

Subset seed automaton

Gregory Kucherov¹, Laurent Noé¹, and Mikhail Roytberg²

¹ LIFL/CNRS/INRIA, Bât. M3 Cité Scientifique, 59655, Villeneuve d'Ascq cedex, France, {Gregory.Kucherov,Laurent.Noé}@lifl.fr

² Institute of Mathematical Problems in Biology, Pushchino, Moscow Region, 142290, Russia, mroytberg@mail.ru

Abstract. We study the pattern matching automaton introduced in [1] for the purpose of seed-based similarity search. We show that our definition provides a compact automaton, much smaller than the one obtained by applying the Aho-Corasick construction. We study properties of this automaton and present an efficient implementation of the automaton construction. We also present some experimental results and show that this automaton can be successfully applied to more general situations.

1 Introduction

The technique of *spaced seeds* for similarity search in strings (sequences) was introduced about five years ago [2,3] and constituted an important algorithmic development [4,5]. Its main applications have been approximate string matching [2] and local alignment of DNA sequences [3,6,7] but the underlying idea applies also to other algorithmic problems on strings [8,9].

Since the invention of spaced seeds, different generalizations have been proposed, such as seeds with match errors [10,11], *daughter seeds* [12], *indel seeds* [13], or *vector seeds* [14]. In [1], we proposed the notion of *subset seeds* and demonstrated its advantages and its usefulness for DNA sequence alignment. In the formalism of subset seeds, an alignment is viewed as a text over some alphabet \mathcal{A} , and a seed as a pattern over a subset alphabet $\mathcal{B} \subseteq 2^{\mathcal{A}}$. The only requirements made is that \mathcal{A} contains a special letter $\mathbf{1}$, \mathcal{B} contains a letter $\# = \{\mathbf{1}\}$, and every letter of \mathcal{B} contains $\mathbf{1}$ in its set. The matching relation is naturally defined: a seed letter $b \in \mathcal{B}$ matches a letter $a \in \mathcal{A}$ iff a belongs to the set b .

For any seed-based similarity search method, including all above-mentioned types of seeds, an important issue is an accurate estimation of the sensitivity of a seed with respect to a given probabilistic model of alignments. For different probabilistic models, this problem has been studied in [15,16,17]. In [1] we proposed a general framework for this problem that allows one to compute the seed sensitivity for different definitions of seed and different alignment models. This approach is based on a finite automata representation of the set of target alignments and the set of alignments matched by a seed, as well as on a representation of the probabilistic model of alignments as a finite-state transducer.

A key ingredient of the approach of [1] is a finite automaton that recognizes the set of alignments matched (or *hit*) by a given subset seed. We call this automaton a *subset seed automaton*. The size (number of states) of the subset seed automaton is crucial for the efficiency of the whole algorithm of [1]. Note that the algorithm of [16] is also based on an automaton construction, namely on the Aho-Corasick automaton implied by the well-known string matching algorithm.

Besides its application to the seeding technique for similarity search and string matching, constructing an efficient subset seed automaton is an interesting problem in its own, as it provides a solution to a variant of the *subset matching problem* studied in literature [18,19,20].

In this paper, we study properties of the subset seed automaton and present an efficient implementation of its construction. More specifically, we obtain the following results:

- we present a construction of subset seed automaton that has $\mathcal{O}(w2^{s-w})$ states, compared to $\mathcal{O}(w|\mathcal{A}|^{s-w})$ implied by the Aho-Corasick construction, where s and w are respectively the *span* and the *weight* of the seed defined in the next Section,
- we further motivate our construction by showing that for some seeds, our construction gives the minimal automaton,
- we prove that our automaton is *always* smaller than the one obtained by the Aho-Corasick construction; we provide experimental data that confirm that for $|\mathcal{A}| = 2$, our automaton is on average about 1.3 times bigger than the minimal one, while the Aho-Corasick automaton is about 2.5 times bigger. For $|\mathcal{A}| = 3$ the difference is much more substantial: while our automaton is still about 1.3 times bigger than the minimal one, the Aho-Corasick automaton turns out to be about 17 times bigger,
- we provide an efficient algorithm that implements the construction of the automaton such that each transition is computed in constant time,
- we show that our construction can be applied to the case of multiple seeds and to the general subset matching problem.

The presented automaton construction is implemented in full generality in HEDERA software package (<http://bioinfo.lifl.fr/yass/hedera.php>) and has been applied to the design of efficient seeds for the comparison of genomic sequences.

2 Subset seed matching

The goal of seeds is to specify short string patterns that, if shared by two strings, have best chances to belong to a larger similarity region common to the two strings. To formalize this, a similarity region is modeled by an alignment between two strings. Usually one considers *gapless alignments* that, in the simplest case, are viewed as sequences of matches and mismatches and are easily specified by binary strings $\{0, 1\}^*$, where 1 is interpreted as “match” and 0 as “mismatch”. A *spaced seed* is a string over binary alphabet $\{\#, _ \}$. The length of π is called its *span* and the number of $\#$ is called its *weight*. A spaced seed $\pi \in \{\#, _ \}^s$

matches (or *hits*) an alignment $A \in \{0, 1\}^*$ at a position p if for all $i \in [1..s]$, $\pi[i] = \#$ implies $A[p + i - 1] = 1$.

