

available at www.sciencedirect.comjournal homepage: www.elsevier.com/locate/cosrev

Survey

DNA-inspired information concealing: A survey

Lukas Kencl^a, Martin Loeb^{b,*}

^a R&D Centre for Mobile Applications (RDC), Czech Technical University in Prague, Czech Republic

^b Department of Applied Mathematics and Institute of Theoretical Informatics (ITI), Charles University, Prague, Czech Republic

ARTICLE INFO

Article history:

Received 9 June 2010

Received in revised form

12 July 2010

Accepted 25 July 2010

Keywords:

Information concealing

Internet security

Web statistics

Repeats in DNA

ABSTRACT

Various research efforts would benefit from the ability to exchange and share information (traces with packet payloads, or other detailed system logs) to enable more data-driven research. Protection of the sensitive content is crucial for extensive information sharing. We present results of Kencl and Loeb (2009) [41] and Blamey et al. (in preparation) [4] about a technique of information concealing, based on introduction and maintenance of families of repeats. The structure of repeats in DNA constitutes an important obstacle for its reconstruction by hybridisation. A large proportion of eukaryotic genomes is composed of DNA segments that are repeated either precisely or in variant form more than once. As yet, no function has been associated with many of the repeats. In the paper by Blamey et al. (in preparation) [4], the authors propose that in eukaryotes the cells have DNA as a depository of concealed genetic information and the genome achieves the self-concealing by accumulation and maintenance of repeats. The protected information may be shared and this is useful for the development of intercellular communication and in the development of multicellular organisms. The results presented here are protected by Czech patent number 301 799 and by US Patent Application number 12/670.

© 2010 Elsevier Inc. All rights reserved.

1. Introduction

Contemporary computer systems may be distributed and may consist of many interconnected processing units or a large number of networked computer subsystems. In addition contemporary digital networks may consist of a large number of end- and intermediate-nodes. In all these systems, information, in the form of the sequences over some alphabet of symbols, is circulating or being stored. The entity controlling a subsystem or a node is often unwilling or prohibited to share these information-sequences with other nodes. However, sharing of some reduced local information might be very useful for purposes of security,

stability and various analyses of the system performance, and for data mining. Such analysis might for example allow the identification of frequently appearing segments by performing approximate statistical analysis on segment frequency, allowing the detection of replicating malicious code-worms. It also allows the identification of segment-markers of computer viral infection, by detecting patterns existing in some database of malicious sequences. Such databases are used e.g. in contemporary intrusion detection systems or spam filters. It has been shown that being able to perform pattern matching against only fixed-length prefixes or substrings of longer sequences can provide approximate hints as to the presence of suspicious content [1]. Likewise, established worm detection techniques such as Autograph [2]

* Corresponding author.

E-mail address: loeb@kam.mff.cuni.cz (M. Loeb).

or EarlyBird [3] are based on counting frequency of small blocks of a fixed size.

Sharing of reduced local information among the members of an interconnected computer system or communication network thus helps to discover attacks earlier. Affected parts may be isolated and further attack spread prevented. The benefits of sharing local information may be reaped in case of existence of a computational information processing, which preserves local information (e.g. all segments of certain maximal length) and makes it impossible to reconstruct longer or sensitive parts of the information sequences.

We call such information processing *concealing*. The systems which conceal information and share the concealed information are likely to possess a competitive advantage in the form of robustness, attack resistance and immunity due to the ability to exchange, publish and protect information. Clearly, any information concealing algorithm needs to address two conflicting goals:

- (1) preserving *presence* and, possibly, *frequency rank* of segments of given size (making spam identification and worm detection still possible), while
- (2) making reconstruction of content longer than the predefined limit computationally hard (e.g. disabling interpretation or understanding of the private content).

The main contribution of this paper is formulation of the information concealing problem and presentation of an information concealing algorithm. The algorithm is based on a principle which we learned in DNA reconstruction. As far as we know this is the first use of this principle outside of nature. The main feature of the method is that the concealed sequence ω_F (the output of the algorithm), which contains the local information from the input sequence ω , gives rise to a set S , $|S|$ exponential, of *feasible reconstructions* of ω . Anyone attempting to reconstruct the original sequence must choose it correctly from S . This is very hard since the elements of S are indistinguishable from the information contained in ω_F , even when knowing the concealing algorithm itself. We analyse the algorithm and under the *consistency assumption* present evidence of the hardness of reconstruction of the input sequence. However, we have not been able to study the problem from the information-theoretic or entropy point of view yet and this remains open to further investigation.

We also want to stress that even though our algorithm is inspired by the repeats in DNA, the output of the algorithm, aside of containing repeats, has no other similarities to the DNA structure as far as we know. The article is organized as follows: in Section 2 we survey the reasons for proposing in [4], along with Jenny Blamey, the conjecture that in eukaryotes the cells have DNA as a depositary of concealed genetic information and the genome achieves the self-concealing by accumulation and maintenance of repeats. In Section 3 we give an overview of related problem areas in the field of information and communication technologies and comment on how they differ from the concealing problem. In Section 4 we pose the concealing problem formally and in Sections 5 and 6 we describe the concealing algorithm and its sub-procedures. In Sections 7 and 8 we analyze the properties of the presented algorithm and prove the hardness of reconstruction, and we conclude in Section 9.

2. Repeats in DNA

2.1. Basic structure of DNA

A *eukaryotic gene* is a unit of DNA segments which facilitates a specific functional gene product that may be either RNA molecules or polypeptides. The segments of a gene include the transcription region (TR) and a large variety of regulatory sequences that flank it. The TR encompasses the coding segments (exons), the intervening segments (introns) and regulatory segments. The TR is copied by a process called transcription into the precursor RNA. The transcribed introns are then removed by a process called splicing. The transcribed exons are joined correctly together during the splicing, and they are transformed into the final specific functional product. Introns vary in size, but the same genes in different species often have the same number of introns at analogous positions, although the length and their sequences differ. The function of the introns has not yet been fully understood.

Economy seems to be less important during the evolution of eukaryotes. Apart from the large part of DNA consumed by introns, eukaryotic genes are separated by long stretches of noncoding DNA sequences. The most striking example is various repeat families. Their existence is a basic feature of eukaryotic DNA. Prokaryotes also contain repeats, but only upto 0.5%. This is in contrast to as much as 50% of repeats for eukaryotes. Telomeric and centromeric regions of chromosomes contain tandemly repeated nucleotide sequences called microsatellites. Aside of tandem repetitions other highly repeated sequences interrupt DNA sequences. The members of these interspread repeat families are found between genes, in introns, within satellite DNA, and even within coding and regulatory regions where they cause mutation. Apart of the repeat families, other examples of non-coding DNA include the intergenic spaces segments (IGS) that separate clusters of ribosomal DNA (rDNA), and the pseudogenes. Pseudogenes are genomic segments similar to specific functional genes but unable to yield functional gene products. They are often found closely linked to the corresponding functional gene and contain corresponding flank sequences that interrupt coding regions. Then, there are processed pseudogenes which are dispersed to distant genomic positions and lack intervening sequences.

It is proposed in [4] that in eukaryotes DNA is formed as a depositary of the concealed genetic information. The concealed genetic information may be shared and this is useful for intercellular and intracellular signalling, communication and in the development of multicellular organisms. The concealing is achieved by repeats accumulation.

2.2. Why DNA is shared and self-concealed

An extensive system of multicellular communication based on signals carried by mollecular messangers enables the multicellular cooperation in the eukaryotes. The basic benefits of such a cooperation are apparent and include optimisation by differentiation into distinct cell types, collective defence, resourses and labor management.

DNA is stable but it is not an inert substance. It is subject to numerous chemical exposures (control of gene expression,

methylation) that modify it. DNA may also be rearranged through the action of transposons, the genetic elements that can move in the genome. Cells have a great variety of DNA damage detection systems and repair mechanisms. This design normally inhibits the formation of unwanted random information and facilitates discarding erroneous information, rendering the cell an efficient information processing device. The view of DNA as a complex computing system has emerged recently. In addition to information, DNA stores energy, available on hybridisation of complementary strands or hydrolysis of its phosphodiester backbones [5]. These systems need to ensure the nearly perfect fidelity of DNA replication. The fidelity is especially important in the multicellular organisms in view of the extend of the cell differentiation. The extensive communication brings up a possibility and a demand for the intercellular monitoring and damage detection.

