

Randomized Data Partitioning with Efficient Search, Retrieval and Privacy-Preservation

M. Oğuzhan Külekci¹[0000-0002-4583-6261]

Dept. of Computer Science, Indiana University Bloomington, IN, USA
okulekci@iu.edu

Abstract. We introduce a new data representation that serves mainly for privacy preserving data storage with efficient search and retrieval capabilities over the distributed systems. The cornerstone of the proposed scheme is based on a novel algorithm that splits an input bit sequence $B[1..n]$ into two as left and right partitions with well-control over the the partition sizes, and the reconstruction of B in absence of either partition is hard to achieve. The algorithm processes the input bit stream in blocks of d -bits, where initially each block is replaced with another d -bit according to a randomly chosen permutation of the set $\{0, 1, \dots, 2^d - 1\}$. Following the replacement, the leftmost bits of each block up until and including the q th set bit are appended to the left and the remaining bits to the right partition. We prove that the expected length of the left partition is $\ell \approx 2qn/d$ bits and the right partition becomes of length $|R| = n - \ell$ bits. Therefore, there is no overhead on the new representation with respect to original input. We also show that due to the randomization step, the input data B is not required to follow any special probability distribution to have the mentioned partitioning ratio $\rho = 2q/d$ and it is possible to tune the parameters d and q to support any desired ratio ρ on the input. We consider recursive application of that splitting algorithm on each partitions, which can be viewed as generating a full binary tree with k -leaves such that at each internal node the data is subject to the proposed splitting operation. Such a construction represents an input bit sequence $B[1..n]$ with k partitions as P_1, P_2, \dots, P_k , where it is hard to reconstruct the original data in absence of any P_i .

Keywords: randomization · data representation · data partitioning · distributed data storage · cloud storage · privacy-preserving data representation

1 Introduction

We introduce a randomized data representation that preserves the privacy of the data, while still supporting efficient access and retrieval. The building block of the proposed scheme is a novel randomized algorithm that splits an input n -bit long sequence into two pieces, which we refer as the *left* and *right* partitions. The algorithm processes the input in blocks of a predetermined length d bits. The initial randomization step replaces each block on the input with its counterpart

according to a randomly generated permutation [8] of d -bit integers. Following the replacement, the leftmost bits of each block up until and including the q th set bit are appended to the left, and the remaining bits to the right partition.

We prove that the randomization step guarantees the *expected* size of the left partition to be $\ell = n \cdot \rho$ bits, while the right partition becomes of length $(n - \ell)$ bits according to the input parameter $0 < \rho < 1$ such that the q value can be chosen to satisfy $\rho = 2q/d$. Access to any block on the original input can be achieved efficiently with the intrinsic addressing in the construction. We show that it is also possible to search for a queried pattern by filtering on one partition and verifying on the other as well.

The privacy of the data in the proposed representation depends on combination of several factors, which can be summarized as follows. First, the random permutation used in the construction is assumed to be secret, which can be maintained by using a secret initial seed in the random number generator [8]. It is known that the permutations are vulnerable to statistical attacks [1, 16], but due to the splitting mechanism it becomes hard to make a frequency analysis on the left and right partitions. On the right partition, each block is represented with variable-length bits, where the code-word boundaries cannot be determined as these codes are not prefix-free, that lacks the possibility of frequency analysis. Actually, the left partition encodes those code-word boundary information. On the left partition, the code-word boundaries can be determined, however, since multiple randomly chosen symbols are represented with the same code, again running a frequency analysis becomes hard as 2^d symbols on the input are mapped to d symbols on the left partition, which introduces again an ambiguity. Therefore, in absence of either partition and the secret permutation it is hard to reconstruct the original input, that provides the privacy in the proposed scheme. It is noteworthy that even the prefix-free codes have been found to be hard to decode without proper information [9, 18, 10], where recent works on combining compression and security [4, 7] considers the hardness of decoding asymmetric numeral codes [6], which produce variable-length non-prefix-free codes.

By applying the proposed split operation recursively to each partition, we observe that an input sequence can be split into more than two partitions. This can be achieved by generating a full binary tree with k -leaves, where the root node is the original input data. Each internal node has two children and the data on internal nodes are subject to the proposed splitting operation. Notice that in such a multiple partition it becomes more and more difficult to analyze and reproduce the original in absence of any partition as deeper the tree the ambiguity in the partitions increases.

There might be several applications of the proposed scheme, where a basic scenario is that a user wants to store a large file on the cloud or on a decentralized network such as the IPFS [3], where keeping data secret from the owners of the storage medium is a serious concern. Splitting data into multiple pieces and encrypting those partitions before saving to the cloud or the distributed network provides the ultimate secrecy. However, search and retrieval, which are vital operations on any data management scheme, becomes difficult to achieve on

such encrypted data. The encryption algorithms, mainly the homomorphic [14] and multi-party computation [5] schemes, support those operations, but their computational load might be prohibiting, particularly on distributed networks or IoT environments with limited resources [11, 19]. On such a case, the proposed scheme would serve as an alternative solution.

Such a scenario have been previously considered and data splitting have been proposed by Li *et al.* [13], where the authors proposed to use a key and create two bit streams from the input by using this key. It becomes hard to decode the original again in absence of the partitions and the key. Each bit stream is of equal length with the input, and hence, total storage requirement increases to $2n$ bits for n bit input. Yet another point is although random access might be supported, the search operation is not possible. Therefore, when compared to [13], the randomized algorithm in this study becomes advantageous in two dimensions by not creating an overhead storage and supporting efficient search and access.

