

UNIVERSITÀ DI PISA
DIPARTIMENTO DI INFORMATICA

TECHNICAL REPORT: TR-08-18

Masking Patterns in Sequences: a New Class of Motif Discovery with Don't Cares

Giovanni Battaglia Davide Cangelosi Roberto Grossi Nadia Pisanti

July 4th, 2008

ADDRESS: Largo B. Pontecorvo 3, 56127 Pisa, Italy. TEL: +39 050 2212700 FAX: +39 050 2212726

Masking Patterns in Sequences: a New Class of Motif Discovery with Don't Cares

Giovanni Battaglia Davide Cangelosi Roberto Grossi Nadia Pisanti

July 4th, 2008

Abstract

In this paper, we introduce a new notion of motifs, called *masks*, that succinctly represent the repeated patterns for an input sequence T of n symbols drawn from an alphabet Σ . We show how to build the set of all maximal masks of length L and quorum q , in $O(2^L n)$ time and space in the worst case. We analytically show that our algorithms perform better than constant-time enumerating and checking all the potential $(|\Sigma| + 1)^L$ candidate patterns in T after a polynomial-time preprocessing of T . Our algorithms are also cache-friendly, attaining $O(2^L \text{sort}(n))$ block transfers, where $\text{sort}(n)$ is the cache oblivious complexity of sorting n items.

1 Introduction

We consider a new class of *motifs with don't cares*, also motivated by sequence analysis in biological data and data mining on sequences. Motifs are repeated patterns, where a pattern is an intermixed sequence of alphabet symbols (*solid* symbols) and special symbols \circ (*don't care* symbols). The don't care symbol found in a position of the pattern specifies that the position may contain any alphabet symbol (so, it does not care which one). For example, $A\circ T\circ\circ C$ is a motif that repeats twice in the sequence $T = AAAATTACCCCATAGT$ at positions 2 and 3 (starting from 0).

Informally, motifs are those patterns that repeat at least a certain number q of times, where q is a user defined positive integer called the *quorum*. Given an input text T of length n , a quorum $q \geq 2$, and a motif length L , *motif discovery* is the problem of finding the motifs of length L in the sequence T . Each motif may have associated the list of the starting positions of its occurrences in the given sequence T . Unfortunately, due to the don't cares, the number of motifs can be exponentially larger for increasing values of L . Potentially, there can be as many as $\Theta((|\Sigma| + 1)^L)$ motifs, where Σ is the alphabet of the distinct symbols in the text T . Even though this number can be smaller for some particular instances, the known algorithms discovering these motifs still require, in the worst case, exponential time and space for increasing values of L . A lot of research has investigated these issues to mitigate the combinatorial explosion of motifs [8, 9, 27, 28, 29, 32].

Problem formulation. In this paper, we follow a new approach based on modeling motifs by using simple binary patterns, called *masks*, that implicitly represent *families of motifs* in T (instead of individual motifs). For example, mask 101001 represents both $A\circ T\circ\circ C$ and $T\circ G\circ\circ A$: each 1 represents a solid symbol while each 0 represents a don't care symbol. A mask is a *motif* if at least one of its represented patterns occurs q or more times in the given sequence T .

As it should be clear from the above informal definition, we can describe interesting repetitions in a sequence, using a description (mask) that is more succinct than before. Therefore, we aim at giving rise to a smaller set of output motifs. Intuitively, consider some patterns that occur at least q times each and that also share the *same structure*, meant as a certain concatenation of solid and don't care symbols. Since they originate from the same mask, we take this mask as a motif. Moreover, any two patterns sharing the same structure but having a different number of occurrences in T (still at least q in number), which were previously considered as different motifs, are now giving rise to the same motif by our definition of mask. Since each mask can be seen as a binary string, we have potentially 2^L masks to examine instead of $(|\Sigma| + 1)^L$ classical motifs.

In summary, our new class of motifs may summarize some regularities in the given sequence T , better than ever before. We study the problem of detecting *maximal masks*, namely, the most specific ones (maximal number of 1s) such that at least one of its represented patterns occurs q times or more (see Section 2 for a formal definition).

For example, given the text $T = \text{AAAATTACCCCATAGT}$, fixing $L = 4$ and $q = 2$, we obtain the maximal masks 1110, 0111, and 1101. Notice that 1110 and 0111 are equivalent since they originate the same patterns (three consecutive solid symbols) ignoring border effects, so we can treat them as the same mask. Therefore the patterns that are represented by the maximal masks are AAA, CCC, and AA◦T, and so the parameter L can equivalently be read as an upper bound on their length. Interestingly, the text positions of the occurrences of the patterns implicitly represented by a mask do not overlap (while the patterns can overlap), so we can always associate a partition of the text positions with the mask.

It is worth noting that motif discovery has many applications in the investigation of properties of biological sequences. In such applications, it is a must to allow distinct occurrences of a motif to show some differences. In other words, we actually infer *approximated* motifs. Such approximation can be realized in several ways, according to the kind of application one has in mind. Motifs of limited length with don't cares can typically model biological object such as transcription factors binding sites, that are characterized by a short length, and a high conservation of their structure. Also, they present a high conservation of the contents in certain positions while for others it does not matter at all. The don't care symbols of our masks indeed aim at *masking* the latter, while the solid character should *unmask* the former.

Moreover, our masks could also be employed as building blocks for longer and flexible motifs, of different kind, allowing also indels. In recent years, there has been a growing interest in *seeds* for several applications (preprocessing filtration prior to a multiple alignment, approximate search task, data base search, BLAST like homology search, profile search, probe design) in bioinformatics ([10, 17, 23, 24, 25, 26, 31]). Among them, many have focused the attention on *gapped seeds*, or *spaced seeds* ([2, 3, 4, 6, 20, 30]). It turns out that gapped seeds can be found using the masks, and thus the algorithms introduced in this paper.

Finally, a new application of finding motifs with don't cares could help to detect structural similarities, with a suitable input sequence. For example, when investigating the folding of a DNA sequence, it can be interesting to rewrite the sequence itself into the alphabet $\{\mathbf{w}, \mathbf{s}\}$ replacing each A and T with \mathbf{w} (weak), and each C and G with \mathbf{s} (strong). The motivation is that in the *base pairing* that assists in stabilizing the DNA structures, adenine (A) binds to thymine (T) via two hydrogen bonds, while cytosine (C) forms three hydrogen bonds with guanine (G). Hence, the latter bond is stronger than the former, and this has an influence on the actual structure of the molecule. In this framework, a motif on such sequence could represent a repeated structure, regardless of the actual DNA bases that form it.

Our results. In this paper, we introduce a new notion of motifs, the *masks*, and investigate how to build the set \mathcal{M} of all *maximal masks* of length L and quorum q , for an input sequence T of n symbols drawn from an alphabet Σ .

We first show that the partition of the text positions induced by an arbitrary mask can be conceptually represented as a pruned trie of the represented patterns, storing the classes that constitute the partition in the leaves of the trie. We then extend the Karp-Miller-Rosenberg doubling scheme [18] to the masks having length an increasing sequence of powers of 2 up to L . Interestingly, the resulting scheme avoids to actually create the tries and just needs scanning and sorting some suitable lists of consecutive pairs and triplets of integers (this is considered to be cache friendly).

However, the above method generates the set \mathcal{Q} of all the masks having quorum q for the given sequence T , where $\mathcal{Q} \supseteq \mathcal{M}$, including those that are not maximal. Maximality checking in \mathcal{Q} to select the masks in \mathcal{M} can be expensive: precisely, it may take $\Theta(|\mathcal{Q}|^2 L)$ time in the worst case (e.g. [13]), yielding a cost of $\Omega(2^{2L} L)$ time. We therefore introduce the crucial notion of *safe masks*, which includes the maximal masks as a special case. We show how to traverse the lattice of 2^L masks of length L touching only safe masks, so that maximal masks can be efficiently detected.

Consequently, the final bound for discovering all the masks belonging to \mathcal{M} is $O(2^L n)$ time and space in the worst case, which is our main result. To evaluate the efficiency of the proposed algorithms, consider the following scenario in which, after preprocessing the text T in polynomial time, we ideally check, in constant time per pattern, if any of the $(|\Sigma| + 1)^L$ candidate patterns has quorum and is maximal. We compare the cost of $O(n^{O(1)} + (|\Sigma| + 1)^L)$ thus obtained against the cost of $O(2^L n)$ that we obtain in this paper.

Since $2^L n = O((|\Sigma| + 1)^L)$ when $L = \Omega(\log_{|\Sigma|} n)$, and $2^L n = O(n^{O(1+1/\log_2 |\Sigma|)})$ otherwise, we can bound our $O(2^L n)$ cost by $O(n^{O(1+1/\log_2 |\Sigma|)} + \min\{2^L n, (|\Sigma| + 1)^L\})$. As a result, the latter is always better than the ideal bound of $O(n^{O(1)} + (|\Sigma| + 1)^L)$ stated above. In other terms, our algorithms perform better than virtually *constant-time enumerating and checking* all the potential $(|\Sigma| + 1)^L$ candidate patterns in T .