We propose that the eukaryotes accomplish this by a system where each cell has its DNA as a storage of its self-concealed genetic information. The information in DNA is concealed for *self-protection* in order to be shared for outside queries. The molecular signals going in and out of the cell and its nucleus carry information about short segments of DNA. The longer or sensitive segments are however protected. We propose that DNA achieves this self-concealing by maintenance of the families of repeats.

The assertion that the repeats are maintained in DNA in a programmed way for self-concealing explains basic puzzling features of repeats: the uniformity along with the polymorphism of the repeated sequences; the freedom of the repeated DNA to adopt quite different primary sequences in closely related species; apparent non-functionality of the precise amount or the precise sequence of the repeats.

The containment of repeats versus DNA sequencing problem is receiving extensive attention of biologists, computer scientists and mathematicians (see [6–8]).

2.3. Repeats versus DNA reconstruction

We explain the basic idea of concealing by repeats in this subsection. Assume we are given a collection \mathcal{K} of segments of DNA. Each segment S from \mathcal{K} is divided into two parts, the initial part $S(I)$ and the terminal part $S(T)$. We thus may write $S = S(I)|S(T)$. This is an artificial assumption imposed only for the clarity of the presentation.

A reconstruction of \mathcal{K} is a sequence of its segments so that the terminal part of each segment agrees with the initial part of the next segment in the sequence. If several of these initial and terminal parts coincide, there may be an exponential number of possible reconstructions.

Let us consider a very simple example. Let \mathcal{K} be the following collection of segments, where the initial and the terminal parts are divided by the vertical line:

$A|B, B|A, A|C, C|A, B|C, C|B$

the following sequences are some of the possible reconstructions:

$ABACBCA, ACABCBA, BACABCB, ABCACBA.$

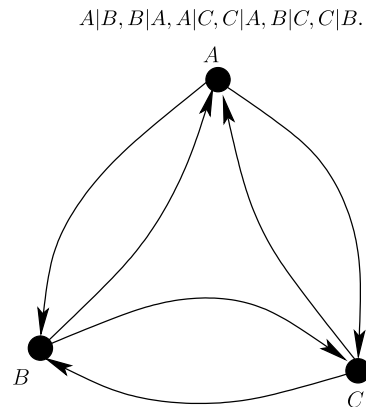


Fig. 1 – De Bruin graph for \mathcal{K} .

In this simple example, even unlimited computational power is useless to anybody who wants to obtain the correct reconstruction from the many possible reconstructions. This phenomenon may well be described in terms of the de Bruin graph (see Fig. 1): this graph has a node for each segment which is an initial or a terminal part of an element of \mathcal{K} . For each segment S of \mathcal{K} there is an arrow (a directed edge) from $S(I)$ to $S(T)$.

The possible reconstructions now correspond to the walks on the de Bruin graph so that each directed edge is traversed exactly once. These walks are usually called Euler walks. If a node of the de Bruin graph has more than one outgoing incident directed edge, then locally there are several independent ways to traverse these edges. The number of the Euler walks of the de Bruin graph is therefore typically exponential in the number of these nodes (see [8] for the calculations).

2.4. Homogeneity and polymorphism of repeat families

An important consequence of the concealing hypothesis is that the repeats in a repeat family should develop in a concerted way; rather than each copy, the whole repeat family has a specific function for the DNA anonymisation. Each intraspecies copy evolves in some sort of communication with the others, a communication that tends to make them all similar. This concerted behaviour has been observed. We include several examples.

Example 1 ([9]). When an activity of an essential gene product becomes insufficient, either by mutation or presence of specific inhibitors, mutant cells that overcome the deficiency can often be found. Many have acquired multiple tandem repeats of the deficient gene, which produces superabundance of mRNA. Repeated units are much longer than the genes themselves; they may contain also sequences from unlinked genomic loci and the individual repeat units are not always identical. Rearrangements within the repeated units can occur continually during extended growth of the cells under conditions favouring amplification. It has been observed that the mutations corresponding to several deficient genes combine. The loss of amplified genes upon removal of selective pressure suggests that

tandem repeats and tandem arrays are dynamic. This fluidity was directly demonstrated in experiments using specially constructed transformed cells. Mouse cells carrying mutations in both thymidine kinase (*tk*–) and adenine phosphoribosyl transferase (*aptr*–) genes were transfected with a recombinant vector containing both a *tk* gene from Herpes simplex virus and an *aptr* gene from hamsters. Cell lines that were both *tk*+ and *aptr*+, containing also tandem amplifications of a DNA segment that includes both genes, were obtained by growth in a suitable selection medium. When selective pressure for amplified *tk* is maintained but the cells are placed in growth media that selects for *aptr*– phenotype, *aptr*– mutants appear at the same frequency as they do in normal cell populations that have a single wild type *aptr* gene. This suggests that all the amplified copies of *aptr*+, 20 in number in these cells, are mutated simultaneously. Thus, the particular mutation varies from one cell line to another, but the overall picture is the same; some mechanism that maintains identity is thus predicted.

In mammals, although different families and subfamilies of interspersed repeats occur, a very small number of families has very high copy number and dominates the genomes. One family up to 5%, together with a small number of families up to 20%. Reciprocal homologous crossing over between nonallelic family members is a possible explanation of homogeneity for tandem arrays but cannot be applied to interspersed repeats without severely jeopardising the integrity of unique sequences that surround the repeats. The homogeneity of these families does not seem to be a result of natural selection acting on each member: if multiple homogeneous copies are functional, their redundancy would tend to minimize the importance of any single mutant copy. However, if most of the copies are nonfunctional, the homogenisation to a species-specific version is even less likely to be influenced by a selective pressure. A natural explanation seems to be that *the repeat family has a concerted (global) function*.

We propose that the natural selection presses for a similarity of the repeats in a family. The similarities in the repetitions and not their actual primal sequences are functional. That is why each single copy in a repeat family is relatively free from selective pressure. This coincides well with another feature of the repeats, the *polymorphism*: the repeats in a repeat family may differ 15%–20%.

Example 2 ([9]). Knobs of heterochromatin are visible at about 23 noncentromeric locations and are associated with a remarkable property: During meiotic divisions, the knobs can form extra functional centromeres. Like maize centromeric heterochromatin, knob heterochromatin is replicated late in S-phase, and the knobs contain long tandem arrays of DNA sequences that are distinct from the sequence of centromeric satellites. Thus, the tandem arrays themselves, not their specific sequence, appear to be associated with the centromere function.

2.5. Mutations and maintenance of repeat families

The necessity of continuous random key generation in some secrecy schemes is consistent with a finding that the repeats mutate more than other genomic regions. Each species has a distinctive set of centromeric satellites, even when compared with closely related species in the same genus. This is a remarkable fact in view of the high degree of sequence preservation elsewhere in genomes. The situation is similar for interspersed repeats: the plasticity is not limited to sequence and copy number. Their locations are not necessarily the same, even in closely related species. Also, otherwise identical alleles in a particular species can differ by the presence or absence of a repeated unit.

It is further proposed in [4] that each cell has a mechanism to maintain the families of repeats. This has been implicitly used in bioengineering. A malfunction of such a mechanism would probably result in the development of several diseases. This may explain a correlation between the increase in the number of repeats present in the genomic DNA and hereditary disorders in humans (see [10–16]). Observations about the changes of repeats in DNA may help to explain other spreading processes in multicellular organisms.

3. Concealing in information and communication technologies

The concept of hiding private or sensitive data but preserving some form of structural information has been studied recently in various sub-domains of ICT. Some techniques concentrate on hiding the originator of information, i.e. *anonymization*, other focus on enabling particular functions over the data that can be shared among multiple partners, such as *private matching*.

3.0.1. Concealing network data

An anonymization scheme over the network packet *IP addresses* called *CryptoPan* [17] preserves the prefix hierarchy of the original addresses, while making them computationally hard to reconstruct by using hashing. This in turn allows the sharing of network traces (with packet headers only), with preservation of the prefix hierarchy.