The outline of the paper is as follows. In section 2 we introduce the proposed randomized data splitting algorithm and provide the proofs of well-control over the sizes of the partitions. The support mechanisms for efficient random access and search operations are explained in section 3, which is followed by the analysis of the privacy aspects in section 4. We describe how to partition a data into multiple splits according to desired proportions in section 5. We conclude by a summary and discussions of possible further research directions.

2 The Split Coding and Its Properties

Definition 1 (Split-Coding). *The split coding of an input d -bit sequence $A = a_1a_2..a_d$ according to the parameter $q, d > q > 0$, is denoted by $S(A, q) \rightarrow \langle L_A, R_A \rangle$ and produces the left partition L_A and right partition R_A as follows*

- i) If A has at least $q-1$ bits, then $L_A = a_1a_2..a_i = (0^*1)^q$ and $R_A = a_{i+1}..a_d = (0|1)^{d-i}$*
- ii) Otherwise, $L_A = a_1a_2..a_d = A$ and $R_A = \emptyset$ (empty).*

*,where (0^*1) denotes any number of 0 bits followed by a 1 bit. $(0|1)$ means either 0 or 1 bit. The superscripts over them represent how many times they are repeated.*

The split-coding that partitions an input d -bit integer into two according to a parameter q is given in Definition 1, where some examples are $S(01001101, 1) \rightarrow \langle 01, 001101 \rangle$, $S(01001101, 2) \rightarrow \langle 01001, 101 \rangle$, $S(01001101, 4) \rightarrow \langle 01001101, \emptyset \rangle$. Let the n -bits long input bit sequence is shown by $B[1..n] = b_1b_2..b_n$ and $d \in (0, \log n)$ is a chosen fixed block-length, which we assume $d|n$, as B can be padded with random bits otherwise. The input B can be shown by $B = D_1D_2..D_m$, where $D_i = B[d(i-1)+1..i \cdot d]$ is a block for $1 \leq i \leq m$, and $m = n/d$. The d -bits long blocks are the binary encoding of the integers from the set $U = \{0, 1, \dots, 2^d-1\}$.

Let $\pi[0..2^d-1] = \{\pi_0, \pi_1, \dots, \pi_{2^d-1}\}$ be a *randomly generated* permutation of the set U . We apply the permutation π on the input $B = D_1D_2..D_m$ that

i	1	2	3	4	5	6	7	8
D_i	100	101	001	000	010	110	111	011
$D'_i = \pi(D_i)$	011	111	010	101	000	100	110	001
L_i	01	1	01	1	000	1	1	001
R_i	1	11	0	01		00	10	

Fig. 1. The split-coding of $B = 100101001000010110111011$ with $d = 3$, $q = 1$, and $\pi = \langle 5, 2, 0, 1, 3, 7, 4, 6 \rangle$ as $P(B, d, \pi, q) \rightarrow (L = 01101100011001, R = 1110010010)$. The decoding process to merge L_B and R_B partitions to restore B can also be traced on the example strings. For example, scanning the initial bits from L returns 2 bits sequence 01 as we hit a set bit. Therefore, we need to append $1 = 3 - 2$ bits from R , which creates $D'_1 = 011$. Applying the inverse permutation $D_1 = \pi^{-1}(3) = 4$ restores the original bits 100.

produces $B' = D'_1 D'_2 \dots D'_m$ such that $D'_i = \pi[D_i]$. The random permutation π can be constructed by using a pseudo-random number generator that is initialized with a selected *seed*. Following the permutation step, we apply split-coding to each d -bits long block $D'_{1 \leq k \leq m}$ that produces the left and right partitions as $S(D'_k, q) \rightarrow \langle L_k, R_k \rangle$ according to the chosen parameter q , $0 < q < d$. For each block, the left partition L_k is the sequence of leftmost bits up until and including the q th set bit, and the right partition R_k is the remaining bits succeeding the q th set bit. As stated in definition 1, if q th set bit appears on the last position or there are less than q th set bits on the block, then the right partition will be empty, and the block is copied into left partition and nothing is written to the right partition. We concatenate the left partitions L_1 to L_m in order and generate the L -partition as $L_B = L_1 L_2 \dots L_m$. Similarly, $R_B = R_1 R_2 \dots R_m$ denotes the R -partition. This complete process is represented by $\text{SplitEncode}(B, d, \pi, q) \rightarrow \langle L, R \rangle$, which is also shown in pseudocode Algorithm ?? in the appendix.

Given the L - and R -partitions, the permutation π , and the parameters q and d , the decoding process reconstructs B by simply scanning the L - and R -partitions from left-to-right. The decoding procedure starts by initializing the read pointers on L_B and R_B to their first positions, and reads bits from L_B -partition until q set bits are observed or total number of bits examined reaches d with less than q set bits. Assuming t bits are read on L_B , the remaining $(d - t)$ bits are read from the R_B -partition in order. The concatenation of these bits creates a d -bits long block, say T , from which the original block is restored via $\pi^{-1}[T]$. This procedure is repeatedly applied until all m blocks are restored as listed on Algorithm ?? in the appendix. A sample split encoding and decoding is depicted on Figure 1 with the parameters $q = 1$, $d = 3$, and $\pi = \langle 5, 2, 0, 1, 3, 7, 4, 6 \rangle$ on a bitmap $B = 100101001000010110111011$.