At this point, we need to make an observation on the model of computation adopted in this paper, which is the standard RAM with word size of w bits. Since we need to address potentially $\Theta(2^L)$ masks, we assume in the rest of the paper that $w = O(L)$. If this is not the case, the total cost must be multiplied by a factor of $O(L/w)$.

Finally, given the scan-and-sort nature of our algorithms, we naturally obtain a cache-oblivious solution to our problem as a byproduct. To our knowledge, this is the first cache-oblivious solution for a motif discovery problem, which is useful for long input sequence(s). Indeed, our algorithms work also in the *cache-oblivious model*, introduced by Frigo et al. [12], to deal with such a situation, where M is the size of the fast memory, and B is the size of the block in each transfer between fast and slow memories. The goal is to minimize the number of block transfers without letting the algorithms know the values of M and B (which are introduced for the purpose of the analysis only). For example, scanning n consecutive elements has a complexity of $\Theta(n/B)$ block transfers while the optimal complexity of sorting is $\text{sort}(n) = \Theta\left(\frac{n}{B} \log_{M/B} \frac{n}{B}\right)$ block transfers [?, ?, ?, 12].

Employing the simple scan and the cache-oblivious sorting in our algorithms, we do not need to further orchestrate their memory accesses. Using this model, we can indeed obtain a complexity of $O(2^L \text{sort}(n))$ block transfers for finding \mathcal{Q} and \mathcal{M} . Along the same lines, our algorithms can run in a distributed setting, such as a cluster of computers, using distributed sorting.

Related problems and state of the art. We are not aware of any previous work introducing our class of motifs. Hence, we relate our results in motif discovery to those of mining frequent itemsets, where more sophisticated techniques have been found over the years. The notion of masks comes naturally into play when performing data mining for frequent itemsets, where the “apriori” algorithm is intensively employed [16]. Here, a set of L items is given, and each transactions (basket) corresponds to a subset of these items, which can be represented as a binary sequence in which the i th symbol is 1 if and only if the i th item is in the basket. A set of baskets can be therefore represented as a set of masks in our terminology. For the lattice of all possible 2^L masks, all possible itemsets should be examined. Note that, instead, our definition of masks has the goal of condensate patterns that have the same sequence of solid and don’t care symbols. Moreover, our traversal of the lattice is different from the apriori, since we start from the top and generate candidates in a different way, namely, using safe masks.

As far as we know, the “dualize and advance” algorithm [14, 15] is the best theoretical approach that can be obtained in terms of running time. It sets up an interesting connection between mining itemsets in the lattice of 2^L masks and finding hypergraph traversals. In our terminology, suppose to have incrementally found some of the maximal masks, say $\mu_1, \mu_2, \dots, \mu_k$. We can build an hypergraph in which the nodes are the L available items numbered from 1 to L , and the hyperedges are the complement of $\mu_1, \mu_2, \dots, \mu_k$ (in other words, the 0s indicate the chosen items). Then, in order to find additional maximal masks (and hence add hyperedges), it suffices to find all the hypergraph traversals as starting points for upward paths in the lattice, where each traversal is a minimal hitting set for the current set of k hyperedges [1].

The problem of finding hypergraph traversals is intimately related to the dualization of monotone Boolean functions [7]. The known algorithms require $O(2^L)$ time in the worst case [1, 19] until the seminal result in [11, 21] showing a subexponential bound proportional to $t(k) = k^{O(\log k)}$ time, when the number of hyperedges k is $o(2^L)$. This algorithm is plugged into the scheme of the “dualize and advance” algorithm, giving a bound of $O(n \times t(|\mathcal{M}| + |Bd^-(\mathcal{M})|))$ as shown in [14], where we include the cost $O(n)$ of verifying the quorum, and $Bd^-(\mathcal{M})$ is a set of non-maximal masks that are “close” inside the lattice to the ones in \mathcal{M} . While $|\mathcal{M}|$ can be subexponential, there are cases in which $|\mathcal{M}| + |Bd^-(\mathcal{M})| = \Theta(2^L)$ (see [14]) and so the final bound can be $\Omega(2^{L^2} n)$.

Surprisingly, this and other approaches based on hypergraph traversals, which are the state of the art theoretically, are slower than our solution in the worst case. We have some preliminary experiments that our solution is also faster in practice. We counted the number of times a function checking the quorum for a given mask is invoked by the algorithms, which is a clear measure of their running time. Here, dualize and advance needs to query many more masks than our safe masks, thus suggesting that the notion of safeness

of masks is crucial to our algorithms. We plan to run an extensive experimental analysis as future work.

Paper organization. This paper is organized as follows. In Section 2, we give a formal definition of our motifs based on masks. We then present efficient algorithms to compute these motifs, reviewing also basic notions as that of maximality in the light of this new class of motifs. In Section 3, we show how to discover the motifs of length L in $O(2^L n)$ time by extending the Karp-Miller-Rosenberg approach to the masks. In Section 4, we represent the space of all possible 2^L masks of length L as a lattice on which it is defined an oracle function for the quorum and introduce the notion of safe masks; we describe how to traverse implicitly this lattice for discovering maximal masks in it, querying the oracle only for safe masks. Finally, we draw our conclusions in Section 5.

2 A New Class of Motifs

In this section, we introduce our new class of motifs with don't cares. Starting from some basic notions, we describe some features of this class that will then be exploited by the algorithms in the rest of the paper. Let T be an input text of size n drawn over the alphabet $\Sigma = \{\sigma_1, \dots, \sigma_p\}$. The sequence T can be seen as an array $T[0 \dots n-1]$ of symbols, where symbol $T[i] \in \Sigma$ is stored into position i , for $0 \leq i \leq n-1$. A substring $T[i]T[i+1] \dots T[j]$ is represented as $T[i \dots j]$.

2.1 Masks and patterns

Given a positive integer L , we call *mask* any binary sequence in $\{0, 1\}^L$ (hence, L is the length of the mask). For a given mask $\mu = \mu[0 \dots L-1]$, we define $S_\mu = \{i \mid \mu[i] = 1, \text{ for } 0 \leq i \leq L-1\}$ as the set of its *solid positions*, and $D_\mu = \{i \mid \mu[i] = 0, \text{ for } 0 \leq i \leq L-1\}$ as the set of its *don't care positions*. For example, mask $\mu = 1010$ has $S_\mu = \{0, 2\}$ and $D_\mu = \{1, 3\}$.

A *pattern* is a regular expression $m \in (\Sigma \cup \{\circ\})^L$, where \circ is the *don't care symbol* that matches any symbol in Σ . We say that a pattern m is an *instance* of a mask μ when, for each position $0 \leq k \leq L-1$, it is $m[k] = \circ$ if and only if $\mu[k] = 0$. For example, given $\Sigma = \{A, C\}$, the mask $\mu = 1010$ has four instances, namely, $m_1 = A\circ A\circ$, $m_2 = C\circ A\circ$, $m_3 = A\circ C\circ$, and $m_4 = C\circ C\circ$.

By exploiting the fact that a mask μ implicitly represents the patterns that are its instances, we define a new relation among the text substrings according to μ , as follows.

Definition 1 (\equiv_μ relation) *Given a mask μ of length L , a text T of length n , and two text positions $0 \leq i, j \leq n-L+1$, we say that $i \equiv_\mu j$ if and only if $T[i+k] = T[j+k]$ for each solid position $k \in S_\mu$ in the mask.*

Definition 1 relates any two text substrings of length L , appearing at positions i and j , when these substrings match in every solid position specified by the mask μ . For example, given $T[i \dots i+L-1] = \text{ACTACT}$ and $T[j \dots j+L-1] = \text{AGTTCT}$, let us consider the two masks $\mu = 101011$ and $\mu' = 110111$. It holds that $i \equiv_\mu j$ because $A\circ T\circ C\circ T$ occurs both at positions i and j of T , while $i \not\equiv_{\mu'} j$ because the two substrings ACTACT and AGTTCT mismatch at the solid positions in $\{1, 3\} \subseteq S_{\mu'}$.

It is easy to see that the relation \equiv_μ is an equivalence relation. Therefore, for a given mask μ , the relation \equiv_μ induces a *partition* of the first $n-L+1$ positions in T . Namely, for each equivalence class C in the partition, we have that $i, j \in C$ if and only if $i \equiv_\mu j$. Hence, each text position i (for $0 \leq i \leq n-L+1$) belongs to exactly one equivalence class. We denote the partition resulting from μ by π_μ . We use $|\pi_\mu|$ to indicate the number of equivalence classes in it, and $|C|$ to indicate the number of elements in a class $C \in \pi_\mu$.

Given an equivalence class C of a partition π_μ , we can associate a pattern m_C of length L with C . Specifically, symbol $m_C[k] = \circ$ if k is a don't care position of the mask μ ($k \in D_\mu$) while $m_C[k] = T[i+k]$ if k is a solid position ($k \in S_\mu$ and for any arbitrary $i \in C$).