Similarly, in [18], structure of the router configuration files and data is preserved, while the actual values are obfuscated.

A technique to process and transform the network packet payload has been proposed in [19]. This method uses dictionaries of important sequences that are valuable from the data mining perspective and should be preserved, while encrypting the rest of the information with cryptographically strong hash function. This technique performs well in terms of data protection, however, it only allows study of content portions pre-determined by a known list, and thus does not allow study of the payload to detect previously unknown content, such as e.g. malicious subsequence.

The popular Bloom filter [20] approach is used in constructing the Hierarchical Bloom Filter *payload attribution* technique [21]. A Bloom filter can store (incompletely but efficiently) input items (which can be substrings) and easily

answer set membership queries. It consists of k hash functions, each associating one of m numbers to each input item. Set membership queries exhibit no false negatives, but can have false positives.

Payload attribution with a Hierarchical Bloom Filter stores segments of network packet payloads with their IP source and destination addresses. Each payload is cut into segments s_1, \dots, s_n . The s_i 's are stored in a Bloom filter of level 0, the pairs s_1s_2, s_3s_4, \dots , in Bloom filter of level 1, quadruples in level 2, and so on. A query on an excerpt of payload, which may consist of several consecutive blocks, may answer the source and destination address by running through consecutive hierarchy levels.

The authors propose deployment at network concentration points. Privacy protection is to be achieved by restricting access of entities that can pose queries, otherwise exhaustive attacks might lead to payload reconstruction.

3.0.2. Private matching

Private matching [22,23] focuses on the problem of two entities trying to find common data elements in their databases, without revealing private information. The basic property (and difference from the general information concealing problem) is that only two parties are involved; a multiparty solution is a future work suggestion. Further problems are asymmetry in the sequence of information exchange among the parties and needed presumption of honesty ('semihonesty' in the paper).

Private matching is a special case of cryptography theory of *multi-party computation*: m parties want to compute function f on their m inputs. In the *ideal model*, where a trusted party exists, the parties give their inputs to the trusted authority, it calculates f , and returns the result to each party. The ideal model assumes an ideal situation: for example, no protocol can prevent a party from changing its input before the communication is started. A secure multiparty computation protocol *emulates* what happens in the ideal model.

Paper [22] also introduces 'data ownership certificates' to modify the private matching protocols to be unspoofable. This technique is shown to be useful in a more practical setting to enable privacy-protecting sharing of e-mail white-lists in [24].

3.0.3. Data masking

Various techniques of masking, sanitizing and obfuscating data have been studied to enable test- or third-party development over sensitive databases (such as the Human Resources data). After sanitization, the database remains usable—the look-and-feel and some relations and distributions are preserved—but the information content is secure. The used techniques include masking, shuffling, substitution, number-variance, encryption etc. [25]. These techniques share a similar goal with information concealing, but focus on structured data without the need of preserving the local information.

3.0.4. Data mining and anonymization

In data mining, anonymization mechanisms (obfuscating the originator or the private part of the data) are currently studied intensively. Privacy mechanisms can be classified into several categories, according to where they are deployed during the life cycle of the data. The mechanism proposed in this paper falls into the category where the individuals trust no one but themselves, and they conceal their respective data before they make them available for sharing. The existing algorithms in this category [26–30] are called local perturbation; they are based on different ideas than the concealing procedure proposed in this article.

In another category, data publishing, data are anonymized at a central server; the individuals are required to trust this server [31]. Anonymization in social networks is studied in [32].

An important theoretical foundation for data anonymity and originator protection was laid in [33]. The k -anonymity model for protecting privacy allows holders to release their private data without being distinguishable from at least $k - 1$ other individuals also in the release.

3.0.5. Steganography

This form of information hiding [34,35] is a related art and science of writing hidden messages in such a way that no one apart from the sender and intended recipient realizes there is a hidden message; this nowadays includes concealment of digital information within computer files. Comparably, steganalysis is the art of detecting the hidden information. Steganography is a mature science, in particular focusing on the domain of Digital Rights Management (DRM), where various 'watermarking' or 'tamper-proofing' techniques may seamlessly embed extra information about the origin of a digital work within itself.

3.0.6. Information retrieval

The 'attacker problem' of concealed string reconstruction (see Section 4) has a strong connection to the problem of information retrieval [36], where probabilistic information about the expected string (e.g. natural text) may be used to derive further information or assist text reconstruction.

3.0.7. Computing functions on encrypted data

In recent years, the ability to delegate processing one's data without giving away access to it has gained much interest (motivated e.g. by virtualisation and cloud computing). A family of the so-called *homomorphic* encryption schemes [37] have been proposed to address this problem. Their purpose is related, but not identical, to concealing: homomorphic encryption allows the execution arbitrary operations on (homomorphically) encrypted data, with the result stored in new ciphertext, still decryptable by the originator party only. No information or data is therefore revealed (or shared), yet the outcome of the operation is not shared either—i.e. a third party may execute an operation on the data, yet the result remains encrypted, and therefore obfuscated, to it.

3.1. Segment shuffling

Finally we mention that our first attempt to solve the anonymization problem [38] was using random permutations of a collection of short overlapping segments. This method however by itself does not lead to concealing the original data information. It is shown in this paper that in order to sufficiently extend the families of repeats of the resulting sequence and make the concealing provably successful, other procedures need to be performed as well. In particular the overlapping segments containing complete local information need to be prolonged by attaching additional short segments to their beginning and/or their end. The shuffling permutation also needs to satisfy some properties. This is described in the rest of the paper.

4. Information concealing problem

We introduce formally the information-sequence concealing problem. Let $|\omega|$ denote the number of symbols (length) of sequence ω . The *sequence concealing* problem is the following: Given a sequence ω and a small positive integer k , we want to transform ω to another sequence ω_F so that:

- I. If s is a segment of ω with $|s| \leq k$, then s is a segment of ω_F .
- II. It is computationally hard to reconstruct sequence ω from ω_F .
- III. The length of ω_F is linear in $|\omega|$.
- IV. It is also desirable that with low probability, a segment not in ω appears in ω_F , and that relative frequency (i.e., frequency rank) of segments of ω of a given length is preserved in ω_F . The precise statement of these two conditions is however strongly application dependent.

Given the statement of the information concealing problem, the key issue is how much information about ω can an attacker deduce from ω_F ; let us call this issue the *attacker problem*.

Clearly, the answer to the attacker problem is application-dependent. If the input sequence ω is very restrictive, e.g. if a short prefix uniquely determines larger part of ω and the k -segments of ω may be distinguished within the larger k -segment superset of ω_F , then inevitably large part of input ω may be reconstructed from ω_F . In quite a number of practical situations (DNA sequence, computer program, sound and video trace, text on non-specific topic), however, this is not the case. Moreover, for restrictive input sequences, we can perform preparatory procedures (as procedure S described below) which make the input sequence less specific.

In an attempt towards the solution of the attacker problem, we make the *consistency assumption* below. Assuming the consistency assumption, we can prove in Section 8 that it is hard for the attacker to reconstruct a large part of the initial sequence ω .

Consistency assumption. The complete input of the attacker problem, i.e. all the useful information an attacker has about ω , for instance obtained by analysing ω_F , is

- The length $|\omega|$, the length k of the preserved segments, and the concealing algorithm used in obtaining ω_F .
- The complete list of the frequencies of repeats of segments of ω_F .

5. Concealing by repeats

The input of the problem is a sequence over an alphabet. We first turn it into a cyclic sequence by connecting its beginning and end.

Next we describe five procedures which are used in the algorithm. The basic pattern of all the procedures is the same and may be described as follows: the input is a cyclic sequence ω . First, ω is partitioned into consecutive disjoint *blocks*. Then the terminal part of the preceding block of length o (the *overlap*) is added in front of each block. The resulting segments contain all the studied local information; depending on the procedure, these segments will also contain some excess information which is vital in a proposed composition of the procedures which forms our concealing algorithm. Next, a segment called *dust* can but need not be added behind each segment. The enhanced blocks are called the *cards*. The last step consists in arranging the cards into the output cyclic sequence ω_F .

