios where DARD can complement encryption:

### E. Data Decoy Strategy

In addition to encrypting sensitive data, organizations can deploy DARD to create decoy repositories or honeypots containing fabricated or low-value information. Attackers who gain unauthorized access may be drawn to these decoy repositories, allowing security teams to detect and respond to intrusions more effectively. Meanwhile, genuine sensitive data remains encrypted, providing an additional layer of protection.

### F. Multi-Layered Defense

Implementing a multi-layered defense approach involves using multiple security measures to protect data. Encryption serves as one layer of defense to safeguard data in transit and at rest. DARD can be deployed as another layer to detect and disrupt malicious activities by deceiving attackers and diverting their attention from genuine data assets. Together, encryption and DARD create a more resilient defense against sophisticated cyber threats.

In summary, the DARD methodology offers a unique and effective approach to enhancing data security and privacy by prioritizing deception over traditional encryption. While encryption and DARD serve different purposes in data security, they can be effectively integrated to provide complementary layers of protection.

## IX. CONCLUSION

This work proposes DARD, a framework of 4 deceptive operations able to manipulate the resulting clusters generated by a document repository. The goal is to deceive adversaries who employ automatic classification approaches to categorize exfiltrated documents. To this end, the deceptive operations replace some of the original keywords in the documents with deceptive keywords through term-replacement operations. We outlined how to apply the term-replacement operations for each deceptive operation and identified the minimum number of terms that needed replacing. Then, we investigate different criteria for selecting the terms to be replaced, highlighting the pros and cons of the different approaches. We show experimentally that our operations can achieve a high level of deception. We conduct our experiments with three different types of adversaries: the Black Box, an adversary who does not know anything about deceptive operations; the Gray Box, who knows how deceptive operations work; and the Enhanced Gray Box, an adversary that can leverage the Oracle Function to discover the potential deceptive keywords in the repository. Our results show that deceptive operations completely deceive adversaries without knowledge of this work (0% ARI). They are very effective (average ARI of 25%) against those adversaries who know how the deceptive operations work (Gray Box), achieving, in the worst-case scenario with the Enhanced Gray Box adversaries, an average ARI of 53.5%. In addition, we analyzed the impact of deceptive operations in the topic modeling task.

We found that when the adversaries perform topic modeling with LDA on the deceptive repositories, LDA only returns deceptive keywords among those most descriptive of the topic. Moreover, we show that our approach against topic modeling successfully deceives also commercial tools such as Amazon Comprehend.