In order to better illustrate the properties of the partition π_μ and the corresponding patterns m_C where $C \in \pi_\mu$, we can use a trie (digital search tree [22]) built on the set of strings $\{m_C \mid C \in \pi_\mu\}$. In this trie, the special symbol \circ can be treated as an ordinary symbol, since all the patterns share the same mask μ . We can arrange the strings in this way since the arcs on the same level of the trie are either all labeled with

Figure 1: Trie for \equiv_μ where $\mu = 11\circ\circ$

Figure 2: Trie for $\equiv_{\mu'}$ where $\mu' = 11\circ 1$.

a solid symbol or all labeled with a don't care symbol. Each root-to-leaf path spells out a distinct pattern m_C , and the corresponding leaf stores the text positions in the class C of the partition π_μ .

Example 1 Consider the alphabet $\Sigma_{DNA} = \{A, T, C, G\}$ and the text $T = AAAATTACCCCATAGT$ of length $n = 16$. For a mask $\mu = 1100$ of length $L = 4$, the partition induced by μ is $\pi_\mu = \{C_0 C_1 \cdots C_6\}$, where $C_0 = \{4\}$, $C_1 = \{5, 12\}$, \dots , $C_6 = \{0, 1, 2\}$ are the classes labeling the leaves of the trie shown in Figure 1. The instances of μ are the patterns $m_{C_0} = TT\circ\circ$, $m_{C_1} = TA\circ\circ$, \dots , $m_{C_6} = AA\circ\circ$.

2.2 Partial order of masks and maximality

We now draw our attention to the masks μ that have the maximum number of solid symbols while inducing a partition π_μ that contains at least one class C such that $|C| \geq q$ for the given *quorum* q . For this, we need to introduce a partial order on the masks.

Definition 2 (\preceq relation) *Given two masks μ and μ' of length L , we say that μ is less specific than μ' (denoted by $\mu \preceq \mu'$) if and only if $\mu[i] \leq \mu'[i]$ for $0 \leq i \leq L - 1$.*

When Definition 2 holds, we also say that μ' is *more specific* than μ , and that μ is a predecessor of μ' , and μ' is a successor of μ . For example, $0001 \preceq 1101$, while $0001 \not\preceq 0010$. The mask μ is an *immediate predecessor* of μ' (denoted $\mu \preceq_1 \mu'$) if $\mu \preceq \mu'$ and they differ in exactly one symbol. For example, $1001 \preceq_1 1101$.

Relation \preceq is a partial order among the masks because it is reflexive, antisymmetric ($1001 \preceq 1101$, but $1101 \not\preceq 1001$) and transitive. Hence, it gives rise to the partially ordered set $\mathcal{L} = \langle \{1, 0\}^L, \preceq \rangle$, which is a finite lattice of 2^L masks. The top mask of \mathcal{L} is $1 \dots 1$ and the bottom mask is $0 \dots 0$. The lattice \mathcal{L} is isomorphic to the power set lattice $\mathcal{P} = \langle \mathcal{P}(\{0, \dots, L-1\}), \subseteq \rangle$ because each mask represents the characteristic vector of a set in the powerset $\mathcal{P}(\{0, \dots, L-1\})$, where $\mu \preceq \mu'$ if and only if $S_\mu \subseteq S_{\mu'}$.

Analogously, we can define the \preceq relation between patterns. Given two patterns m and m' of length L , we say that m is *less specific* than m' (written $m \preceq m'$) if and only if either $m[i] = m'[i]$ or $m[i] = \circ$ for $0 \leq i \leq L - 1$. For example, $\circ\circ\circ A \preceq A T \circ A$ while $\circ\circ\circ A \not\preceq A T \circ C$ and $\circ C \circ A \not\preceq A \circ \circ A$. Note that \preceq is a partial order also for the patterns. However, it gives rise to a lattice of $(|\Sigma| + 1)^L$ patterns, and so in the rest of the paper we adopt the binary lattice \mathcal{L} of 2^L masks.

At this point, we may wonder what is the connection between the partitions induced by the masks (Section 2.1) and the lattice \mathcal{L} formed by the masks. The following result shows interesting properties of two masks regarding inclusive relation between their equivalence classes.

Proposition 1 *For any two masks μ and μ' such that $\mu \preceq \mu'$, the following properties hold:*

- (i) *For each class $C \in \pi_\mu$, there exist classes $C_0, C_1, \dots, C_{s-1} \in \pi_{\mu'}$ that form a partition of C .*
- (ii) *For each class $C' \in \pi_{\mu'}$, there exists a class $C \in \pi_\mu$ such that $C' \subseteq C$.*

Example 2 (continued) Consider mask $\mu' = 1101$, for which $\mu \preceq \mu'$, where $\mu = 1100$. Consider the equivalence classes of π_μ (Figure 1) and $\pi_{\mu'}$ (Figure 2). We have that $C_0 = \{4\} = C'_0$, $C_1 = \{5, 12\} = \{12\} \cup \{5\} = C'_1 \cup C'_2$, $C_2 = \{7, 8, 9\} = \{7\} \cup \{8\} \cup \{9\} = C'_3 \cup C'_4 \cup C'_5$, $C_3 = \{10\} = C'_6$, $C_4 = \{3, 11\} = \{3\} \cup \{11\} = C'_7 \cup C'_8$, $C_5 = \{6\} = C'_9$, and $C_7 = \{0, 1, 2\} = \{1, 2\} \cup \{0\} = C'_{10} \cup C'_{11}$. Therefore, both properties (i) and (ii) hold.

In fact, $\pi_{\mu'}$ is a refinement of π_μ whenever $\mu \preceq \mu'$. We can easily extend the above properties to the patterns, obtaining the following corollary.

Corollary 2 *Let μ and μ' be any two masks such that $\mu \preceq \mu'$. For any class $C \in \pi_\mu$, there exists a class $C' \in \pi_{\mu'}$ such that their corresponding patterns satisfy $m_C \preceq m_{C'}$. (The reverse holds as well, namely, for any class $C' \in \pi_{\mu'}$, there exists a class $C \in \pi_\mu$ such that $m_C \preceq m_{C'}$.)*

2.3 Problem statement

We are given a text T of length n and a length L for the masks, along with the quorum $q \geq 2$. We say that a mask μ has *quorum* if there is an equivalence class C in the partition π_μ such that $|C| \geq q$. The relation \preceq and the quorum q are the key ingredients to find interesting masks.

Let $\mathcal{Q}(L, T, q)$ the set of all masks of length L that have quorum q . A mask μ is *maximal* if it has quorum and no other mask μ' of the same length has quorum and is more specific than μ (i.e. $\mu \preceq \mu'$). We denote by $\mathcal{M}(L, T, q) \subseteq \mathcal{Q}(L, T, q)$ the set of all maximal masks (\mathcal{M} for short), and address the following problem in this paper.

Problem 1 (Maximal Masks Problem) Given a text T , a mask length L , and a quorum q , find all the maximal masks in $\mathcal{M}(L, T, q)$.

Example 3 (continued) Using the same text $T = \text{AAAATTACCCCATAGT}$, for $L = 4$ and $q = 2$, the set of masks with quorum is $\mathcal{Q}(4, T, 2) = \{0000, 0001, 0010, 0011, 0100, 0101, 0110, 0111, 1000, 1001, 1010, 1100, 1101, 1110\}$, while the maximal masks are only those of the set $\mathcal{M}(4, T, 2) = \{0111, 1101, 1110\}$. The patterns that are instances of the maximal masks are $\circ\text{AAA}$, $\text{AA}\circ\text{T}$, and $\circ\text{CCC}$. Notice that masks 0111 and 1110, that are both in \mathcal{M} , actually represent the same patterns, except possibly border effects, because one is the shift of the other obtained only changing the number of \circ at the sides. We will show in Section 3 how we will remove this sort of redundancy, that will actually achieve obtaining also an important reduction of the search space, as a positive side effect.

We can see that checking whether a mask μ has quorum q corresponds to evaluate a Boolean predicate $P_T(\mu)$ that returns true if and only if there exists a pattern m that is an instance of μ and that has at least q matches in T . Note that $P_T(\mu)$ is anti-monotone.¹ Hence, Problem 1 can be equivalently restated as checking the predicate $P_T(\mu)$ on a binary lattice of 2^L masks. In the rest of the paper, we will describe several approaches to solve this problem, and some of them will use this view of Problem 1.

3 The KMR Approach for Masks with Quorum

We now describe our first algorithm to solve Problem 1 and build the sets $\mathcal{Q}(L, T, q)$ and $\mathcal{M}(L, T, q)$ in $O(2^L n)$ time. Conceptually, we want to build the partition π_μ for each mask $\mu \in \{0, 1\}^L$, check whether μ has quorum (i.e. there exists a class $C \in \pi_\mu$ such that $|C| \geq q$), and that no mask μ' with $\mu \preceq \mu'$ has quorum. Note that a straightforward way of checking whether a mask has quorum requires to scan the text T and build a trie like that shown in Figure 2, in $\Theta(Ln)$ time. This gives a total cost of $\Omega(L2^L n)$ time since there are $\Omega(2^L)$ masks to check.