We now analyze the expected bit lengths of the partitions created by split-coding process on a n bits long input B and prove that the expected length of the left partition is $n \cdot 2q/d$ bits long, which leaves the right partition to be $\approx n \cdot (d - 2q)/d$ bits long.

Proposition 1 (Length Invariance). *The sum of the lengths of the L- and R- partitions created by the split-coding $\text{SplitEncode}(\mathbf{B}, \mathbf{d}, \pi, \mathbf{q}) \rightarrow \langle \mathbf{L}_B, \mathbf{R}_B \rangle$ is equal to the length of the input bit string \mathbf{B} .*

Proof. Each \mathbf{d} -bits long block \mathbf{D}_i on \mathbf{B} is split into two bit strings \mathbf{L}_i and \mathbf{R}_i without any insertion or deletion of bits. Therefore, per each block, the lengths of its corresponding bits on \mathbf{L}_B and \mathbf{R}_B partitions sum up to \mathbf{d} as $\mathbf{L}_i + \mathbf{R}_i = \mathbf{d}$ for all $i = 1$ to \mathbf{m} . Thus, $\sum_{1 \leq i \leq \mathbf{m}} (\mathbf{L}_i + \mathbf{R}_i) = \mathbf{m} \cdot \mathbf{d} = |\mathbf{B}|$ holds. \square

Lemma 1 (Expected bit lengths of partitions for $\mathbf{q} = 1$). *For $\mathbf{q} = 1$ and $\mathbf{d} \geq 1$, the **expected** bit lengths of the L and R partitions of a randomly chosen integer from $\mathbf{U} = \{0, 1, \dots, 2^{\mathbf{d}} - 1\}$ are $\mathbf{E}(|\mathbf{L}|) = 2 - \frac{1}{2^{\mathbf{d}-1}}$ and $\mathbf{E}(|\mathbf{R}|) = \mathbf{d} - 2 + \frac{1}{2^{\mathbf{d}-1}}$.*

Proof. A randomly selected \mathbf{d} -bit integer i can take any value in $\{0, 1, \dots, 2^{\mathbf{d}} - 1\}$ with probability $1/2^{\mathbf{d}}$. For all values of $i \geq 2^{\mathbf{d}-1}$, the L partition will be equal to 1, and thus, $|\mathbf{L}| = 1$. Since there are $2^{\mathbf{d}-1}$ such i values, the probability that the length of the L partition will be one is $2^{\mathbf{d}-1}/2^{\mathbf{d}} = 1/2$. Similarly, if $2^{\mathbf{d}-2} \leq i < 2^{\mathbf{d}-1}$, then L will be 01 that will dictate $|\mathbf{L}| = 2$ with probability $2^{\mathbf{d}-2}/2^{\mathbf{d}} = 1/4$, which can be generalized to $|\mathbf{L}| = \ell$ with probability $1/2^\ell$ for $1 \leq \ell \leq (\mathbf{d} - 1)$. When $i = 0$ or $i = 1$, the L partitions will completely include i , and thus, the corresponding length $|\mathbf{L}|$ will be \mathbf{d} with probability $2/2^{\mathbf{d}}$. Therefore, the expected length $\mathbf{E}(|\mathbf{L}|)$ by the split coding can be calculated as follows.

$$\mathbf{E}(|\mathbf{L}|) = \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \dots + \frac{1}{2^{\mathbf{d}-1}} \cdot (\mathbf{d} - 1) + \frac{2}{2^{\mathbf{d}}} \cdot \mathbf{d} = \frac{2}{2^{\mathbf{d}}} \cdot \mathbf{d} + \sum_{\ell=1}^{\mathbf{d}-1} \ell \cdot \frac{1}{2^\ell} \quad (1)$$

The algebraic identity $\sum_{i=1}^n i \cdot r^i = \frac{r - r^{n+1}(1+n-nr)}{(1-r)^2}$ with $r=1/2$ and $n=\mathbf{d}-1$, reduces Equation 1 to

$$\mathbf{E}(|\mathbf{L}|) = \frac{2\mathbf{d}}{2^{\mathbf{d}}} + 4 \cdot \left(\frac{1}{2} - \frac{\mathbf{d}+1}{2^{\mathbf{d}+1}} \right) = \frac{2\mathbf{d}}{2^{\mathbf{d}}} + 2 - \frac{2\mathbf{d}+2}{2^{\mathbf{d}}} = 2 - \frac{1}{2^{\mathbf{d}-1}} \quad (2)$$

For any integer i the sum of the lengths of L and R sum up to \mathbf{d} as $|\mathbf{L}| + |\mathbf{R}| = \mathbf{d}$ by definition and $\mathbf{E}(|\mathbf{R}|) = \mathbf{E}(\mathbf{d} - |\mathbf{L}|) = \mathbf{d} - \mathbf{E}(|\mathbf{L}|) = \mathbf{d} - 2 + 1/2^{\mathbf{d}-1}$ holds. \square