## REFERENCES

- [1] Deloitte. (2021). *Impact of COVID-19 on Cybersecurity*. [Online]. Available: <https://www2.deloitte.com/ch/en/pages/risk/articles/impact-covid-cybersecurity.html>
- [2] Interpol. (2021). *COVID-19 Cyberthreats*. [Online]. Available: <https://www.interpol.int/Crimes/Cybercrime/COVID-19-cyberthreats>
- [3] J. Tidy. *The Three Russian Cyber-attacks the West Most Fears*. Accessed: Apr. 21, 2023. [Online]. Available: <https://www.bbc.com/news/technology-60841924>.
- [4] M. L. Morgia, A. Mei, E. N. Nemmi, S. Raponi, and J. Stefa, "Nationality and geolocation-based profiling in the Dark(Web)," *IEEE Trans. Services Comput.*, vol. 15, no. 1, pp. 429–441, Jan. 2022.
- [5] M. L. Morgia, A. Mei, A. M. Mongardini, and J. Wu, "It's a trap! Detection and analysis of fake channels on telegram," in *Proc. IEEE Int. Conf. Web Services (ICWS)*, Jul. 2023, pp. 97–104.
- [6] CISA. *Impacket and Exfiltration Tool Used To Steal Sensitive Information From Defense Industrial Base Organization*. Accessed: Apr. 21, 2023. [Online]. Available: <https://www.cisa.gov/uscert/ncas/alerts/aa22-277a>
- [7] S. AG. (2020). *Software Ag Adhoc: Disruption of Services Due To Malware Attack*. [Online]. Available: <https://www.softwareag.com/encorporate/company/news>
- [8] T. Group. (2020). *Toll It Systems Update*. [Online]. Available: <https://www.tollgroup.com/toll-it-systems-updates>
- [9] Verizon. (2021). *Dbir 2021 Data Breach Investigations Report*. [Online]. Available: <https://www.verizon.com/business/resources/reports/dbir/>
- [10] Wikipedia. (2020). *United States Federal Government Data Breach*. [Online]. Available: <https://en.wikipedia.org/wiki/2020UnitedStatesfederalgovernmentdatabreach>
- [11] G. Pagnotta, F. De Gaspari, D. Hitaj, M. Andreolini, M. Colajanni, and L. V. Mancini, "DOLOS: A novel architecture for moving target defense," *IEEE Trans. Inf. Forensics Security*, vol. 18, pp. 5890–5905, 2023.
- [12] Intel. *Intel Software Guard Extensions*. Accessed: Apr. 21, 2023. [Online]. Available: <https://www.intel.com/content/www/us/en/developer/tools/softwareguardextensions/overview.html>
- [13] Amazon. (2021). *Topic Modeling*. [Online]. Available: <https://docs.aws.amazon.com/comprehend/latest/dg/topic-modeling.html>
- [14] M. Steinbach, G. Karypis, and V. Kumar, "A comparison of document clustering techniques," Dept. Comput. Sci. Eng., Univ. Minnesota, Tech. Rep. 00-034, 2000.
- [15] M. Marciniuk et al., "Text document clustering: Wordnet vs. TF-IDF vs. word embeddings," in *Proc. 11th Global Wordnet Conf.*, 2021, pp. 207–214.
- [16] *Open-Access Archive. Arxiv*. Accessed: Apr. 21, 2023. [Online]. Available: <https://arxiv.org/>
- [17] M. S. G. Karypis, V. Kumar, and M. Steinbach, "A comparison of document clustering techniques," in *Proc. TextMining Workshop (KDD)*, 2000, pp. 1–2.
- [18] T. Hofmann, "Unsupervised learning by probabilistic latent semantic analysis," *Mach. Learn.*, vol. 42, nos. 1–2, pp. 177–196, 2001.
- [19] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet allocation," *J. Mach. Learn. Res.*, vol. 3, pp. 993–1022, Jan. 2003.
- [20] A. K. Uysal and S. Gunal, "The impact of preprocessing on text classification," *Inf. Process. Manage.*, vol. 50, no. 1, pp. 104–112, Jan. 2014.
- [21] T. Joachims, "A probabilistic analysis of the Rocchio algorithm with tfidf for text categorization," in *Proc. 14th Int. Conf. Mach. Learn. (ICML)*, vol. 97, Jul. 1997, pp. 143–151.
- [22] R. Tibshirani, G. Walther, and T. Hastie, "Estimating the number of clusters in a data set via the gap statistic," *J. Roy. Stat. Soc. Ser. B, Stat. Methodol.*, vol. 63, no. 2, pp. 411–423, Jul. 2001.
- [23] G. W. Milligan and M. C. Cooper, "An examination of procedures for determining the number of clusters in a data set," *Psychometrika*, vol. 50, no. 2, pp. 159–179, Jun. 1985.
- [24] P. J. Rousseeuw, "Silhouettes: A graphical aid to the interpretation and validation of cluster analysis," *J. Comput. Appl. Math.*, vol. 20, pp. 53–65, Nov. 1987.
- [25] T. Calinski and J. Harabasz, "A dendrite method for cluster analysis," *Commun. Statist. Simul. Comput.*, vol. 3, no. 1, pp. 1–27, 1974.
- [26] D. L. Davies and D. W. Bouldin, "A cluster separation measure," *IEEE Trans. Pattern Anal. Mach. Intell.*, vols. PAMI-1, no. 2, pp. 224–227, Apr. 1979.
- [27] J. A. Hartigan and M. A. Wong, "Algorithm AS 136: A K-means clustering algorithm," *Appl. Statist.*, vol. 28, no. 1, p. 100, 1979.
- [28] L. v. d. Maaten and G. Hinton, "Visualizing data using t-SNE," *J. Mach. Learn. Res.*, vol. 9, pp. 2579–2605, Nov. 2008.
- [29] U. Cambridge, *Online Edition (C) 2009 Cambridge UP An Introduction to Information Retrieval Christopher D.* Cambridge, U.K.: Cambridge Univ. Press, 2009.
- [30] K. Wagstaff, C. Cardie, S. Rogers, and S. Schroedl, "Constrained k-means clustering with background knowledge," in *Proc. Int. Conf. Mach. Learn.*, 2001, pp. 577–584.
- [31] Y. Liu, Z. Li, H. Xiong, X. Gao, and J. Wu, "Understanding of internal clustering validation measures," in *Proc. IEEE Int. Conf. Data Mining*, Dec. 2010, pp. 911–916.
- [32] L. Hubert and P. Arabie, "Comparing partitions," *J. Classification*, vol. 2, no. 1, pp. 193–218, Dec. 1985.
- [33] L. L. Wang et al., "CORD-19: The COVID-19 open research dataset," in *Proc. 1st Workshop NLP COVID-19 ACL*, Jul. 2020. [Online]. Available: <https://www.aclweb.org/anthology/2020.nlpCOVID19-acl.1>
- [34] M. E. Eren, N. Solovyev, E. Raff, C. Nicholas, and B. Johnson, "COVID-19 Kaggle literature organization," in *Proc. ACM Symp. Document Eng.*, Sep. 2020, pp. 1–4.
- [35] F. Pedregosa, "Scikit-learn: Machine learning in Python," *J. Mach. Learn. Res.*, vol. 12, pp. 2825–2830, Aug. 2011.
- [36] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," 2018, *arXiv:1810.04805*.
- [37] Y. Liu et al., "RoBERTa: A robustly optimized BERT pretraining approach," 2019, *arXiv:1907.11692*.
- [38] J. G. Dutrisac and D. B. Skillicorn, "Hiding clusters in adversarial settings," in *Proc. IEEE Int. Conf. Intell. Secur. Informat.*, Jun. 2008, pp. 185–187.
- [39] B. Biggio, I. Pillai, S. R. Bulò, D. Ariu, M. Pelillo, and F. Roli, "Is data clustering in adversarial settings secure?" in *Proc. ACM workshop Artif. Intell. Secur.*, Nov. 2013, pp. 87–98.