We reduce the above cost to $O(2^L n)$ time using a different and simple approach that avoids to explicitly build the trie for each mask. Our idea is to maintain the partitions induced by the masks of increasing lengths as follows. Assuming without loss of generality that L is a power of 2, we first compute the partitions induced by the masks of length 1; inductively, given the partitions induced by the masks of length 2^i , we show how to compute the partitions induced by the masks of length 2^{i+1} in a way that does not explicitly need the tries, even though we will implicitly refer to them during the description of our algorithm.

We implement our idea using the *KMR approach* proposed by Karp, Miller and Rosenberg [18]. The KMR approach addresses the problem of identifying exact repeated substructures of fixed size in a given combinatorial structure. It applies to finding repeated substrings in strings, repetead subtrees in trees, and repeated segments in arrays. For strings, KMR uses a relation E_k according to which two substrings of length k beginning at positions i and j of the text T are *k-equivalent*, written $i E_k j$, if and only if they are identical in every position. Given this, KMR provides a characterization of $E_{k+k'}$ in terms of E_k and $E_{k'}$, so that it constructs inductively the sets $\{E_2, E_4, E_8, \dots, E_L\}$ by setting $k = k'$. That is, it doubles the length of the substrings by means of a concatenation of two substrings from the previous iteration. KMR starts out from the set E_1 (obtained by a simple scan of the input sequence T), and ends at the required length L . Each iteration takes time $O(n)$, where n is the length of the text T . Since the number of iterations is $O(\log L)$, the overall complexity of KMR is $O(n \log L)$ time.

¹We recall that, given a partial order \preceq , a predicate p is *anti-monotone* if, for any x and y such that $x \preceq y$, we have that $p(y) = \text{true}$ implies $p(x) = \text{true}$. Conversely, p is *monotone* if $p(x) = \text{true}$ implies $p(y) = \text{true}$.

Figure 3: Partitions for masks $\mu = 11$ (left) and $\mu' = 01$ (right).

In our case, Problem 1 differs from that solved by KMR since we use masks as one further level of abstraction. We do not apply KMR directly to the text substrings; instead, we double the length of the masks and, as a side effect, we double the length of the induced equivalent substrings in the text (i.e they match for all the solid positions of the given mask). In this way, the original KMR can be seen as a solution for the special case in which the mask is made by all solid symbols $11 \dots 1$.

We observe that, for any given mask, we can employ KMR when the relation E_k is replaced by the \equiv_μ relation introduced in Definition 1. We also replace the concatenation performed by KMR at each iteration by the concatenation operation among masks. Given two masks μ and μ' , we indicate their concatenation by $\mu\mu'$. The following result relates the \equiv_μ equivalence relation to the mask concatenation operation, showing how the KMR paradigm can be generalized to our case.

Lemma 3 *Given a string T of length n , two masks μ, μ' , two positions i, j in T then $i \equiv_{\mu\mu'} j$ if and only if $i \equiv_\mu j$ and $(i + |\mu|) \equiv_{\mu'} (j + |\mu|)$.*

Proof We start by showing how $i \equiv_{\mu\mu'} j$ implies $i \equiv_\mu j$ and $(i + |\mu|) \equiv_{\mu'} (j + |\mu|)$: If $i \equiv_{\mu\mu'} j$, then by definition we have that

$$T[i + k] = T[j + k] \quad \forall k \in S_{\mu\mu'}. \quad (1)$$

Notice that $S_{\mu\mu'} = S_\mu \cup (S_{\mu'} + |\mu|)$. Hence, (1) implies that (i) $T[i + k] = T[j + k] \quad \forall k \in S_\mu$, and (ii) $T[i + k] = T[j + k] \quad \forall k \in (S_{\mu'} + |\mu|)$. Observe that (i) exactly matches the definition of $i \equiv_\mu j$, which is then proved. On the other hand, (ii) implies that $T[i + k] = T[j + k]$ holds for $k = |\mu| + k'$ for all $k' \in S_{\mu'}$, that is to say $T[i + |\mu| + k'] = T[j + |\mu| + k'] \quad \forall k' \in S_{\mu'}$, and hence that $(i + |\mu|) \equiv_{\mu'} (j + |\mu|)$. In order to show that if $i \equiv_\mu j$ and $(i + |\mu|) \equiv_{\mu'} (j + |\mu|)$, then $i \equiv_{\mu\mu'} j$, it is enough to observe that all steps above can be inverted, and hence the result is proved.

3.1 Partition construction and generation of masks

The implementation of Lemma 3 requires an efficient procedure that computes all the related equivalence classes in a partition. In this section we describe an algorithm that solves this problem.

Given a quorum $q \geq 2$ and two partitions π_μ and $\pi_{\mu'}$, where μ is not necessarily different from μ' , the algorithm returns a new partition $\pi_{\mu\mu'}$ built according to Lemma 3, possibly filtered in order to satisfy the quorum constraint. The key point is illustrated in Figure 3, where we are given two partitions π_μ and $\pi_{\mu'}$ and we want to obtain the new partition $\pi_{\mu\mu'}$ shown in Figure 2. (Note that the text is the same $T = \text{AAAATTACCCCATAGT}$, that $q = 2$, and that the tries are shown for the sake of presentation since we do not actually employ them in our implementation.) In order to concatenate two masks μ and μ' of length ℓ , the main steps can be summarized as follows (we refer to Figure 2 and Figure 3 as an example).

1. We are given masks μ and μ' and their induced partitions π_μ and $\pi_{\mu'}$. We consider only the equivalence classes C such that $|C| \geq q$, and number these classes so that each class has its own *class name* inside its partition. (In our example, $\pi_\mu = [\{5, 12\}_0, \{7, 8, 9\}_1, \{3, 11\}_2, \{0, 1, 2\}_3]$ and $\pi_{\mu'} = [\{3, 4, 11, 14\}_0, \{6, 7, 8, 9\}_1, \{1, 2, 5, 10, 12\}_2]$, where we ignore classes with less than q elements and report the numbering of each relevant class as its subscript.)
2. We create a (multiset) list LP of pairs as follows. First, for each class $C \in \pi_\mu$, we add the pairs $\langle i, n_C \rangle$ to LP for all positions $i \in C$, where n_C is the number assigned to C in step 1. Second, for each class $C' \in \pi_{\mu'}$, we add the pairs $\langle i' - |\mu|, n_{C'} \rangle$ to LP for all positions $i' \in C'$ such that $i' \geq |\mu|$, where $n_{C'}$ is the number assigned to C' in step 1. (We obtain $LP = [\langle 5, 0 \rangle, \langle 12, 0 \rangle, \langle 7, 1 \rangle, \langle 8, 1 \rangle, \langle 9, 1 \rangle, \langle 3, 2 \rangle, \langle 11, 2 \rangle, \langle 0, 3 \rangle, \langle 1, 3 \rangle, \langle 2, 3 \rangle, \langle 1, 0 \rangle, \langle 2, 0 \rangle, \langle 9, 0 \rangle, \langle 12, 0 \rangle, \langle 4, 1 \rangle, \langle 5, 1 \rangle, \langle 6, 1 \rangle, \langle 7, 1 \rangle, \langle 3, 2 \rangle, \langle 8, 2 \rangle, \langle 10, 2 \rangle]$ in our example since $|\mu| = 2$.)
3. We sort the list LP in a *stable* way according to the first component of each pair in it. We drop from the list the pairs $\langle i, j \rangle$ such that no other pair has i as its first component in the list. (We obtain $LP = [\langle 1, 3 \rangle, \langle 1, 0 \rangle, \langle 2, 3 \rangle, \langle 2, 0 \rangle, \langle 3, 2 \rangle, \langle 3, 2 \rangle, \langle 5, 0 \rangle, \langle 5, 1 \rangle, \langle 7, 1 \rangle, \langle 7, 1 \rangle, \langle 8, 1 \rangle, \langle 8, 2 \rangle, \langle 9, 1 \rangle, \langle 9, 0 \rangle, \langle 12, 0 \rangle, \langle 12, 0 \rangle]$.)