Lemma 2 (Expected bit lengths of the partitions on a \mathbf{d} -bit block). *For any $0 < \mathbf{q} < \mathbf{d}/2$, the **expected** bit lengths of the L- and R- partitions of a randomly chosen \mathbf{d} -bit integer from $\mathbf{U} = \{0, 1, \dots, 2^{\mathbf{d}} - 1\}$ are*

$$\mathbf{q} \leq \mathbf{E}(|\mathbf{L}|) \leq 2\mathbf{q} - \frac{\mathbf{q}}{2^{\mathbf{d}-1}} \quad \text{and} \quad (\mathbf{d} - \mathbf{q}) > \mathbf{E}(|\mathbf{R}|) > \mathbf{d} - 2\mathbf{q} + \frac{\mathbf{q}}{2^{\mathbf{d}-1}}$$

Proof. Let $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_q$ denote the first \mathbf{q} set bit positions on the randomly chosen integer $i \in \{0, 1, \dots, 2^{\mathbf{d}} - 1\}$, and $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_q$ represent the differences in between the positions such that $\mathbf{x}_1 = \mathbf{p}_1$, and $\mathbf{x}_j = \mathbf{p}_j - \mathbf{p}_{j-1}$. For instance, if the positions

of the first four set bits are, say 2, 5, 6, 8, then the $\langle x_1, x_2, x_3, x_4 \rangle = \langle 2, 3, 1, 2 \rangle$. The length of the L partition is then $|L| = x_1 + x_2 + \dots + x_q$.

We observe that x_j is the length of the L-partition of a $d_j = d - \sum_{r=1}^{j-1} x_r$ bits long sequence. Due to Lemma 1, $E(x_j) = 2 - 1/2^{d_j-1}$. Since the lengths of the sequences decrease at each step as $d_1 > d_2 > \dots > d_q$, the expected values decline accordingly $E(x_1) > E(x_2) > \dots > E(x_q)$, e.g., for $d=8$, $E(|L|) = 2 - 1/2^7 = 1.992$, where for $d=6$, $E(|L|) = 2 - 1/2^7 = 1.984$. That defines $E(|L|)$ to be upper bounded by $q \cdot E(x_1)$ due to $E(x_1) > E(x_2) > \dots > E(x_q)$ as shown in equation 5. Therefore, the expected value $E(|L|)$ can be represented as in equation 3 due to the linearity of expectations.

$$E(|L|) = E(x_1 + x_2 + \dots + x_q) = \sum_{j=1}^q E(x_j) \quad (3)$$

$$= 2q - \left(\frac{1}{2^{d_1-1}} + \frac{1}{2^{d_2-1}} + \dots + \frac{1}{2^{d_q-1}} \right) \quad (4)$$

$$\leq q \cdot E(x_1) = 2q - \frac{q}{2^{d-1}} \quad (5)$$

By definition 1, the length of the L-partition is at least one bit, and thus, $x_i \geq 1$ holds for each i , which provides the lower bound for $E(|L|) \geq q$. Again by using $|L| + |R| = d$, $E(|R|) = d - E(|L|)$, the upper- and lower-bounds for $E(|R|)$ is obtained. \square

Theorem 1. *The split-coding of a given n bits long B with a randomly chosen permutation π , block length d , and parameter $0 < q \leq d/2$ generates left partition L_B with an expected length of $E(|L_B|) \leq n \cdot \left(\frac{2q}{d} - \frac{q}{d \cdot 2^{d-1}} \right)$. Accordingly, the right partition is expected to be $E(|R_B|) > n \cdot \left(1 - \frac{2q}{d} + \frac{q}{d \cdot 2^{d-1}} \right)$ bits long.*

Proof. Split-coding operates on the n -bits long input sequence by n/d blocks, each of which is a d -bit integer. The blocks are replaced by their corresponding integers in the set $\{0, 1, 2, \dots, 2^d - 1\}$ with the the randomly generated permutation π . Let X_1, X_2, \dots, X_{2^d} denote the lengths of the L-partitions of randomly assigned integers to each distinct d -bit integers and f_1, f_2, \dots, f_{2^d} represent the frequencies of each integer. The expected length $E(|L(B)|)$ for the L-partition of input B can be expressed as

$$\begin{aligned} E(|L_B|) &= E\left(\sum_{i=1}^{2^d} f_i \cdot X_i\right) = \sum_{i=1}^{2^d} f_i \cdot E(X_i) && \text{by linearity of expectations} \\ &\leq \sum_{i=1}^{2^d} f_i \cdot \left(2q - \frac{q}{2^{d-1}}\right) && \text{by substituting } E(X_i) \text{ due to Lemma 2} \\ &\leq \left(2q - \frac{q}{2^{d-1}}\right) \cdot \frac{n}{d} && \text{since } \sum_{i=1}^{2^d} f_i = \frac{n}{d} \text{ by definition} \\ &\leq n \cdot \left(\frac{2q}{d} - \frac{q}{d \cdot 2^{d-1}}\right) \end{aligned}$$

Since $|L_B| + |R_B| = n$, $E(|R_B|) = n - E(|L_B|) > n \cdot \left(1 - \frac{2q}{d} + \frac{q}{d \cdot 2^{d-1}}\right)$. \square

Corollary 1. *The lengths of the L_B and R_B partitions of an input n -bits long B approach respectively to $2nq/d$ and $n(d - 2q)/d$ bits as the term $q/d \cdot 2^{d-1}$ in Theorem 1 becomes insignificant especially on larger values of d , e.g., assuming $d = 8$ and $q = 1$, $q/d \cdot 2^{d-1} = 1/(8 \cdot 2^7) = 1/2^{10} = 0.00097$, which makes $E(|L(B)|) = 0.24903n$ and $E(|R(B)|) = 0.75097n$ bits that can be assumed to be $0.25n$ and $0.75n$.*

3 The Random Access and Search Operations

In this section we investigate the random access and search mechanisms of the split-coding scheme.