The first procedure S has a preparatory character in the concealing algorithm. Several runs of S have the role of breaking the local sequential order in the input sequence.

5.1. Procedure $S(\omega, o, lb, ub)$

Its input is a cyclic sequence ω , and it has parameters o, lb, ub ; o stands for the size of the overlap, lb is for lower bound of the length of a block, and ub is for the upper bound of the length of a block. The procedure $S(\omega, o, lb, ub)$ is defined by 1–4 below.

1. We partition (sometimes we say that we cut) ω into consecutive disjoint *blocks* P_1, \dots, P_m such that the length of each P_i is chosen at random between lb, ub .
2. We add overlap of length o in front of each block. The overlapping segments thus contain all the original sub-segments of length up to $o + 1$.
3. The blocks enhanced by the overlaps now start and end with the corresponding overlaps. If these were arranged into a cyclic sequence, the overlaps would neighbor. This may help an attacker in reconstruction. To break the neighborhood relationship of the overlaps, we may add dust (a randomly chosen segment) behind each block. Adding dust is optional and application dependent. A natural restriction is that the dust is a segment of the input sequence and that the average length of dust is $1/2(lb + ub) - o$ to match the average length of the segments complementing the overlaps. However, depending on applications, and the stringency of condition [IV] of the sequence concealing problem, length of dust may be different and the dust need not be a segment of the input sequence.
4. We arrange the resulting cards randomly into a cyclic sequence.

As an illustration we perform S on an example input sequence: 'theEaimEofEthisEpaperEisEtoEpresentEanEinformationEconcealingEalgorithm' where in order to make the presentation easier, we replace the symbol for the empty-space by symbol 'E'. Hence, in the example below, ignore the empty-spaces.

Example 1: Procedure $S(\omega, o, lb, ub)$

Input $\omega = \text{theEaimEofEthiSEpaperEisEtoEpresentEanEinformationEconcealingEalgorithm}$, parameters $o = 3, lb = 4, ub = 6$.

1. First, the input sequence is partitioned randomly into blocks of length 4, 5 or 6. The blocks are divided by '+' below:
 $\text{theEai+mEofEt+hisE+pape+rEisE+toEpr+esent+EanEi+nfor+matio+nEco+ncea+linge+algor+ithm+}$
2. Next we add overlap (of length $o = k - 1 = 3$) in front of each block:
 $\text{thmtheEai+EaimEofEt+fEthiSE+isEpape+aperEisE+isEtoEpr+Epresent+entEanEi+nEinformatio+nEco+Econcea+cealingE+ngEalgor+gorithm+}$
3. Next we add the dust behind each block (of length approximately 2), and we get the cards:
 $\text{thmtheEaip+EaimEofEtim+fEthiSEcon+isEpapeEin+aperEisEa+isEtoEproEp+Epresentese+entEanEilgo+nEinformatiofo+tionEcoE+EconceaEci+cealingEpa+ngEalgorEp+gorithmap+}$
4. Finally the output is given by arranging the cards in a random order (here we use the order 14, 9, 10, 13, 5, 3, 12, 1, 6, 4, 7, 11, 8, 15, 2):
 $\text{ngEalgorEpnEinformiofocealingEpaaperEisEafEthiSEconEconceaEcihmtheEaipisEtoEproEpisEpapeEinEpresentesetionEcoEentEanEilgogorithmapEaimEofEtim}$

5.2. Procedure $S^1(\omega, lb, ub)$

Procedure $S^1(\omega, lb, ub)$ is as S but the overlap is always the whole preceding block—typically exceeding the size needed to preserve the studied local information (this excess is used in the composition of the procedures forming our concealing algorithm). Hence, if the blocks are

$$\omega = P_1 P_2 P_3 \cdots P_m,$$

then the cards of S^1 are $P_1 P_2, P_2 P_3, \dots, P_m P_1$.

Each P_i appears once as initial segment and once as terminal segment of each card. Hence, the cyclic consecutive order of the cards of S^1 may be described by a permutation π of $1, \dots, m$; for further discussions it turns out useful to define such permutation so that it assigns, to each terminal block of a card, the initial block of the next card. By permutation of $1, \dots, m$ we mean a bijection from set $\{1, \dots, m\}$ onto itself. If π is a permutation then π^{-1} denotes the inverse permutation ($\pi(x) = y$ if and only if $\pi^{-1}(y) = x$). Hence, in our formalism, card $P_{i-1} P_i$ is followed by card $P_{\pi(i)} P_{\pi(i)+1}$.

The output of S^1 thus always has form

$$P_1 P_2 P_{\pi(2)} P_{\pi(2)+1} \cdots P_{\pi^{-1}(1)-1} P_{\pi^{-1}(1)}.$$

For instance, if we have $m = 3$ then the cards are $P_1 P_2, P_2 P_3, P_3 P_1$ and a shuffling which results in sequence $P_1 P_2 P_3 P_1 P_2 P_3$ is described by permutation $\pi(1) = 2, \pi(2) = 3, \pi(3) = 1$.

5.2.1. Acceptable permutations

For our purposes, not all permutations π are acceptable; let us formally denote by \mathcal{A} the set of all the acceptable permutations.

To define \mathcal{A} , we first introduce an auxiliary bipartite graph $G(\pi)$.

Definition 5.1. Graph $G(\pi)$ has vertex-set $V = V_1 \cup V_2$ where $V_1 = \{u_1, \dots, u_m\}$ and $V_2 = \{v_1, \dots, v_m\}$. The edge-set of $G(\pi)$ is the union of three disjoint perfect matchings of the vertex-set, namely:

1. The perfect matching M_1 consisting of the edges $\{u_i, v_i\}$.
2. The perfect matching M_2 consisting of the edges $\{u_{i+1}, v_i\}$.
3. The perfect matching M_3 consisting of the edges $\{u_{\pi(i)}, v_i\}$.

Definition 5.2. We construct a directed graph $G'(\pi)$ from $G(\pi)$ by first directing each edge of $M_2 \cup M_3$ from V_2 to V_1 , and then contracting each edge of M_1 .

Definition 5.3 (Of Set \mathcal{A} of All Acceptable Permutations). Permutation π is acceptable ($\pi \in \mathcal{A}$) if and only if the following two conditions are satisfied:

1. The directed graph $G'(\pi)$ has a directed Eulerian closed walk where the edges of M_2 and M_3 alternate. This condition is equivalent to saying that permutation π describes a rearrangement of the cards of S^1 into a sequence.
2. In the auxiliary graph $G(\pi)$, the union of the perfect matchings $M_2 \cup M_3$ contains many (at least m/c where $c \geq 2$ is a small constant) cycles. This condition is added in order to make the reconstruction of the input sequence hard; see the sections below.

The following observation about the graph $G(\pi)$ will be used in the analysis of the attacker problem.

Observation 5.4. Let $G(\pi)$ be as in Definition 5.1. For $v_i \in V_2$ let $s(v) = P_i$ be its associated segment. Then we have the following equality between cyclic sequences:

$$P_1 P_2 P_3 \cdots P_m = s(M_2(1)) s(M_2(2)) \cdots s(M_2(m)),$$

where $M_2(i)$ denotes the vertex of V_2 connected with $u_i \in V_1$ by an edge of M_2 .

For illustration we perform S^1 on the output sequence of the previous example (which would be the natural use of S^1 , as described later). Again in this example ignore the empty-spaces.

5.3. Procedure $S^{1+}(\omega, lb, ub)$

If the input of the procedure S^1 comes from several runs of the preparatory procedure S described above, then we need to modify S^1 in order to make its output generic, that is to intentionally preserve the attacker-confusing overlaps. This modified procedure is called S^{1+} .