4. The actual concatenation between μ and μ' takes place. Indeed, there is an occurrence of a pattern m (instance of mask $\mu\mu'$) in position i if and only if $\langle i, j \rangle$ and $\langle i, j' \rangle$ are consecutive pairs in LP for some $0 \leq j, j' \leq n-1$. Hence, we generate a triplet $\langle j, j', i \rangle$ from these two pairs (note that $\langle i, j \rangle$ precedes $\langle i, j' \rangle$ in LP , thus forming a list LT of triplets. (We obtain $LT = [\langle 3, 0, 1 \rangle, \langle 3, 0, 2 \rangle, \langle 2, 2, 3 \rangle, \langle 0, 1, 5 \rangle, \langle 1, 1, 7 \rangle, \langle 1, 2, 8 \rangle, \langle 1, 0, 9 \rangle, \langle 0, 0, 12 \rangle]$.)
5. We lexicographically sort the list LT according to the first two components of each triplet in it. We drop from the list the triplets $\langle j, j', i \rangle$ such that there are less than q triplets in the list having j and j' as their first component in the list, since they do not reach the quorum. (We obtain $LT = [\langle 3, 0, 1 \rangle, \langle 3, 0, 2 \rangle]$.)
6. We start from an empty partition $\pi_{\mu\mu'}$. For each maximal run of consecutive triplets $\langle i, j, k_1 \rangle, \langle i, j, k_2 \rangle, \dots, \langle i, j, k_r \rangle$ in LT ($r \geq q$), we add the class $\{k_1, k_2, \dots, k_r\}$ to $\pi_{\mu\mu'}$. We return $\pi_{\mu\mu'}$ after completing the scan of LT . (In our example, $\pi_{\mu\mu'} = [\{1, 2\}]$ since only one class contains at least q elements in Figure 2.)

As for the sorting in steps 3 and 5, we actually observe that it suffices a stable pairing, meaning that pairs (triplets) having the same first (two) component(s) should be consecutive in the resulting list. Note that this does not imply that the list should be sorted (e.g. Google's `mapReduce`).

Lemma 4 *Given partitions π_μ and $\pi_{\mu'}$ and an integer $q \geq 2$, steps 1–6 correctly compute the set $\{C \in \pi_{\mu\mu'} \mid |C| \geq q\}$ in $O(n)$ time and space.*

Proof: We start by proving that the method is correct. Given two partitions π_μ and $\pi_{\mu'}$ we start by labeling, at Step 1, each class in π_μ and $\pi_{\mu'}$ with distinct class names. In Step 2 π_μ is rewritten as sequence of pairs $\langle i, n_C \rangle$, while $\pi_{\mu'}$ is rewritten as sequence of pairs $\langle i' - |\mu|, n_{C'} \rangle$. In this way, for each index $i \in C$ such that $C \in \pi_\mu$ and $i + |\mu| \in C'$ such that $C' \in \pi_{\mu'}$, any $j \in C$ such that $i \equiv_\mu j$ and $i + |\mu| \equiv_{\mu'} j + |\mu|$ has now become i in this rewriting. Step 3 drops pairs all $\langle pos, class \rangle$ corresponding to a position pos where there can be no instance of the mask $\mu\mu'$ because either μ occurs there but μ' does not occur in $pos + |\mu|$, or the opposite holds. Before doing so, a stable order is performed on the pairs in order to facilitate the detection of these positions, but also to simplify the execution of Step 4. Indeed, at the beginning of Step 4, pairs are sorted in such a way that two pairs that start with the same positions are consecutive, and, thanks to the rewriting of Step 2, they represent positions where there exists an instance of the mask $\mu\mu'$. Hence, for the remaining pairs, a triplet $\langle j, j', i \rangle$ is created for consecutive pairs $\langle i, j \rangle$ and $\langle i, j' \rangle$ storing index i for the new class that can be built concatenating j and j' . Step 5 drops, among the resulting triplets, those that do not satisfy the quorum. This is done in a simple way thanks to the preliminar lexicographical sorting of the list done according to the n_C and $n_{C'}$ components of each triplet. In this way, indeed, distinct occurrences of the same class are now consecutive elements of the sorted list of triplets. Step 6 then constructs the classes of $\pi_{\mu\mu'}$ storing the set of indexes in the consecutive triplets with the same class names that have remained (because they reach the quorum). Then $\pi_{\mu\mu'}$ is returned. Thanks to Lemma 3, the above steps, correctly compute the partition $\pi_{\mu\mu'}$.

We prove that steps 1–6 take $O(n)$ time and space. Since there can be at most $O(n)$ positions, the list generated at Step 1 has size in $O(n)$. For the same reason, its sorting in Step 2 can be done in $O(n)$, for example, using radix sort. Thanks to this sorting, the detection of pairs to be dropped at Step 2, as well as that of triplets to be generated at Step 3 can also be done in linear time, because pairs that start with the same position are now consecutive. Each newly generated list is either a permutation of the previous one, or even a subset of it, and thus the size remains in $O(n)$. Therefore, also the sorting of Step 5 can be done in linear time using radix sort because also the number of distinct classes cannot be larger than n . This sorting allows to detect in linear time the triplets to be dropped at the same step. Similarly, it allows the final detection of maximal runs to be done in linear time as well. Therefore, the overall time and space complexity is in $O(n)$. \square

Notice that the elimination, at each iteration, of masks that do not satisfy the quorum, actually results in a practical important reduction of the search space. Notice also that if we build the classes of a new mask obtained by the overlapping of two shorter masks (rather than their concatenation), then the very same procedure can be applied (having the same complexity) with the only difference that at Step 2 one should build pairs $\langle i' - \delta, n_{C'} \rangle$ instead of $\langle i' - |\mu|, n_{C'} \rangle$, where δ is the size of the overlap.

Figure 4: Trie for $\pi_{\mu'}$ where $\mu' = 0110$

Figure 5: Trie for $\pi_{\mu''}$ where $\mu'' = 0011$.

3.2 Equivalent masks

We now list some interesting properties that allow us to generate half of all possible masks that are at most 2^L . Although these properties do not improve over the worst-case complexity, they are useful optimizations for the practical behavior of our algorithm.

We first need to introduce some notation. Intuitively, consider the partitions shown in Figures 1, 4, and 5, respectively, for masks $\mu = 1100$, $\mu' = 0110$, and $\mu'' = 0011$. The classes in these partitions are reported in Table 1, so that their mutual dependence is highlighted. Ignoring border effects in the first and the last $L - 1$ positions of the text, we can say that μ , μ' , and μ'' conceptually represent the same set of patterns: those that have only two solid that are adjacent. Each row of Table 1 puts the classes into a one-to-one correspondence, in which a class can be obtained from another by adding an integer d to the positions, such that $|d|$ is exactly the amount of shifted symbols in their masks needed to make them equal. For example, class $\{5, 6, 7\}$ can be obtained from $\{7, 8, 9\}$ by adding the integer -2 to the positions of the latter, since μ'' can be transformed into μ shifting its symbols by 2 positions to the left.

Given two masks μ and μ' of the same length, we say that they are *equivalent* if they have the same number of 1s and μ can be obtained from μ' by a shift of the symbols, as long as only 0 exceed the mask border and are then added circularly on the opposite side (also the vice versa holds). We define the following notations for sets of positions. Given two sets C and C' of text positions in $[0 \dots n - 1]$, we say that $C \equiv_d C'$ for an integer d if two conditions hold:

1. for each $i' \in C'$ such that $0 \leq i' + d \leq n - 1$, we have that $i' + d \in C$;
2. for each $i \in C$ such that $0 \leq i - d \leq n - 1$, we have that $i - d \in C'$.

It is easy to see that, removing border effects, there is a one-to-one correspondence between the classes of the partitions induced by two equivalent masks that actually coincides with the relation \equiv_d , where d is the size of the shift. In other words, the following result clearly holds.

Lemma 5 *For any two equivalent masks μ and μ' , there exists an integer d inducing a one-to-one correspondence between the classes of the partitions π_μ and $\pi_{\mu'}$ as follows: for each class $C \in \pi_\mu$, we have a unique class $C' \in \pi_{\mu'}$ such that $C \equiv_d C'$ (and vice versa).*

It is worth noting that π_μ and $\pi_{\mu'}$ have the same number of equivalence classes, the same number of elements and any position shifted by an integer $d > 0$. We can ignore the first and the last $L - 1$ text positions or, alternatively, extend the text with $L - 1$ null symbols to its left and its right. Lemma 5 makes partitions $\pi_{\mu'}$ redundant with respect to π_μ .

Consequently, for each group of equivalent masks, we choose as representative the one having 1 as its first symbol. (Such a mask always exists except for the mask $00 \dots 0$, which forms a trivial singleton group whose partition contains just one class made up of all the text positions.) We then eliminate the other masks in the class from our computation (steps 1–6), since we can always recover their partition by adding a suitable integer d . For example, with $L = 4$, we now build explicitly only the following masks when applying Lemma 4, which illustrates the fact that the number of representatives is *at most half* the number of equivalent masks: only those that start with a 1.

$$\begin{aligned} F_1 &= \{1\}, \\ F_2 &= \{11, 10\}, \\ F_4 &= \{1111, 1110, 1101, 1100, 1011, 1010, 1001, 1000\}. \end{aligned}$$

3.3 Algorithm KMR for masks

We now have all the ingredients for describing our algorithm that applies KMR to solve Problem 1. Given a string T , a length L and a quorum q , we first compute the partition for $F_1 = \{1\}$ (since that for the mask

Figure 1	Figure 4	Figure 5
{4}	{3}	{2}
{5,12}	{4,11}	{3,10}
-	-	{12}
{7,8,9}	{6,7,8}	{5,6,7}
{10}	{9}	{8}
{3,11}	{2,10}	{1,9}
-	{12}	{11}
{6}	{5}	{4}
{0,1,2}	{0,1}	{0}

Table 1: The partitions shown in Figures 1, 4, and 5.