Lemma 3 (Random access support). *Let $Q = \binom{d}{0} + \binom{d}{1} + \dots + \binom{d}{q-1}$ denote the number of distinct d -bit integers that has less than q set bits. Given the L_B and R_B partitions of an input sequence B , if the number of **unique** d -bit blocks on B is less than $2^d - Q$, then the random access to any block D_i of $B = D_1 D_2 \dots D_m$ can be achieved in $O(\log |L_B|)$ -time by using an extra space of $o(|L_B|)$ bits. Otherwise, it takes $O(q \log(n/d))$ -time at the expense of $O((n/d) \cdot (Q/2^d) \cdot \log n)$ additional bits space.*

Proof. Once the L_B and R_B partitions are given, the retrieval of any d -bit block D'_i on $B' = D'_1 D'_2 \dots D'_m$ can be done by using the observation that $D'_1 D'_2 \dots D'_{i-1}$ is $w = d \cdot (i - 1)$ bits long, where ℓ of those bits have been deposited on L_B and the rest $(w - \ell)$ on R_B . Therefore, given ℓ , we know the starting bit position of the left and right partitions of queried D'_i on L_B and R_B , respectively.

We start reading the bits on L_B from the first position and keep scanning until we read either d bits or encounter q set bits, whichever happens first. Once, one of these cases occurs, then the bits are actually the left partition of the first block of the input. Re-initiating the same procedure starting from the next bit position will retrieve the left partition of the second block, and so on. After repeating this operation $(i - 1)$ times, we reach the first bit of the left partition L_i of the queried D'_i . If we have scanned ℓ bits until we reach the beginning of L_i , then it is immediate that the R_i starts just after $(w - \ell)$ bits on R_B . Therefore, once we retrieve L_i from L_B , we read R_i by depositing the remaining $(d - |L_i|)$ bits from the starting bit position $(w - \ell + 1)$ of the R_B . Combining the extracted partitions will reconstruct $D'_i = L_i R_i$ and deploying the inverse permutation would return the original block $D_i = \pi^{-1}(D'_i)$.

The time complexity of that random access procedure depends on finding ℓ , since once ℓ is known, rest is simple constant time operation. Let's first assume each block in $D'_1 D'_2 \dots D'_{i-1}$ includes at least q set bits. Then, the left partitions of these blocks have exactly q set bits and each end with a set bit on L_B , e.g., the initial bits of L_B up until and including the q th set bit are the left partition of D'_1 . Therefore, $\ell + 1$ is the position of the first bit after the $q \cdot (i - 1)$ th set bit on L_B . Finding the position of the k th set bit on a static bitmap of n bits can

be achieved in constant or logarithmic time by constructing a dictionary structure (the *rank/select*(R/S) dictionary) that occupies $o(n)$ bits space [15, 20, 17]. Therefore, the position of the $q \cdot (i - 1)$ th set bit on L_B can be returned by creating that R/S data structure in expense of $o(|L_B|)$ additional space. In this case, detecting the ℓ value can be achieved in $O(\log |L_B|)$ in the worst case, and thus, the random access. On the other hand, possible occurrences of the blocks with less than q set bits, which are directly copied to the L_B , corrupts the above procedure. That requires maintaining the indices of such blocks separately. Assume S_z is the sorted list of block ids that have z set bits for all $z < q$. On list S_z , let the number of items less than the queried i is represented by c_z . Therefore, the total number of set bits of $D'_1 D'_2 \dots D'_{i-1}$ on L_B is $\ell = q \cdot (i - \sum_{z=0}^{q-1} c_z) + \sum_{z=0}^{q-1} (z \cdot c_z)$. As a last step to find the initial position of D'_i , it is necessary to check whether the immediately preceding blocks are in S_0 , and the correct position is computed accordingly. Algorithm 1 lists the complete pseudo-code to retrieve D_i .

Algorithm 1: RandomAccess($i, L_B, R_B, d, \pi, q, S_0, S_1, \dots, S_{q-1}$)

input : i is the queried block index, L_B and R_B are the partitions, d, q, π are the parameters of the coding, S_z is the list of block indices that has $z < q$ set bits.

output: The original d -bit block D_i .

```

1  $C \leftarrow 0$ 
2  $r \leftarrow 1$ 
3 for  $z = 0$  to  $(q - 1)$  do
4   |  $t \leftarrow$  number of indices less than  $i$  in  $S_z$ 
5   |  $C \leftarrow C + z \cdot t$ 
6   |  $r \leftarrow r + t$ 
7 end
8  $\ell \leftarrow q \cdot (i - r) + C$ 
9  $p \leftarrow 1 +$  The position of the  $\ell$ th set bit on  $L_B$ 
10  $r \leftarrow i - 1$ 
11 while  $(r \in S_0)$  do  $p \leftarrow p + d, r \leftarrow r - 1$ 
12  $detected \leftarrow 0, step \leftarrow 0$ 
13 while  $((detected < q) \wedge (step < d))$  do
14   | if  $L_B[p + step] == 1$  then  $detected++$ 
15   |  $step++$ 
16 end
17  $L \leftarrow L_B[p..p + step - 1]$ 
18  $w \leftarrow d \cdot (i - 1) - p + 1$ 
19  $R \leftarrow R_B[w..w + (d - step) - 1]$ 
20  $D'_i \leftarrow L \cap R$ 
21 return  $\pi^{-1}(D'_i)$ 