We recall that S^1 repeats the whole blocks P_i , i.e. the output of S^1 is the cyclic sequence

$$P_1 P_2 P_{\pi(2)} P_{\pi(2)+1} \cdots P_{\pi^{-1}(1)-1} P_{\pi^{-1}(1)}.$$

Example 2: Procedure $S^1(\omega, lb, ub)$
 Input $\omega = \text{ngEalgorEpnEinfioriEformatiofocealingEpaape}$
 $\text{EisEafEthisEconEconceaEcithmtheEaipisEtoEproEpiEpa}$
 $\text{EinEpresentesetionEcoEentEanEilgogorithmapEaimEofEtim}$
 parameters $lb = 6$ and $ub = 8$.
 First, the input sequence is partitioned randomly into blocks of length 6, 7 or 8. The blocks are divided by '+' below:
 $\text{ngEalgo+rEpnEinf+oriEfor+matiofo+cealing+Epaape+}$
 $\text{rEisEa+fEthisEc+onEconcea+aEcithm+theEaiPi+sEtoEpro}$
 $\text{+EpiEpa+apeEin+Epresent+esetionE+coEentE+anEilg+}$
 $\text{ogorith+mapEai+mEofEtim+}$
 Next we add overlap in front of each block. For procedure S^1 the overlap is always the whole preceding block. We get the following cards; to make the example easier to understand we indicate by '*' the division of each card into two blocks:
 $\text{mEofEtim*ngEalgo+ngEalgo*rEpnEinf+rEpnEinf*oriEfor+}$
 $\text{oriEfor*matiofo+matiofo*cealing+cealing*Epaape+}$
 $\text{Epaape*rEisEa+rEisEa*fEthisEc+fEthisEc*onEconcea+}$
 $\text{onEconcea*aEcithm+aEcithm*theEaiPi+theEaiPi*sEtoEpro+}$
 $\text{sEtoEpro*EpiEpa+EpiEpa*apeEin+apeEin*Epresent+}$
 $\text{Epresent*esetionE+esetionE*coEentE+coEentE*anEilg+}$
 $\text{anEilg*ogorith+ogorith*mapEai+mapEai*mEofEtim+}$
 Finally the output is given by rearranging the cards by an acceptable permutation, i.e. by a permutation whose corresponding bipartite graph consists of a lot of cycles. The smallest length of a cycle is 4. It is not difficult to see that the following permutation π creates nine 4-cycles and one 6-cycle. In the following description of π , the cycles are grouped together; for instance the first 4-cycle has edges $(v_1, u_{10}), (v_9, u_2), (v_{11}, u_2), (v_9, u_{10})$. The first two of them belong to perfect matching M_3 , the last two belong to perfect matching M_2 .
 $[\pi(1) = 10, \pi(9) = 2]; [\pi(2) = 6, \pi(5) = 3]; [\pi(3) = 9, \pi(8) = 4]; [\pi(7) = 13, \pi(12) = 8]; [\pi(14) = 11, \pi(10) = 15]; [\pi(11) = 18, \pi(17) = 12]; [\pi(19) = 14, \pi(13) = 20]; [\pi(16) = 21, \pi(20) = 17]; [\pi(15) = 19, \pi(18) = 16]; [\pi(21) = 7, \pi(4) = 5, \pi(6) = 1]$.
 Hence the final sequence (for ease of understanding we preserve the separation symbols '*', which in reality would not be present) is:
 $\text{ngEalgo*rEpnEinfEpaape*rEisEaEpiEpa*apeEintheEaiPi*}$
 $\text{sEtoEproEthisEc*onEconceaEpnEinf*oriEforonEconcea*}$
 $\text{aEcithmEpresent*esetionEmEofEtim*ngEalgoaEcithm*}$
 $\text{theEaiPianEilg*ogorithapeEin*Epresentogorith*}$
 $\text{mapEaicoEentE*anEilgesetionE*coEentEsEtoEpro*}$
 $\text{EpiEpaEai*mEofEtimrEisEa*fEthisEc*matiofo*}$
 $\text{cealingioriEfor*matiofocealing*Epaape}$

We assume that the input ω of S^{1+} comes from repeated runs of procedure S and so ω contains a lot of segments of length o (the overlaps of runs of S) repeated at least twice; let us denote by R the set of all these segments.

Procedure S^{1+} starts as S^1 by partitioning of ω into blocks P_1, P_2, \dots, P_m .

The blocks of S^{1+} cut some of the segments from R . To reflect this, we write $P_i = r_{i-1}^T Q_i r_i^T$ where

- Segment r_{i-1}^T is an empty segment or a terminal segment of an element of R cut by the partition between blocks P_{i-1} and P_i .

- Segment r_i^T is an empty segment or an initial segment of an element of R cut by the partition between blocks P_i and P_{i+1} .

Summarizing this notation we write

$$P_1 P_2 \dots P_m = Q_1 r_1 Q_2 r_2 Q_3 r_3 \dots r_{m-1} Q_m r_m,$$

where each r_i is such an element of R that is cut by the blocks of S^{1+} , or an empty segment. Each $P_i = r_{i-1}^T Q_i r_i^T$ where $r_i = r_{i-1}^T r_i^T$.

The first difference of S^1 and S^{1+} is that the overlaps of S^{1+} are not the whole preceding blocks. Instead, the overlap added in front of block P_{i+1} is $Q_i r_i^T$. Hence, block P_{i+1} with the overlap added in front of it has form $Q_i r_i Q_{i+1} r_{i+1}^T$.

To make the cards of S^{1+} more generic (see the same step in the description of Procedure S), we change each such $Q_i r_i Q_{i+1} r_{i+1}^T$ into $Q_i r_i Q_{i+1} r'_{i+1}$ where r'_{i+1} is obtained from r_{i+1}^T by adding a segment so that r'_{i+1} has length o and is repeated elsewhere in ω .

Summarising, the output of S^{1+} has form

$$Q_1 * Q_2 * Q_{\pi(2)} * Q_{\pi(2)+1} * \dots * Q_{\pi^{-1}(1)-1} * Q_{\pi^{-1}(1)}^*,$$

where each $*$ stands for a segment of length o which is repeated (at least) twice in this output, or the empty string. More specifically, if $*$ follows segment Q_i then it is equal to r_i or to r'_i .

5.4. Procedure $S^2(\omega, o)$

Let $S^2(\omega, o)$ be as follows: we assume its input is an output of S^1 , i.e. it is the cyclic sequence

$$P_1 P_2 P_{\pi(2)} P_{\pi(2)+1} \dots P_{\pi^{-1}(1)-1} P_{\pi^{-1}(1)}.$$

Note that in this sequence, each block P_i appears twice. Procedure S^2 first cuts each P_i randomly into P_i^1, P_i^2 so that length of P_i^1 is at least o , i.e. the whole overlap of length o , which we denote by o_i , is contained in P_i^1 . The trick of the concealing algorithm is that both copies of each P_i are cut in the same way! Let $o_i P_i^2$ denote P_i^2 with the added overlap.

For example, if P_i is equal to 'abcdefghijkl' and $o = 3$ then a possible cut of S^2 is 'abcde+ fghijkl'; P_i^1 is equal to 'abcde', P_i^2 is equal to 'fghijkl' and $o_i P_i^2$ is equal to 'cdefghijkl'.

We may describe the set of the cards of S^2 as the disjoint union of two sets $C_1 \cup C_2$, where

$$C_1 = \{o_1 P_1^2 P_1^1, o_2 P_2^2 P_2^1, \dots, o_m P_m^2 P_m^1\}$$

and

$$C_2 = \{o_1 P_1^1 P_{\pi(1)}^1, o_2 P_2^1 P_{\pi(2)}^1, \dots, o_m P_m^1 P_{\pi(m)}^1\}.$$

We remark here that the cards of C_1 correspond to the edges of perfect matching M_2 of graph $G(\pi)$ and the cards of C_2 correspond to the edges of perfect matching M_3 of $G(\pi)$ (see Definition 5.1).

Finally S^2 arranges $C_1 \cup C_2$ into a random cyclic sequence.

For illustration we perform S^2 on the output sequence of the previous Example 2 (which would be the natural use of S^2 , as described later). Again in this example ignore the empty-spaces.