0 is trivial). Next, we compute F_2 , F_4 and so on, as described in Section 3.1–3.2. In particular, we just compute the representative for equivalent masks, when running steps 1–6: for each pair of mask $\mu, \mu' \in F_\ell$ (not necessarily distinct), we build $F_{2\ell}$ by processing the concatenation of μ with all possible shifts of μ' (whose partitions we can quickly recover from that of μ'). Summing up, we actually perform the following two steps:

1. Scan T to construct the partition induced by relations \equiv_1 and build F_1 .
2. Use Lemma 3 (and the ideas in Section 3.2) to construct, successively, $\{F_2, F_4, F_8, \dots, F_{2^r}, F_L\}$, where r is the largest value such that $2^r < L$.

Theorem 6 *Using the KMR approach, we can build the set $\mathcal{Q}(L, T, q)$ of masks for Problem 1, along with their induced partitions, in $O(2^L n)$ time and space.*

Proof: Given a length L , we have at most $\sum_{\ell=1}^r 2^\ell + 2^L < 2^{L+1}$ different masks μ to consider, and so as many equivalence relations \equiv_μ , during the execution of the procedure Lemma 4 refers to. Since this latter guarantees that such execution takes $O(n)$ per mask, we have that KMR requires a total of $O(2^L n)$ time and space. \square

Corollary 7 *Using the KMR approach, we can build $\mathcal{Q}(L, T, q)$ in external memory with $O(2^L \text{sort}(n))$ block transfers, where $\text{sort}(n)$ is the I/O complexity of sorting n integers in $[0 \dots n^{O(1)}]$.*

In order to solve Problem 1, we need to select the maximal masks that form the set $\mathcal{M}(L, T, q) \subseteq \mathcal{Q}(L, T, q)$. This maximality check on the set $\mathcal{Q}(L, T, q)$ computed in Theorem 6 can be done *a posteriori* using the LESS algorithm proposed in [13], with a cost that is linear on the average and quadratic in the worst-case with the size of the input. Given that we have up to 2^L masks of length L , such method has average time complexity in $O(L2^L)$, and worst case complexity in $O(L2^{2L})$.

Corollary 8 *Using the KMR approach, we can solve Problem 1 in $O(L2^L n)$ average time and space, and $O(L2^{2L} n)$ worst-case time and space.*

One drawback of the algorithm in Corollary 8 is that if a maximal mask μ is discovered and has a certain number k of 1s in it, we have to generate anyway all the $\Theta(2^k)$ masks μ' such that $\mu' \preceq \mu$. In Section 4, we show how to implant a suitable branch-and-bound method to KMR, so as to obtain a theoretically faster $O(2^L n)$ solution to solve Problem 1 (see Theorem 10).

4 Adaptive KMR for Maximal Masks

In this section, we describe how to improve the worst-case complexity $O(L2^{2L} n)$ of Corollary 8. The reason for this bound is that, when using the KMR approach described in Section 3, we first generate all the masks of length L having quorum and, then, select *a posteriori* those being maximal too. Our task can be better viewed in terms of the lattice $\mathcal{L} = \langle \{1, 0\}^L, \preceq \rangle$ of 2^L masks, introduced in Section 2.2 and illustrated in Figure 6. Given a mask μ having quorum, we would like to avoid to compute the partitions for mask μ' , such that $\mu' \preceq \mu$. Recall that μ' is called predecessor of μ , and the latter is called successor of μ' . In general, given

Figure 6: Two binomial (spanning) trees B_L , rooted at the top and at the bottom of lattice \mathcal{L} .

a mask, its successors are those masks that are reachable going upward in the lattice \mathcal{L} and its predecessors are those reachable going downward.

Our idea can be summarized as follows. Suppose that we compute the set $\mathcal{Q}(L/2, T, q)$ of all masks of length $L/2$ that have quorum q , using Theorem 6, in $O(2^{L/2}n)$ time and space. (We store these masks using perfect hashing [5], so each can be retrieved in $O(1)$ worst-case time; also, we assume without loss of generality that L is an even number.) Instead of considering all $O(2^{L/2}) \times O(2^{L/2})$ possible concatenations of two masks in $\mathcal{Q}(L/2, T, q)$, we perform concatenation on demand, thus obtaining an *adaptive KMR approach*. The masks of length $L/2$, which we decide to concatenate in $O(n)$ time using Lemma 4, are chosen according to a suitable traversal of the lattice \mathcal{L} . We can employ two different pruning strategies, analogously to the *a priori* algorithm [16], where μ is the current mask of length L in a traversal of \mathcal{L} :

1. if μ has the quorum, we do not check its predecessors in \mathcal{L} since they have quorum but cannot be maximal (since they are less specific);
2. if μ has not the quorum, we do not check its successors because they cannot have quorum either (since having quorum is an anti-monotone property).

In the former case, we simulate the traversal of \mathcal{L} by enumerating the masks of length L starting from $1 \dots 1$ and proceeding to less specific masks, which are downward in \mathcal{L} , while in the latter case we start from $0 \dots 0$ and go upward looking for more specific masks. In both cases, whenever we need to build the partition π_μ for the current mask μ , we split it as $\mu = \mu_1\mu_2$ such that $|\mu_1| = |\mu_2| = L/2$. If both $\mu_1, \mu_2 \in \mathcal{Q}(L/2, T, q)$, we apply Lemma 4 on them to check whether μ has quorum and compute its partition π_μ . Otherwise, we declare that μ does not have quorum in T . We denote this checking operation by the predicate $P_T(\mu)$.

Since we apply Lemma 4 to at most 2^L masks of length L , the complexity of the algorithm is $O(2^L n)$. It remains to see how to enumerate the masks avoiding those that are non-maximal.

Since the mask lattice \mathcal{L} is isomorphic to the powerset $\mathcal{P}(\{0, \dots, L-1\})$ (see Section 2.2), our implicit traversal of \mathcal{L} can be obtained by enumerating all the subsets of $\{0, \dots, L-1\}$ visiting the corresponding *binomial tree* [33], which is also a spanning tree for \mathcal{L} . We recall that a binomial tree B_k is an ordered tree representing the set $\mathcal{P}(\{0, \dots, k-1\})$, and can be defined recursively as follows: B_0 consists of a single node and, for $k > 0$, B_k consists of two binomial trees B_{k-1} , where the root of the former is added as the rightmost child to the root of the latter. Figure 6 shows two binomial trees B_L , both spanning \mathcal{L} . The first tree, shown on the left, is rooted at the top of the lattice and can be employed to implement the first pruning strategy, avoiding to visit the predecessors of the mask μ (downward in the lattice). The second tree, shown on the right, is rooted at the bottom of the lattice and can be employed to implement the second pruning strategy, avoiding to visit the successors of μ . Since we want to identify the maximal masks, we opt for the first tree, starting from the top of the lattice \mathcal{L} .

Algorithm 1 shows the main steps when starting on top. Line 1 checks if mask $1 \dots 1$ has quorum and, if this is so, that mask is the only one returned since any other masks would be less specific. Otherwise, create a queue D containing $1 \dots 1$. As long as D is not empty, the main loop on line 1 selects a mask μ and generates one of its immediate predecessors μ' that are not yet visited. (In order to obtain μ' , function `next_immediate_predecessor` systematically switches a 1 into 0 in μ at a time, so $\mu' \preceq \mu$ and they differ in one bit.) If μ' does not exist, then it means that all of μ predecessors have been already visited and μ is dequeued; otherwise, line 1 checks if μ' has quorum. If this is so, μ' is added to \mathcal{M} if and only if a more specific mask is not already in it (lines 1–1). Otherwise, μ' is not enqueued, because anti-monotonicity of $P_T(\mu)$ guarantees that all of its predecessors have also quorum but they are less specific. Instead, if μ' has not quorum, it is enqueued because one of its predecessors could have quorum.

Implementing D as a queue actually leads to a breadth-first visit of the binomial tree, because we visit a node of level $i+1$ when all the node of level i have been removed from the queue. Conversely, implementing D as a stack leads to a depth first visit of the tree, but we have to pay more attention when adding a new mask to \mathcal{M} (i.e. all of the predecessors of the mask must be removed).

Although in the worst case scenario almost all the masks of \mathcal{L} must be checked, the pruning strategies can affect heavily the performance of the algorithms. Proceeding top-down with Algorithm 1, if all the

Algorithm 1 *Top-Down Binomial-Tree Traversal of the Lattice \mathcal{L}*

Input: The predicate $P_T(\mu)$ for checking if μ has quorum q in text T .