```

The space overhead for the random access support by Algorithm 1 is the sum of the space occupied by the S_z lists and the additional space used by rank/select

(R/S) data structure for efficient execution of line 9. The number of unique d -bit sequences that have less than q set bit is $Q = \binom{d}{0} + \binom{d}{1} + \dots + \binom{d}{q-1}$. Assuming all distinct 2^d blocks are equally likely among observed (n/d) blocks of the input, the size of the S_z list is expected to be $\approx (n/d) \cdot (Q/2^d)$. Therefore, the space required to store the S_z lists becomes around $\approx (n/d) \cdot (Q/2^d) \cdot \log n$ bits, where the space overhead of the (R/S) dictionaries can be kept around $o(|L_B| \approx 2qn/d)$ with the state-of-the-art implementations. The time complexity of Algorithm 1 depends on detecting the blocks with less than q set bits via the `for` loop in between lines 3 and 7 and also the `select` query run on line 9. The `for` loop takes roughly $O(q \log(n/d))$ -time where the select query can be executed in $O(1)$ or $O(\log(2qn/d))$ time according to the R/S data structure used. Thus, the total time complexity is upper bounded by $O(q \log(n/d))$. \square

It is important to note that if the number of distinct d -bit symbols on the input is less than $2^d - Q$, then the random permutation over integers $\{Q, Q + 1, \dots, 2^d - 1\}$ can be used for the observed blocks, which guarantees to have at least q set bits per block on L_B , and thus, removes the necessity of maintaining S_z lists. In such a case, the time and space complexity of the random access is determined solely by the integrated R/S dictionary data structure, where the state of the art solutions provide $O(\log |L_B|)$ -time with $o(|L_B|)$ extra space.

Lemma 4 (Search on split encoded sequences). *When the number of distinct d bit blocks on input is less than $2^d - Q$, it is possible to efficiently detect all occurrences of a queried pattern P on a split encoded sequence B by maintaining a rank-support on left partition L_B .*

Proof. Let $P = p_1 p_2 \dots p_m$ be a binary sequence that we would like to find its all occurrences on the split encoded B . We apply split encoding to P with the same parameters as $\text{SplitEncode}(P, d, \pi, q) \rightarrow (L_P, R_P)$. The basic idea of the search is to find the matching positions of L_P on L_B , and then verifying each match by comparing R_P with the corresponding position of the candidate on R_B . The first step is to search L_P on L_B . Notice that not all matching positions are valid left partitions as a match can appear bridging two consecutive blocks. Hence, while scanning for L_P , we should accept only completely matching left partitions on L_B . Since we assume the number of distinct d bit blocks on input is less than $2^d - Q$, the permutation π can be constructed over integers $\{Q, Q + 1, \dots, 2^d - 1\}$, which guarantees that each block has exactly q set bits on L_B . Therefore, the validity of a matching position can be confirmed by checking whether *i*) the bit preceding the match position on L_B is a set bit, and *ii*) total number of previous set bits is a multiple of q . For each such valid matching left partition on L_B , we can find the corresponding positions on R_B to extract its right partition via Algorithm 1. We then extract its corresponding right partition from R_B and verify whether the constructed pattern matches the query. The performance of the search process greatly depends the efficiency of the filtering phase, where short P sequences would result in larger set of verification points, and longer ones less. \square

4 The Privacy Aspects

Each d -bit block of input B is represented on R_B with variable number of bits whose lengths are ranging from 0 to $d - 1$. We observe that these variable-length blocks are not prefix-free, and thus, the code-word boundaries on R_B are uncertain, which can only be determined by using the L_B partition, assuming the parameters d and q are known. For instance, in Figure 1, the code-word boundaries on $R_B = 1110010010$ can be detected by using the $L_B = 01101100011001$. One can simply use the random access function Algorithm 1 to detect the bits of a queried i th block on R_B . It is important to address that even after the construction the d -bit block, the permutation function π should be known to reach the original d -bit on B by finding the inverse permutation mapping of the block. Therefore, in absence of the L_B partition and the permutation π , it seems hard to reconstruct the original data B even we assume all other parameters d, q, n are known in advance. We investigate this hardness by analyzing the information released by each partition about the other partition by counting the number of possible bit sequences of length n that produce the same left or right partition. Ideally, given the left partition, all bit sequences of length $n - |L_B|$ should be a valid right partition and decode a distinct B sequence. Therefore, the correct right partition is indistinguishable in the search space of size $2^{n-|L_B|}$. Lemma 5 proves that this is indeed the case.