Example 3: Procedure $S^2(\omega, o)$

Input $\omega = \text{ngEalgo}^* \text{rEpnEinfEpaape}^* \text{rEisEaEpiE}^* \text{apeEintheEaipi}^* \text{sEtoEprofEthisE}^* \text{conEconcerEpnEinf}^* \text{oriEforonEconce}^* \text{aEcithmEpre}^* \text{esetionEmEofEtim}^* \text{ngEalgoaEci}^* \text{thm}^* \text{theEaipiEilg}^* \text{ogorithape}^* \text{Ein}^* \text{Epre}^* \text{sentogorith}^* \text{mapEaicoEent}^* \text{anEilgesetionE}^* \text{coEentEsEtoEpro}^* \text{EpiEpmE}^* \text{mapEai}^* \text{mEofEtimrEisEa}^* \text{fEthisEcmatiofo}^* \text{cealingoriEfor}^* \text{matiofocealing}^* \text{Epaape}$,
parameter $o = 3$.

A consistent partitioning into blocks is indicated below:
 $\text{ngEa+lgo}^* \text{rEpnEi+nEpa+ape}^* \text{rEis+EaEpiE}^* \text{ape+E}^* \text{EintheEa+ipi}^* \text{sEtoE+profEthisE+c}^* \text{onEco+nEconcerEpnEi+n}^* \text{f}^* \text{ori+EforonEco+n}^* \text{ce}^* \text{aEci+thmEpre+sent}^* \text{eseti+onEmEofE+tim}^* \text{ngEa+lgoaEci+thm}^* \text{theEa+ipianEi+l}^* \text{g}^* \text{ogori+thape+Ein}^* \text{Epre+sentogorith}^* \text{mapE+aicoEe+ntE}^* \text{anEi+lgeseti+onE}^* \text{coEe+ntEsEtoE+pro}^* \text{EpiE+pmE}^* \text{mapE+ai}^* \text{mEofE+timrEis+Ea}^* \text{fEthisE+cmati+ofo}^* \text{ceal+ingori+Efor}^* \text{matio+ofoceal+ing}^* \text{Epa+ape}$

Next we add overlap (of length o) in front of each block (and we delete the 'helpful symbol' *):

$\text{apengEa+gEalgorEpnEi+nEinfEpa+ape}^* \text{rEis+EisEaEpiE}^* \text{pisEpa}^* \text{apeEintheEa+eEaEpiEtoE+toEprofEthisE+isEconEco+EconcerEpnEi+nEinf}^* \text{oriEforonEco+Econce}^* \text{aEci+EcithmEpre+presenteseti+etionEmEofE+ofEtimngEa+gEalgoaEci+EcithmtheEa+eEaipiEi+nEilgogori+orithape+apeEinEpre+presentogorith}^* \text{mapE+apEaicoEe+oEentEanEi+nEilgeseti+etionEcoEe+oEentEsEtoE+toEproEpiE+pisEpmE+apEaimEofE+ofEtimrEis+EisEafEthisE+isEcmati+ofoceal+ealingori+oriEformati+atiofoceal+ealingEpa+Epaape}$

Finally we rearrange the cards in a random order. The resulting sequence is as follows:

$\text{apengEagEalgorEpnEi+nEinfEpaatiofocealaperEisEisEaE}^* \text{pispisEpaEisEafEthisEapeEintheEanEilgesetieEapisEtoE}^* \text{toEprofEthisEisEconEcoEconcerEpnEiEpaapenEinf}^* \text{orisEcmatiofocealealingorioriorEforonEcopisEpmE}^* \text{EconceEciE}^* \text{cithmEpre}^* \text{presentesetiofEtimngEagEalgoaEciEoriEformatiE}^* \text{cithmtheEaeEaipiEi+nEilgogoriorithapeapeEinE}^* \text{pre}^* \text{presentogoriorithmapEetionEmEofEapEaicoEeealingE}^* \text{paoEentEanEisetionEcoEeoEentEsEtoEtoEproEpiEapEaimE}^* \text{ofEofEtimrEisEisEcmatiofocealealingori}$

5.5. Procedure $S^{2+}(\omega, o)$

We assume its input is an output of S^{1+} . This procedure is defined analogously as S^2 with the only difference that the cuts are performed to segments Q_i instead of segments P_i .

6. The concealing algorithm

Let the input string be ω , and the length of the preserved segments be k . We consider two scenarios, *weak concealing* and *strong concealing*, depending on the nature of the input. We perform the *weak concealing algorithm* if the input is nonspecific, i.e., short segments have many possible alternative prolongations, or there does not exist any outside knowledge about the likelihood of presence of some segments in the input (e.g. an English text).

The *weak concealing algorithm* may be described as

$$\omega_F = S^2(S^1(\omega, 3k/2, 2k), k - 1).$$

We choose to have the block length in S^1 longer and to overlap the whole blocks in S^1 since we want to ensure that the cuts of S^2 may be done in the same way in each of the two copies of the blocks P_i .

The *strong concealing algorithm* may be written as

$$\omega_F = S^{2+}(S^{1+}(S \dots S(\omega, k - 1, k, 3k/2)), 3k/2, 2k, k - 1,$$

where the number of repetitions of procedure S is application specific.

7. Analysis of the concealing algorithm

Observation 7.1. The concealing algorithm preserves all segments of length k present in the input sequence ω within the output sequence ω_F .

This observation is straightforward as whenever any of the above procedures cuts the input string, an overlap of length $o \geq k - 1$ is added in front of the segment following the cut, thus preserving all sub-segments of length k which would otherwise be partitioned by the cut.

It is also straightforward that both weak and strong concealing algorithms are linear in $|\omega|$ if we have

- Access to a generator of random permutations of the numbers less than $|\omega|$.
- Access to a generator of random elements of \mathcal{A} (see Definition 5.3).

A random permutation may be generated in linear time (see [39]). We will not discuss the complexity of generating random elements of \mathcal{A} . Instead, we specify a large subset \mathcal{B} of \mathcal{A} such that generating a random element of \mathcal{B} may be reduced to generating a random permutation of a number less than $|\omega|$.

Each element of \mathcal{B} may be constructed as follows: we take any permutation π of $m/2$ (we assume m is even) and we consider the pairing $P(\pi)$ of $\{1, 2, \dots, m\}$ given by $(1, \pi(1) + m/2), \dots, (m/2, \pi(m/2) + m/2)$. This pairing may be looked at as an involution $i(\pi)$ (a permutation α is involution if $\alpha(\alpha(x)) = x$ for each x) on m .

Finally, we construct $\beta = \beta(\pi)$ by shifting $i(\pi)$ by 1, i.e., by letting $\beta(a) = i(\pi)(a) + 1$ modulo m . We let β belong to \mathcal{B} if the following additional condition is satisfied: if $O(a)$ denotes $\beta(a + 1)$ then $O^j(1) \neq 1$ for $j < m$. This condition makes sure that the first condition of the definition of the acceptable permutation (Definition 5.3) is satisfied. The following observation is straightforward.

Observation 7.2.

$$|\mathcal{B}| \geq (m/2 - 1)!$$

Further, the graphs defined by a permutation from \mathcal{B} are disjoint unions of $m/2$ cycles of length 4. Generating a random element from \mathcal{B} is as hard as choosing a random permutation of $m/2$.

The following observation is also straightforward.

Observation 7.3. The length of the output of each of the procedures applied to input ω is linear in $|\omega|$. For example, for S and S^1 it is $2|\omega|$.

Remark 7.4. Preliminary computational experiments with English text suggest that the algorithm behaves well with respect to the property IV of the information concealing problem.

8. Hardness of the attacker problem

We recall that the attacker problem introduced in Section 4 reads:

How much information about ω can an attacker deduce from $|\omega|$, k , the knowledge of the concealing algorithm, and the complete list of the frequencies of repeats of segments of ω_F ?

For instance, the attacker can try to get all the overlaps of S^2 since assuming ω_F has no accidental repeats these overlaps appear exactly four times in ω_F and no other segment is like that. The attacker may partition ω_F into cards as indicated by all these overlaps. She gets a collection of cards, with $(k-1)$ -length segments marked in the beginning and the end of each card. The attacker wants to overlap these marked segments. Depending on whether ω_F has accidental repeats, the attacker possibly cuts in more places than were the original cards used in the algorithm. Hence, in her collection of cards some overlaps should not have been considered, and some segments have overlaps with more than one other card. These considerations naturally specify the *domino* and *donkey* problems.