Out: The set $\mathcal{M}(L, T, q)$ of all maximal masks with quorum q in T .

```
1:  $\mathcal{M} := \emptyset$ 
2: if  $P_T(1 \dots 1)$  then
     $\{1 \dots 1\}$ 
3: end if
4: Queue  $D := \{1 \dots 1\}$ 
5: while not empty( $D$ ) do
6:    $\mu := \text{top}(D)$ 
7:    $\mu' := \text{next\_immediate\_predecessor}(\mu)$ 
8:   if  $\mu' = \text{null}$  then
9:     remove_top( $D$ )
10:  else
11:    if  $P_T(\mu')$  then
12:      if  $\mathcal{M}$  does not contains  $\mu'$  s.t.  $\mu' \preceq \mu''$  then
13:         $\mathcal{M} := \mathcal{M} \cup \{\mu'\}$ 
14:      end if
15:    else
16:      insert( $\mu', D$ )
17:    end if
18:  end if
19: end while. $\mathcal{M}$ 
```

interesting masks are closed to the top, they are quickly found visiting only a small fraction of the 2^L masks of the lattice. For example, if only $1 \dots 1$ has quorum, the predicate $P_T(\mu)$ is evaluated only once. (Similar considerations hold for the bottom-up traversal of \mathcal{L} .)

Unfortunately, Algorithm 1 cannot yet obtain $O(2^L n)$ time, since checking the condition in line 1 can take time $O(L |\mathcal{M}|) = O(L 2^L)$ per mask, thus giving a total cost of $O(2^L (n + L 2^L))$. We therefore discuss how to refine the breadth-first traversal of Algorithm 1 to get $O(2^L n)$ time.

Consider the lattice \mathcal{L} and call level 0 the top mask $1 \dots 1$, level 1 its predecessors, and so on, up to level L , which is the bottom mask $0 \dots 0$. Also, consider the maximal masks in \mathcal{L} (for the given predicate $P_T(\mu)$) and their predecessors. We call a mask μ on level i *safe*, where $0 \leq i \leq L$, if μ is not predecessor of any maximal mask on levels $0, 1, \dots, i-1$. Note that a safe mask μ itself can be maximal (i.e. $P_T(\mu)$ holds), which is consistent with our definition of safeness.

Our goal is to modify Algorithm 1, so that it runs $P_T(\mu)$ only for safe masks μ on each level $i = 0, 1, \dots, L$. The rationale is that, having traversed the first i levels of the lattice \mathcal{L} and having found the maximal masks on these levels, we cannot exclude any safe mask μ on level i without first testing $P_T(\mu)$ on it. We show how to find safe masks on each level. Initially, for $i = 0$, the mask $1 \dots 1$ is trivially safe.

During the top-down (breadth-first) traversal of \mathcal{L} , let us call S_i the set of safe masks on level i . We enforce the invariant that S_i is indeed the set of mask that are in queue D on level i . The other masks on level i are not of interest to us, since they surely have a maximal successor. We show how to produce S_{i+1} , so that the traversal can move the masks from S_{i+1} to the queue D for the next level $i+1$. We need the crucial lemma below.

Lemma 9 *Let M_i be the set of maximal masks on level i , where $0 \leq i \leq L$. Then, the following properties hold for each mask μ :*

- (i) $\mu \in M_i$ if and only if $\mu \in S_i$ and $P_T(\mu)$ holds (hence, $M_i \subseteq S_i$).
- (ii) $\mu \in S_{i+1}$ if and only if all the immediate successors of μ are in $S_i - M_i$.

Proof:

- (i) Since μ is maximal, it is interesting by definition while no one of its successors can be, hence it is

safe. The converse also holds because if μ is interesting and no one of its successors is maximal (hence interesting), then it is maximal by definition.

- (ii) A mask μ at level i can be maximal ($\mu \in M_i$), or interesting but not maximal or safe but not maximal ($\mu \in S_i - M_i$): no other cases are possible. In particular μ cannot be both not safe and not interesting, because it would have at least a maximal mask among its successors forcing μ to be interesting (since the anti-monotonicity of $P_T(\mu)$). The first implication (\Rightarrow) is equivalent to say that if there exist μ' immediate successor of μ such that $\mu' \notin S_i - M_i$ then $\mu \notin S_{i+1}$. Since the observation above $\mu' \notin S_i - M_i$ implies that μ' is interesting, hence μ' of one of its successors must be maximal and μ cannot be safe having a maximal mask among its successors. To prove the second implication (\Leftarrow) it suffices to observe that if all the immediate successors μ' of μ are in $S_i - M_i$, they are not interesting and since the anti-monotonicity of $P_T(\mu)$ no one of their successors up in the lattice can be, hence μ is safe.

□

We now show how to exploit Lemma 9 during the traversal. First, since the queue D stores the set S_i , we examine the masks $\mu \in S_i$ and perform the check with $P_T(\mu)$ by Lemma 9(i). Second, we remove M_i from the queue D , which now stores the set $S_i - M_i$. In order to apply Lemma 9(ii), we recall that each of the immediate predecessors and immediate successors of a mask μ differs from μ in exactly one position. We generate all immediate predecessors of the masks in the queue $D = S_i - M_i$. In this way, we create a superset of S_{i+1} from which we select only the masks that have all their *immediate successors* in the queue. We detail more this task and state its complexity.

Recall that D denotes the set $S_i - M_i$ (this is indeed the content of the queue after the removal of the maximal motifs in it). We generate all the immediate predecessors of the masks in D as follows. Given a mask $\mu \in D$ of arbitrary length L , for any position j of a symbol 1 in μ , we generate the immediate predecessor of μ that have all symbols equal to those in μ except that it contains symbol 0 in position j . Let P_D be the multiset of immediate predecessors so built, which represent the predecessors of the masks in D . By Lemma 9, we have that P_D is a superset of S_{i+1} and a mask $\mu \in S_{i+1}$ if and only if μ has multiplicity $i+1$, that is, μ occurs $i+1$ times in the multiset P_D . For example, supposing $S_2 = \{10011, 11001, 10101\}$, we have $S_3 = \{10001\}$ since it appears three times in $P_D = \{00011, 10001, 10010, 01001, 10001, 11000, 00101, 10001, 10100\}$.

We can proceed in several ways for this checking. Either we sort the multiset P_D and output the masks that appear (consecutively) $i+1$ times in the sorted multiset (e.g. Google's `mapReduce`) or we build a trie on the strings in P_D and output those stored in the leaves with multiplicity $i+1$. From a theoretically point of view, we can build a perfect hash function $f()$ on the distinct values in P_D in $O(|P_D|)$ time and space [5]. In this way, given any two masks μ and μ' , we have that $f(\mu) = f(\mu')$ implies $\mu = \mu'$. We then use an array of counters C initially set to zero. For each mask $\mu \in P_D$, we increment $C[f(\mu)]$ by one. At the end, with a further scan of P_D , for each mask μ , if $C[f(\mu)] = i+1$ then we output μ and reset $C[f(\mu)]$ to zero.

The overall cost for finding the sets S_i for all levels i can be bounded by $O(\sum_{i=0}^{L-1} (|S_i| L))$ time since $|P_D| \leq |S_i| (L-i)$ for level $i+1$. Using the fact the $|S_i| \leq \binom{L}{i}$, we obtain a cost of $O(\sum_{i=0}^{L-1} \binom{L}{i} L) = O(L2^L)$ for all the masks (instead of paying this cost for a single mask as before).

We can now apply Algorithm 1 in which we do *not* run the test in line 1 but, rather, we follow the traversal indicated by sets S_i on each level i of the lattice of \mathcal{L} . In this way the cost is dominated by $O(2^L n)$ due to checking $P_T(\mu)$ for the masks $\mu \in \cup_{i=0}^L S_i$, since the cost of generating the sets S_i is $O(2^L L) = O(2^L n)$. The required space depends on the chosen visit strategy. Since we use the breadth-first traversal, the queue D can contain all the $\binom{L}{i}$ masks on level i (i.e. it may happen $|S_i| = \binom{L}{i}$). We have thus proved our main result.

Theorem 10 *Using the adaptive KMR approach, we can solve Problem 1 computing $\mathcal{M}(L, T, q)$ in $O(2^L n)$ worst-case time and space.*

In order to evaluate the cost stated in Theorem 10, consider the following scenario in which, after preprocessing the text T in polynomial time, we ideally check, in constant time per pattern, if any of the $(|\Sigma| + 1)^L$ candidate patterns has quorum and is maximal. We compare the cost of $O(n^{O(1)} + (|\Sigma| + 1)^L)$ thus obtained against the cost of $O(2^L n)$ stated in Theorem 10. Since $2^L n = O((|\Sigma| + 1)^L)$ when $L =$

$\Omega(\log_{|\Sigma|} n)$, and $2^L n = O(n^{O(1+1/\log_2 |\Sigma|)})$ otherwise, we can conclude that the cost in Theorem 10 can be upper bounded by $O(n^{O(1+1/\log_2 |\Sigma|)} + \min\{2^L n, (|\Sigma| + 1)^L\})$. The latter is always better than the ideal bound of $O(n^{O(1)} + (|\Sigma| + 1)^L)$, that is, than constant-time enumerating and checking all the potential $(|\Sigma| + 1)^L$ patterns in T . As we discussed in Section 1, more sophisticated techniques that are the state of the art cannot improve, in the worst case, over the bound given in Theorem 10.