Lemma 5 (Number of distinct sequences with the same L_B). *Given the left partition L_B of an input n bits long sequence B , there are exactly $2^{n-|L_B|}$ different bit sequences of length n that generates the same L_B .*

Proof. Let's assume we have chosen a random binary string B of length n and generated the left partition L_B and right-partition R_B . Since there is no restriction on the possible bits of the right partition, any binary string X of length $n - |L_B|$ will be valid and produce a different n -bit sequence by the decoding process. In other words, the L partition of all these $2^{n-|L_B|}$ sequences will be equal, where the correct right partition is only one of them. \square

However, we observe interestingly that the same does not hold for the right partition by Lemma 6, which means the knowledge of right partition narrows the possibilities of left partition. Although this exhibits an information leakage, it still provides a reasonable privacy with large enough search space.

Lemma 6 (Number of distinct sequences with the same R_B). *Given the right partition R_B of an input n bits long sequence B , there are less than $2^{n-|R_B|}$ different bit sequences of length n that generates the same R_B . It is also expected that the number of bit sequences sharing the same R_B is also more than $2^{2qn/d}/(2qn/d + 1)$.*

Proof. There are n/d blocks and the length of the L_B partition is $(n - |R_B|)$ bits. Each block has at least zero and at most q set bits on L_B , which means the left partition should include certain number of set bits. Therefore, the binary

strings of length $(n - |R_B|)$ that include less number of set bits than dictated by the original input are not valid and can not be decoded properly, which means the number of possible n bit input strings generating the known R_B is less than $2^{n-|R_B|}$.

If we assume all blocks are represented by q set bits on L_B (in other words, the number of observed distinct blocks on the input is less than $2^d - Q$ and we used only blocks with at least q set bits in the permutation step), then the number of $|L_B|$ bits long binary strings with qn/d set bits are all valid left partitions for the given right partition and each of those strings spell a different B . According to Theorem 1, the length of L_B is expected to be $\approx 2qm$ bits with $m = n/d$. The number of sequences of length $2qm$ bits with qm set bits is greater than $2^{2qm/d}/(2qm/d + 1)$ due to the simple bounds of the central binomial coefficient $4^x = (1 + 1)^{2x} = \sum_{k=0}^{2x} \binom{2x}{k}$. \square

Due to Lemma 5, given the left partition L_B of a split encoded bit sequence, there is no information released about the right partition R_B , since any bit sequence of length $n - |L_B|$ is a valid R_B partition and correctly decodes with the given L_B . On the other hand, due to Lemma 6, once the right partition R_B is given, not all bit sequences of length $n - |R_B|$ are valid left partitions, and thus, the knowledge of R_B narrows down the search space of L_B . However, this reduction is limited and on large input sequences, still it can provide some privacy. It is noteworthy that the above analyses assume the permutation π is known, where actually keeping that π secret improves security significantly.

5 Splitting Beyond Two Partitions

We observe that recursive application of the splitting mechanism on each partition will generate further splits. As long as the depth of this partitioning increases the partitions gets smaller and it gets harder to reconstruct the data in absence of the one piece. Such a cascaded splitting operation might make sense as in distributed storage systems or recent block-chain style storage like IPFS [3]. We generalize our observation and provide a mechanism such that an input bit stream is divided into k partitions, where any ratio on the size of the partitions can be supported, e.g., creating 5 partitions that will occupy approximately $4/16, 5/16, 1/16, 3/16, 3/16$ of the input size, respectively.

Lemma 7 (Splitting into multiple pieces). *Let $\{r_1, r_2, \dots, r_k\}$ represent the given ratios such that $\sum_{i=1}^k r_i = 1$ and $k > 2$. An input bit string of length n can be split into k pieces p_1, p_2, \dots, p_k by recursive application of the split coding, where $|p_i| \approx r_i \cdot n$.*

Proof. Assume a *full* binary tree of $s = 2k - 1$ elements, which is shown by an array $A[1..s]$. In this tree, the leaf nodes represent the final partitions we aim to generate. We place the ratios to the last k elements of the array, i.e., $A[s - i] = r_{k-i}$ for $i = 0..(k - 1)$. We compute the values of the vacant positions $A[t]$ in the array starting from $t = s - k$ down to the root $t = 1$ by

summing up their corresponding children, where the indices of the children of $A[t]$ are $A[2t]$ and $A[2t + 1]$. Notice that the already filled positions are the leaf nodes and thus, the internal nodes can be computed from them. As an example, for $(r_1, r_2, r_3, r_4, r_5) = \langle 4/16, 5/16, 1/16, 3/16, 3/16 \rangle$, the array will end up with $A[1..9] = [1, 10/16, 6/16, 6/16, 4/16, 5/16, 1/16, 3/16, 3/16]$.

Considering that the root node $A[1]$ corresponds to the input bit string, we apply the split coding at each internal node with the most appropriate q and d parameters to support the ratios mentioned in its children. For instance, in the example case, since $A[2] = 10/16$ and $A[3] = 6/16$, choosing $d = 16$ and $q = 5$ provides us the left partition to be close to $n \cdot 10/16$ bits as desired and right partition automatically scales to the remaining $n \cdot 6/16$. We keep applying the split coding to the internal nodes until we produce all the leaves. In our sample case, for instance, we need to split the $A[2]$ according to its children's ratios as $A[4] = 6/16$ and $A[5] = 4/16$, which dictates the left partition $A[4]$ is desired to be around $6/10$ of its parent $A[2] = 10/16$. This can be achieved by setting $d = 10$ and $q = 3$. The bits deposited at the leaf nodes finalize the splitting construction, where reconstruction of the original input requires traversing the tree from the leaves to the root and applying split decode procedure for each internal node. \square

6 Conclusions and Further Studies

We have presented a randomized data splitting algorithm that provides well-control on the sizes of the partitions with efficient search and access mechanisms. Reconstruction of the original input is hard without having access to both partitions and the permutation used in the construction.