8.1. Donkey problem

In the *donkey* problem we assume that ω_F has no accidental repeats. What does the attacker get? There are two versions of the algorithm. Let us first consider the *strong concealing* where the preliminary step is performed.

(1) As described above, using the 4-repeats of length $k-1$ of ω_F , the attacker gets the cards of S^{2+} , i.e. $C_1 \cup C_2$, where

$$C_1 = \{o_1 Q_1^2 r_1 Q_2^1, o_2 Q_2^2 r_2 Q_3^1, \dots, o_m Q_m^2 r_m Q_1^1\}$$

and

$$C_2 = \{o_1 Q_1^2 r_1' Q_{\pi(1)}^1, o_2 Q_2^2 r_2' Q_{\pi(2)}^1, \dots, o_m Q_m^2 r_m' Q_{\pi(m)}^1\}.$$

(2) The attacker also gets each Q_i^1 and each $o_i Q_i^2$ since these are exactly maximal initial and terminal segments of the cards above which are repeated twice in ω_F .

(3) By matching the overlaps, the attacker gets each pair $Q_i^1 Q_i^2$ since the overlap o_i in $o_i Q_i^2$ is a terminal segment of Q_i^1 and we may assume that these cannot be misinterpreted.

(4) What does the attacker get from the initial applications of procedure S ? Each of their overlaps (of length $k-1$) appears at least twice in the input of S^{1+} . Moreover most of the cuts of the procedures S are different. Let us recall here that among these overlaps may be also the dust. Procedures S^{1+} and S^2 cut into some of these. Those cut will remain 2-repeats, those not cut may gain repeats. Moreover, S^{1+} introduces dust in the border of each card: this adds 2-repeats of strings of length $k-1$ indistinguishable from the 2-repeats coming from initial procedures S .

In the case where the weak concealing algorithm is applied, the attacker has 1., 2., 3. where $Q_i^2 r_i$ and $Q_i^2 r_i'$ are replaced by P_i^2 and Q_i^1 is replaced by P_i^1 .

The next proposition summarises the possible types of repeats introduced by the algorithm.

Proposition 8.1. All the repeats of ω_F generated by the weak or strong concealing algorithm are those described in (1), (2), (3) and (4).

Corollary 8.2. All the useful information for the attacker problem is $|\omega|$, k , and (1), (2), (3) and (4).

The information (1), (2), and (3) may be described by the auxiliary bipartite graph $G(\pi)$ defined in Definition 5.1.

If the weak concealing algorithm is applied, information (4) does not exist. The attacker problem is thus reduced to the following:

The donkey-decision problem. The input is a bipartite graph G where the vertices in both parts V_1, V_2 are ordered. Let $V_1 = \{u_1, \dots, u_m\}$ and $V_2 = \{v_1, \dots, v_m\}$. Moreover a segment $s(v)$ of length at least $3k/2$ is associated with each element of V_2 . The set of the edges of G is formed by a disjoint union of two perfect matchings M_2, M_3 . The attacker needs to reconstruct string

$$s(M_2(1))s(M_2(2)) \dots s(M_2(m)),$$

where $M_2(i)$ is the vertex of V_2 connected with $u_i \in V_1$ by an edge of M_2 .

The difficulty of the donkey-decision problem is the following: bipartite graph G is a union of two edge-disjoint perfect matchings. Each vertex of G thus has degree 2 and G is a union of disjoint cycles. To solve the donkey-decision problem, one needs to choose the correct perfect matching independently in each of these cycles (namely, the perfect matching induced by M_2). This is impossible, and the list of all the possibilities is almost always exponential in the number of the cycles, since each of the cycles has two perfect matchings. This is analysed precisely below, when we speak about the *feasible solutions*.

Next we argue that, when the strong concealing algorithm is applied, the attacker problem is reduced to the donkey-decision problem too. The attacker is left with the statistics of the repeats of ω_F . Here comes the reason why we introduced the dust in S^{1+} : it is to make sure that the 2-repeats appear symmetric for both matchings M_2, M_3 . This hides the repeats introduced by the initial applications of procedure S . The information of [4.] is thus useless. We obtain:

Proposition 8.3. The attacker problem for both strong and weak concealing is reduced to the analysis of the donkey-decision problem.

A *feasible solution* to the donkey-decision problem is any sequence $s(M(1))s(M(2)) \dots s(M(m))$, where M is any perfect matching of the input bipartite graph G . In order to solve the donkey-decision problem, one needs to choose, from the pool of these feasible solutions, the unique correct one. Next we argue that unless the input to our problem is extremely restrictive, there is an exponential number of the competitive solutions.

The bipartite graphs G coming from \mathcal{A} have at least $2^{m/c}$ perfect matchings. The output sequences of two perfect

matchings M, N may still be equal: if the cycle has length 4, this happens if and only if the two vertices v_i, v_j of V_2 in each 4-cycle in which M, N differ have the same associated segment ($s(v_i) = s(v_j)$) as defined in [Observation 5.4](#).

For instance, if all the vertices of V_2 have the same associated segment, then there is only one competitive solution. This extreme situation may happen if the input ω is a sequence of repetitions of one symbol only.

If two symbols may appear in the segments (of length at least $3k/2$) associated with the vertices of V_2 , then the probability that in a 4-cycles the corresponding pairs of strings are indistinguishable is $2^{-3ka/2}$. Hence with only exponentially small probability there is less than an exponential number of feasible solutions.

8.2. Domino problem

In a more realistic situation the attacker does not know the correct list of cards of S^2 and hence she needs to choose which 4-repeats to ignore. We may assume that she has some hints as to which overlaps are ‘likely’ ok. This is the situation we model by the following problem.

Shortest domino row problem (SDRP). Assume we are given a collection of dominoes (domino will mean a rectangle partitioned vertically into two squares, where one is initial and the other one is terminal), and we are also given a graph on the squares. This graph should be interpreted as the graph of hints. We want to put all the dominoes into a row, so that if two consecutive squares are connected by an edge of the graph, we can put one square on top of the other (i.e., identify them). The aim is to make the resulting row as short as possible, i.e. to satisfy as many hints as possible.

Let us define the (de Bruijn-type) graph $G = (V, E)$ where V is the set of all the squares, and E is the set of the dominoes: edge e_i connects the squares of domino Q_i . The following observation is straightforward.

Observation 8.4. *There is a natural bijection between the set of the Euler circuits (Eulerian closed walks) of G and the set of all the circular sequences consistent with the overlapping dominoes Q_1, \dots, Q_m .*

Theorem 1. *The SDRP is search-NP-complete.*

Proof. Assume that in the auxiliary graph, there is an edge between two squares if they are equal, but not all such edges are there. This is exactly consistent with our interpretation. Now, in the reformulation with the de Bruijn graph and the Euler circuit, this corresponds to the problem that we are given a graph, with some transitions between neighboring edges recommended, and we want to find an Euler circuit with as many recommended transitions as possible. A particular instance is that some transitions are forbidden, and we want to find out whether an Euler circuit where all the transitions are allowed exists. This is known to be NP complete [40].

We have in fact a *search instance* of this problem: we know that such an Euler circuit exists, and we want to find it. There is a standard trick which shows that the decision problem is polynomial if the search problem is polynomial:

Assume there is a polynomial algorithm A that solves the search version, and let its running time be n^{10} , say. To solve the decision problem, we apply A to an input. It either finds the right Euler circuit and then the answer is YES, or it runs longer than n^{10} , and then the answer is NO. \square

9. Conclusion

We define the information concealing problem and propose an algorithm to solve it. It is based on a seminal observation coming from the difficulties of DNA reconstruction by hybridisation. The authors along with Jenny Blamey proposed in [4] that in eukaryotes the cells have DNA as a depositary of concealed genetic information and the genome achieves the self-concealing by accumulation and maintenance of repeats. The protected information may be shared and this is useful for the development of intercellular communication and in the development of multicellular organisms. It is proposed further in [4] that the multicellular organisms have a mechanism to maintain the families of repeats. Observations about the development of the repeats may help to explain some other spreading mechanisms in multicellular organisms, and disease.