In order to get an external-memory algorithm, we notice that the both the construction of the safe masks and the concatenation of two masks of length $L/2$ needed to check $P_T(\mu)$ on each safe mask, require scanning and sorting. These are well studied algorithms in the two-level model of computation, and the sorting is the dominating cost.

Corollary 11 *Using the adaptive KMR approach, we can build $\mathcal{M}(L, T, q)$ in external memory with $O(2^L \text{sort}(n))$ block transfers, where $\text{sort}(n)$ is the I/O complexity of sorting n items.*

Finally, we conclude by observing that our bounds hold also for the case in which Problem 1 has the additional requirements of reporting one representative mask for each equivalence class of masks. We recall from Section 3.2 that any two masks μ and μ' of the same length are *equivalent* if they have the same number of 1s and μ can be obtained from μ' by a shift of the symbols (and vice versa). In our introductory example, 1110 and 0111 are equivalent and we take the leftmost shift as the representative.

We can implement this extension with minor variations in our algorithms. In particular, we append $L - 1$ copies of a new special symbol $\$$ that is an endmarker for the text T , and so it does not belong Σ . We then run our algorithms on the resulting text $T\$ \$ \dots \$$ and traverse the lattice \mathcal{L} by considering only the masks having 1 as first symbol (see the subtree induced by these masks in each of the binomial trees shown in Figure 6). In this way, whenever we consider a mask, it is always the leftmost shift of its equivalence class, and we cover all this kind of masks.

Corollary 12 *Using the adaptive KMR approach, we can solve Problem 1 modulo the equivalence between the masks, by selecting the leftmost shifts of the maximal masks in $\mathcal{M}(L, T, q)$ in $O(2^L n)$ worst-case time and space.*

5 Conclusions

We have introduced a new notion of motifs that succinctly represent the repeated patterns for an input sequence. We have described how to build the set of all maximal masks of length L in $O(2^L n)$ time and space in the worst case. This bound is better than constant-time enumerating and checking all the potential $(|\Sigma| + 1)^L$ candidate patterns in T after a polynomial-time preprocessing of T . It is also better than the bound obtained by the algorithms for mining frequent itemsets using hypergraph traversals. Interestingly, our algorithms can work also in external memory and in a distributed setting since they hinge on scanning and sorting sequential data. We plan to make a systematic experimental study to find practical incarnations of our ideas using real-life data sets. It is an open problem to obtain an output-sensitive algorithm, whose time complexity is of the form $o(2^L) + \Theta(|\mathcal{M}|)$ in the worst case.

Acknowledgments

We are in debt to Romeo Rizzi for having brought to our attention the connection of the problem to hypergraph traversals. We thank Elias C. Stavropoulos for having kindly sent us his software.

References

- [1] C. Berge. *Hypergraphs: Combinatorics of Finite Sets*, volume 45. Elsevier, 1989.
- [2] B. Brejová, D. Brown, and T. Vinar. Vector seeds: an extension to spaced seeds allows substantial improvements in sensitivity and specificity. In *WABI*, volume 2812 of *LNCS*, pages 39–54. Springer, 2003.

- [3] S. Burkhardt and J. Kärkkäinen. Better filtering with gapped q-grams. In *CPM*, volume 2089 of *LNCS*, pages 73–85. Springer, 2001.
- [4] K.P. Choi, F. Zeng, and L. Zhang. Good spaced seeds for homology search. *Bioinformatics*, 20(7):1053–1059, 2004.
- [5] Z. J. Czech, G. Havas, and B. S. Majewski. Perfect hashing. *Theoret. Comput. Sci.*, 182(1-2):1–43, 1997.
- [6] J.-E. Duchesne, M. Giraud, and N. El-Mabrouk. Seed-based exclusion method for non-coding RNA gene search. In *COCOON*, volume 4598 of *LNCS*, pages 27–39. Springer, 2007.
- [7] T. Eiter, K. Makino, and G. Gottlob. Computational aspects of monotone dualization: A brief survey. *Discr. Appl. Math.* to appear. doi:10.1016/j.dam.2007.04.017, 2007.
- [8] R. Eres, G.M. Landau, and L. Parida. A combinatorial approach to automatic discovery of cluster-patterns. In *WABI*, volume 2812 of *LNCS*, pages 139–150. Springer, 2003.
- [9] E. Eskin. From profiles to patterns and back again: a branch and bound algorithm for finding near optimal motif profiles. In *RECOMB*, pages 115–124, 2004.
- [10] S. Feng and E. Tillier. A fast and flexible approach to oligonucleotide probe design for genomes and gene families. *Bioinformatics*, 23(10):1195–1202, 2007.
- [11] M.L. Fredman and L. Khachiyan. On the complexity of dualization of monotone disjunctive normal forms. *J. Algorithms*, 21(3):618–628, 1996.
- [12] M. Frigo, C. E. Leiserson, H. Prokop, and S. Ramachandran. Cache-oblivious algorithms. In *40th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 285–297, 1999.
- [13] P. Godfrey, R. Shipley, and J. Gryz. Algorithms and analyses for maximal vector computation. *VLDB Journal: Very Large Data Bases*, 16(1):5–28, October 2006.
- [14] D. Gunopulos, R. Khardon, H. Mannila, S. Saluja, H. Toivonen, and R. Sewak Sharma. Discovering all most specific sentences. *ACM Trans. Database Syst.*, 28(2):140–174, 2003.
- [15] D. Gunopulos, H. Mannila, R. Khardon, and H. Toivonen. Data mining, hypergraph transversals, and machine learning (extended abstract). In *PODS*, pages 209–216, 1997.
- [16] J. Han and M. Kamber. *Data Mining: Concepts and Techniques*. Morgan Kaufmann, 2001.
- [17] L. Ilie and S. Ilie. Fast computation of good multiple spaced seeds. In *WABI*, volume 4645 of *LNBI*, pages 346–358. Springer, 2007.
- [18] R.M. Karp, R.E. Miller, and A.L. Rosenberg. Rapid identification of repeated patterns in strings, trees and arrays. In *STOC*, pages 125–136, 1972.
- [19] D.J. Kavvadias and E.C. Stavropoulos. An efficient algorithm for the transversal hypergraph generation. *J. Graph Algorithms and Applications*, 9(2):239–264, 2005.
- [20] U. Keich, M. Li, B. Ma, and J. Tromp. On spaced seeds for similarity search. *Discr. Appl. Math.*, 138(3):253–263, 2004.
- [21] L. Khachiyan, E. Boros, K. Elbassioni, and V. Gurvich. An efficient implementation of a quasi-polynomial algorithm for generating hypergraph transversals and its application in joint generation. *Discr. Appl. Math.*, 154(16):2350–2372, 2006.
- [22] D.E. Knuth. *The Art of Computer Programming: Sorting and Searching*, volume 3. Addison Wesley, 2nd edition, 1997.
- [23] G. Kucherov, L. Noé, and M.A. Roytberg. Multiseed lossless filtration. *IEEE/ACM Trans. Comput. Biology Bioinform.*, 2(1):51–61, 2005.

- [24] G. Kucherov, L. Noé, and M.A. Roytberg. A unifying framework for seed sensitivity and its application to subset seeds. *J. Bioinform. and Comput. Biology*, 4(2):553–570, 2006.
- [25] G. Kucherov, L. Noé, and M.A. Roytberg. Subset seed automaton. In *CIAA*, pages 180–191, 2007.
- [26] M. Li, B. Ma, D. Kisman, and J. Tromp. Patternhunter ii: Highly sensitive and fast homology search. *J. Bioinform. and Comput. Biology*, 2(3):417–440, 2004.
- [27] M. Michael, F. Nicolas, and E. Ukkonen. On the complexity of finding gapped motifs. *CoRR*, abs/0802.0314, 2008.
- [28] L. Parida, I. Rigoutsos, A. Floratos, D. Platt, and Y. Gao. Pattern Discovery on Character Sets and Real-valued Data: Linear Bound on Irredundant Motifs and Efficient Polynomial Time Algorithm. In *SIAM SODA*, 2000.
- [29] N. Pisanti, M. Crochemore, R. Grossi, and M.-F. Sagot. Bases of motifs for generating repeated patterns with wild cards. *IEEE/ACM Trans. Comput. Biology Bioinform.*, 2(1):40–50, 2005.
- [30] Y. Sun and J. Buhler. Designing multiple simultaneous seeds for dna similarity search. *J. Comput. Biology*, 12(6):847–861, 2005.
- [31] Y. Sun and J. Buhler. Designing patterns for profile hmm search. *Bioinformatics*, 23(2):42–43, 2007.
- [32] E. Ukkonen. Structural analysis of gapped motifs of a string. In *MFCS*, pages 681–690, 2007.
- [33] J. Vuillemin. A data structure for manipulating priority queues. *CACM*, 21(4):309–315, 1978.