Such a privacy-preserving data splitting scheme may find application areas in secure and searchable massive data storage on the cloud. For example, the user may keep the permutation secret, and then store the left and right partitions on different remote servers. Due to the privacy aspects, the admins of those remote locations will not be able to extract the content, while the user will still be able to achieve search and retrieval with the mechanisms described in this study. It is also possible to maintain the left partition on-premise and save right partition encrypted on a cloud to make it completely secure. The user can still perform search operations on the encrypted cloud storage by filtering the candidates from the data on-premise, and then verify the candidates by fetching and decrypting them from the cloud. Obviously, such applications require further analysis of some different possible architectures. Similar applications might be considered for other distributed storage schemes including block-chain storage systems[12]. The recursive application of the proposed split coding may be used to privacy-preserving partitioning of the data into some desired proportions. This can also serve as an alternative in secret sharing schemes [2].

References

1. Bard, G.V., Ault, S.V., Courtois, N.T.: Statistics of random permutations and the cryptanalysis of periodic block ciphers. *Cryptologia* **36**(3), 240–262 (2012)

2. Beimel, A.: Secret-sharing schemes: A survey. In: Coding and Cryptology: Third International Workshop, IWCC 2011, Qingdao, China, May 30-June 3, 2011. Proceedings 3. pp. 11–46. Springer (2011)
3. Benet, J.: Ipfs-content addressed, versioned, p2p file system. arXiv preprint arXiv:1407.3561 (2014)
4. Camtepe, S., Duda, J., Mahboubi, A., Morawiecki, P., Nepal, S., Pawłowski, M., Pieprzyk, J.: Compcrypt–lightweight ans-based compression and encryption. *IEEE Transactions on Information Forensics and Security* **16**, 3859–3873 (2021)
5. Du, W., Atallah, M.J.: Secure multi-party computation problems and their applications: a review and open problems. In: Proceedings of the 2001 workshop on New security paradigms. pp. 13–22 (2001)
6. Duda, J.: Asymmetric numeral systems. arXiv preprint arXiv:0902.0271 (2009)
7. Duda, J., Niemiec, M.: Lightweight compression with encryption based on asymmetric numeral systems. arXiv preprint arXiv:1612.04662 (2016)
8. Durstenfeld, R.: Algorithm 235: Random permutation. *Commun. ACM* **7**(7), 420 (jul 1964). <https://doi.org/10.1145/364520.364540>
9. Fraenkel, A.S., Klein, S.T.: Complexity aspects of guessing prefix codes. *Algorithmica* **12**, 409–419 (1994)
10. Gillman, D.W., Mohtashemi, M., Rivest, R.L.: On breaking a huffman code. *IEEE Transactions on Information theory* **42**(3), 972–976 (1996)
11. Kaaniche, N., Laurent, M.: Data security and privacy preservation in cloud storage environments based on cryptographic mechanisms. *Computer Communications* **111**, 120–141 (2017)
12. Li, R., Song, T., Mei, B., Li, H., Cheng, X., Sun, L.: Blockchain for large-scale internet of things data storage and protection. *IEEE Transactions on Services Computing* **12**(5), 762–771 (2018)
13. Li, Y., Gai, K., Qiu, L., Qiu, M., Zhao, H.: Intelligent cryptography approach for secure distributed big data storage in cloud computing. *Information Sciences* **387**, 103–115 (2017). <https://doi.org/https://doi.org/10.1016/j.ins.2016.09.005>, <https://www.sciencedirect.com/science/article/pii/S0020025516307319>
14. Martins, P., Sousa, L., Mariano, A.: A survey on fully homomorphic encryption: An engineering perspective. *ACM Computing Surveys (CSUR)* **50**(6), 1–33 (2017)
15. Okanojima, D., Sadakane, K.: Practical entropy-compressed rank/select dictionary. In: 2007 Proceedings of the Ninth Workshop on Algorithm Engineering and Experiments (ALENEX). pp. 60–70. SIAM (2007)
16. Plackett, R.L.: The analysis of permutations. *Journal of the Royal Statistical Society Series C: Applied Statistics* **24**(2), 193–202 (1975)
17. Raman, R., Raman, V., Satti, S.R.: Succinct indexable dictionaries with applications to encoding k-ary trees, prefix sums and multisets. *ACM Transactions on Algorithms (TALG)* **3**(4), 43–es (2007)
18. Rubin, F.: Cryptographic aspects of data compression codes. *Cryptologia* **3**(4), 202–205 (1979)
19. Sharma, P., Jindal, R., Borah, M.D.: Blockchain technology for cloud storage: A systematic literature review. *ACM Computing Surveys (CSUR)* **53**(4), 1–32 (2020)
20. Vigna, S.: Broadword implementation of rank/select queries. In: Experimental Algorithms: 7th International Workshop, WEA 2008 Provincetown, MA, USA, May 30-June 1, 2008 Proceedings 7. pp. 154–168. Springer (2008)