- [40] S. Afroz, M. Brennan, and R. Greenstadt, "Detecting hoaxes, frauds, and deception in writing style online," in *Proc. IEEE Symp. Secur. Privacy*, May 2012, pp. 461–475.
- [41] E. Arabnezhad, M. L. Morgia, A. Mei, E. N. Nemmi, and J. Stefa, "A light in the dark web: Linking dark web aliases to real Internet identities," in *Proc. IEEE 40th Int. Conf. Distrib. Comput. Syst. (ICDCS)*, Nov. 2020, pp. 311–321.
- [42] T. Chakraborty, S. Jajodia, J. Katz, A. Picariello, G. Sperli, and V. S. Subrahmanian, "A fake online repository generation engine for cyber deception," *IEEE Trans. Depend. Secure Comput.*, vol. 18, no. 2, pp. 518–533, Mar. 2021.
- [43] A. Abdibayev, D. Chen, H. Chen, D. Poluru, and V. S. Subrahmanian, "Using word embeddings to deter intellectual property theft through automated generation of fake documents," *ACM Trans. Manage. Inf. Syst.*, vol. 12, no. 2, pp. 1–22, Jun. 2021.



**Alberto Maria Mongardini** received the Laurea (summa cum laude) and Ph.D. degrees in computer science from the Sapienza University of Rome, Rome, Italy, in 2020 and 2024, respectively. From April 2023 to October 2023, he was a Visiting Scholar with the Department of Cyber Security, George Mason University, Fairfax, VA, USA, hosted by Prof. Giuseppe Ateniese. He is currently a Post-Doctoral Researcher with the Department of Computer Science, Sapienza University of Rome. He specializes in studying frauds and scams in

online social networks (OSNs) and analyzing market manipulation in blockchains. His research interests include security and privacy.



**Massimo La Morgia** received the Laurea (summa cum laude) and Ph.D. degrees in computer science from the Sapienza University of Rome, Rome, Italy, in 2014 and 2019, under the supervision of Alessandro Mei. From 2014 to 2015, he was a Consultant at different Italian companies. He was a Visiting Scholar with the Center for Secure Information Systems, George Mason University, Fairfax, VA, USA, from October 2023 to November 2023, hosted by Prof. Sushil Jajodia. He is currently an Assistant Professor (RTD-A) with the Computer Science

Department, Sapienza University of Rome. He is involved in several technology transfer activities regarding the IoT systems, machine learning, mobile technology, and proximity payments. His research interests include computer systems, blockchain, security and privacy, blockchain and cryptocurrencies, applied machine learning, computer science, and human behavior. He won the 'Avvio alla Ricerca 2017' and 'Avvio alla Ricerca 2021' awards from the Sapienza University of Rome. In 2017, he obtained the Google Associate Android Developer Certification. In 2023, he was awarded by the European Commission with the Seal of Excellence.



**Sushil Jajodia** (Life Fellow, IEEE) is currently a University Professor, a BDM International Professor, and the Director of the Center for Secure Information Systems, George Mason University. Prior to joining George Mason University, he held permanent positions at NSF, NRL, and the University of Missouri-Columbia. He has sustained a highly active research agenda spanning database and cyber security for over 30 years. He has authored or coauthored seven books, edited 52 books and conference proceedings, and published more than 500 technical articles in refereed journals and conference proceedings. He holds 23 U.S. patents and has received a number of prestigious awards in recognition of his research accomplishments. According to Google Scholar, he has over 54 000 citations, and his H-index is 116.



**Luigi Vincenzo Mancini** received the Ph.D. degree in computer science from Newcastle University, Newcastle upon Tyne, U.K., in 1989. He is currently a Full Professor with the Dipartimento di Informatica, Sapienza University of Rome. His contributions extend beyond research; he has actively participated in program committees for several prominent international conferences. He recognized the importance of education in cybersecurity, he established a series of Master's degree programs in information and network security at the Sapienza University of Rome,

which continues to attract numerous students. His commitment to advancing knowledge in the field is evident through his leadership as the principal investigator in numerous national and international research projects focusing on security and privacy. He has published more than 140 scientific papers in international cybersecurity conferences and journals.



**Alessandro Mei** (Member, IEEE) received the Laurea degree (summa cum laude) in computer science from the University of Pisa, Italy, in 1994. He is currently a Full Professor with the Computer Science Department, Sapienza University of Rome, Rome, Italy. He is a member of the ACM. He was a Marie Curie Fellow from 2010 to 2012. He was a past Associate Editor of IEEE TRANSACTIONS ON COMPUTERS from 2005 to 2009 and the General Chair of IEEE IPDPS 2009, Rome.