In other words, we propose that DNA contains repeats in order to solve the information concealing problem. We want to stress that even though our algorithm is inspired by the repeats in DNA, the output of the algorithm, aside of containing repeats, has no other similarities to the DNA structure as far as we know. The ‘concealing’ text fraction of the output of the algorithm is higher than the 50% in eukaryotes. This may reflect the fact that our algorithm is a general method.

The presented algorithm may be efficiently implemented. In analysing the amount of information leaked by the concealing algorithm to an attacker (this is called the *attacker problem* in the paper), it is shown that with very high probability there is an exponential number of feasible solutions which are indistinguishable from the available information, among which the attacker needs to choose the correct one.

REFERENCES

- [1] R. Ramaswamy, L. Kencl, G. Iannaccone, Approximate fingerprinting to accelerate pattern matching, in: Proceedings of the ACM Internet Measurement Conference, IMC, Rio de Janeiro, Brazil, 2006.
- [2] H.-A. Kim, B. Karp, Autograph: toward automated, distributed worm signature detection, in: Proceedings of the 13th Usenix Security Symposium, Security 2004, San Diego, CA, August 2004.
- [3] S. Singh, C. Estan, G. Varghese, S. Savage, Automated worm fingerprinting, in: Proceedings of the ACM/USENIX Symposium on Operating System Design and Implementation, OSDI, San Francisco, CA, USA, 2004.
- [4] J. Blamey, L. Kencl, M. Loebel, DNA self-concealing by repeats (2009) (in preparation).
- [5] Y. Benenson, R. Adar, T. Paz-Elizur, et al., Proceedings of the National Academy of Sciences 100 (5) (2003) 2191–2196.
- [6] R. Guy-Franck, F. Paques, EMBO Reports 1, 2000, pp. 122–126.

- [7] P.A. Pevzner, *Computational Molecular Biology: An Algorithmic Approach*, The MIT Press, Cambridge, Massachusetts, London, England, 2000.
- [8] R. Arratia, B. Bollobas, D. Coppersmith, G.B. Sorkin, Euler circuits and DNA sequencing by hybridization, *Discrete Applied Mathematics* 104 (2000) 63–96.
- [9] M. Singer, P. Berg, *Genes and Genomes*, University Science Books Sausalito, California, 1991.
- [10] S. Mirkin, Expandable DNA repeats and human disease, *Nature* 447 (2007) 932–940.
- [11] J. Buard, A.J. Jeffreys, Big, bad minisatellites, *Nature Genetics* 5 (1997).
- [12] J.D.H. Stead, A.J. Jeffreys, Allele diversity and germline mutation at the insulin minisatellite, *Human Molecular Genetics* 9 (2000) 713–723.
- [13] S. Mirkin, Molecular models for repeat expansions, *Chemtracts Biochemistry and Molecular Biology* 17 (2004) 639–662.
- [14] C.E. Pearson, K. Nichol Edamura, J.D. Cleary, Repeat instability: mechanisms of dynamic mutations, *Nature Reviews Genetics* 6 (2005) 729–742.
- [15] L.P. Ranum, T.A. Cooper, RNA-mediated neuromuscular disorders, *Annual Review of Neuroscience* 29 (2006) 259–277.
- [16] A. Abu-Baker, G.A. Rouleau, in: R.D. Wells, T. Ashizawa (Eds.), *Genetic Instabilities and Neurological Diseases*, Elsevier, Amsterdam, 2006, pp. 487–513.
- [17] J. Xu, J. Fan, M.H. Ammar, S.B. Moon, Prefix-preserving IP address anonymization: measurement-based security evaluation and a new cryptography-based scheme, in: *Proceedings of the IEEE International Conference on Network Protocols, ICNP, 2002*.
- [18] D. Maltz, J. Zhan, G. Xie, H. Zhang, G. Hjalmtysson, J. Rexford, A. Greenberg, Structure preserving anonymization of router configuration data, in: *Proceedings of the ACM Internet Measurement Conference, IMC, Taormina, Italy, 2004*.
- [19] R. Pang, V. Paxson, A high-level programming environment for packet trace anonymization and transformation, in: *Proceedings of ACM SIGCOMM, 2003*.
- [20] B. Bloom, Space/time tradeoffs in hash coding with allowable errors, in: *CACM*, p. 422.
- [21] K. Shanmugasundaram, H. Broennimann, N. Memon, Payload attribution via hierarchical Bloom filters, in: *Proceedings of the ACM Conference on Computer and Communications Security, CCS, Washington, DC, USA, October 2004*.
- [22] Y. Li, J. Tygar, J.M. Hellerstein, Private matching, in: *Proceedings of the 30th VLDB Conference, Toronto, Canada, 2004*.
- [23] R. Agrawal, A. Evfimievski, R. Srikant, Information sharing across private databases, in: *Proceedings of the ACM SIGMOD International Conference on Management of Data, San Diego, CA, USA, 2003*.
- [24] S. Garriss, M. Kaminsky, M.J. Freedman, B. Karp, D. Mazires, H. Yu, Re: Reliable email, in: *Proceedings of the Symposium on Networked Systems Design and Implementation, NSDI, San Jose, CA, USA, 2006*.
- [25] Net 2000 Ltd., Data masking whitepapers, 2005. <http://www.datamasker.com/>.
- [26] G. Miklau, D. Suci, A formal analysis of information disclosure in data exchange, in: *Proc. of ACM SIGMOD Intl. Conf. on Management of Data, 2006*.
- [27] R. Agrawal, R. Srikant, Privacy-preserving data mining, in: *Proc. of ACM SIGMOD Intl. Conf. on Management of Data, 2000*.
- [28] S. Agrawal, J.R. Haritsa, A framework for high-accuracy privacy-preserving mining, in: *Proceedings of the 21st International Conference on Data Engineering ICDE, 2005*.
- [29] A. Evfimievski, J. Gehrke, R. Srikant, Limiting privacy breaches in privacy preserving data mining, in: *Proc. of ACM Symp. on Principles in Database Systems PODS, 2003*.
- [30] N. Mishra, M. Sandler, Privacy via pseudorandom sketches, in: *Proc. of ACM Symp. on Principles in Database Systems PODS, 2006*.
- [31] V. Rastogi, D. Suci, S. Hong, The boundary between privacy and utility in data publishing, in: *Proc. of Intl. Conf. on Very Large Data Bases VLDB, 2007*.
- [32] M. Hay, M. Miklau, G.D. Jensen, S. Weiss, et al., Anonymizing social networks, University of Massachusetts, Amherst, Technical Report, 2007.
- [33] L. Sweeney, *k*-anonymity: a model for protecting privacy, *International Journal on Uncertainty, Fuzziness and Knowledge-Based Systems* 10 (5) (2002) 557–570.
- [34] N.F. Johnson, Z. Duric, S. Jajodia, *Information Hiding: Steganography and Watermarking—Attacks and Countermeasures*, Springer, 2000.
- [35] I. Cox, M. Miller, J. Bloom, J. Fridrich, T. Kalker, *Digital Watermarking and Steganography*, second ed., in: *The Morgan Kaufmann Series in Multimedia Information and Systems, 2007*.
- [36] C.D. Manning, P. Raghavan, H. Schütze, *Introduction to Information Retrieval*, Cambridge University Press, 2008.
- [37] C. Gentry, Computing arbitrary functions of encrypted data, *Communications of the ACM* 53 (3) (2010) 97–105.
- [38] L. Kencl, J. Zamora, M. Loebl, Packet content anonymization by hiding words, in: *Demo at IEEE INFOCOM, 2006*.
- [39] D.E. Knuth, *The Art of Computer Programming, III ed.*, Addison-Wesley, Reading, Massachusetts, 1997.
- [40] M.R. Garrey, D.S. Johnson, *Computers and intractability: a guide to the theory of NP-completeness*.
- [41] L. Kencl, M. Loebl, DNA-inspired information concealing, [arXiv:0904.4449](https://arxiv.org/abs/0904.4449).