In [1], we proposed a generalization of this basic framework, based on the idea to distinguish between different types of mismatches in the alignments. This leads to representing both alignments and seeds as words over larger alphabets. In the general case, consider an alignment alphabet \mathcal{A} of arbitrary size. We always assume that \mathcal{A} contains a symbol 1 , interpreted as “match”. A *subset seed* is defined as a word over a *seed alphabet* \mathcal{B} , such that

- each letter $b \in \mathcal{B}$ denotes a subset of \mathcal{A} that contains 1 ($b \in 2^{\mathcal{A}} \setminus 2^{\mathcal{A} \setminus \{1\}}$),
- \mathcal{B} contains a letter $\#$ that denotes subset $\{1\}$.

As before, s is called the *span* of π , and the *#-weight* of π is the number of $\#$ in π . A subset seed $\pi \in \mathcal{B}^s$ *matches* an alignment $A \in \mathcal{A}^*$ at a position p iff for all $i \in [1..s]$, $A[p + i - 1] \in \pi[i]$.

Example 1. For DNA sequences over the alphabet $\{A, C, G, T\}$, in [21] we considered the alignment alphabet $\mathcal{A} = \{1, h, 0\}$ representing respectively a match, a transition mismatch ($A \leftrightarrow G, C \leftrightarrow T$), or a transversion mismatch (other mismatch). In this case, the appropriate seed alphabet is $\mathcal{B} = \{\#, @, _\}$ corresponding respectively to subsets $\{1\}$, $\{1, h\}$, and $\{1, h, 0\}$. Thus, seed $\pi = \#@_\#$ matches alignment $A = 10h1h1101$ at positions 4 and 6. The span of π is 4, and the #-weight of π is 2.

One can view the problem of finding seed occurrences in an alignment as a special string matching problem. In particular, it can be considered as a special case of *subset matching* [18] where the text is composed of individual characters. It is also an instance of the problem of matching in indeterminate (degenerate) strings [19,20]. Therefore, an efficient automaton construction that we present in the following sections applies directly to these instances of string matching. One can also freely use the string matching terminology by replacing words “seed” and “alignment” by “pattern” and “text” respectively.

3 Subset Seed Automaton

Let us fix an alignment alphabet \mathcal{A} , a seed alphabet \mathcal{B} , and a seed $\pi = \pi_1 \dots \pi_s \in \mathcal{B}^*$ of span s and #-weight w . Denote $r = s - w$ and let R_π , $|R_\pi| = r$, be the set of all non-# positions in π . Throughout the paper, we identify each position $z \in R_\pi$ with the corresponding prefix $\pi_{1..z} = \pi_1 \dots \pi_z$ of π , and we interchangeably regard elements of R_π as positions or as prefixes of π .

We now define an automaton $S_\pi = \langle Q, q_0, Q_F, \mathcal{A}, \psi : Q \times \mathcal{A} \rightarrow Q \rangle$, $q_0 \in Q$, $Q_F \subseteq Q$, that recognizes the set of all alignments matched by π . The states Q are defined as pairs $\langle X, t \rangle$ such that $X = \{x_1, \dots, x_k\} \subseteq R_\pi$, $t \in [0 \dots s]$, $\max\{X\} + t \leq s$. The automaton maintains the following invariant condition. Suppose that S_π has read a prefix $a_1 \dots a_p$ of an alignment A and has come to a state $\langle X, t \rangle$. Then t is the length of the longest suffix of $a_1 \dots a_p$ of the form 1^i , $i \leq s$, and X contains all positions $x_i \in R_\pi$ such that prefix $\pi_{1..x_i}$ matches a suffix of $a_1 \dots a_{p-t}$.

$$\begin{array}{ll}
 (a) \pi = \# \textcircled{0} \# _ \# \# _ \# \# \# & (c) \begin{array}{l} a_9 t \\ 111\text{h}1011\text{h}\bar{1}\bar{1} \\ \pi_{1..7} = \# \textcircled{0} \# _ \# \# _ \\ \pi_{1..4} = \# \textcircled{0} \# _ \\ \pi_{1..2} = \# \textcircled{0} \end{array} \\
 (b) A = 111\text{h}1011\text{h}11 &
 \end{array}$$

Fig. 1. Illustration to Example 2

Example 2. In the framework of Example 1, consider a seed π and an alignment prefix $A = a_1 \dots a_p$ of length $p = 11$ given in Figure 1(a) and (b) respectively. The length t of the last run of 1's of A is 2. The last non-1 letter of A is $a_9 = \text{h}$. The set R_π of non-# positions of π is $\{2, 4, 7\}$ and π has 3 prefixes belonging to R_π (Figure 1(c)). Prefixes $\pi_{1..2}$ and $\pi_{1..7}$ do match suffixes of $a_1 a_2 \dots a_9$, but prefix $\pi_{1..4}$ does not. Thus, the state of the automaton after reading $a_1 a_2 \dots a_{11}$ is $\langle \{2, 7\}, 2 \rangle$.

The initial state q_0 of S_π is the state $\langle \emptyset, 0 \rangle$. Final states Q_F of S_π are all states $q = \langle X, t \rangle$, where $\max\{X\} + t = s$. All final states are merged into one state $\langle \rangle$.

The transition function $\psi(q, a)$ is defined as follows. If q is a final state, then $\forall a \in \mathcal{A}$, $\psi(q, a) = q$. If $q = \langle X, t \rangle$ is a non-final state, then

- if $a = 1$ then $\psi(q, a) = \langle X, t + 1 \rangle$,
- otherwise $\psi(q, a) = \langle X_U \cup X_V, 0 \rangle$ with
 - $X_U = \{x \mid x \leq t + 1 \text{ and } a \in \pi_x\}$
 - $X_V = \{x + t + 1 \mid x \in X \text{ and } a \in \pi_{x+t+1}\}$

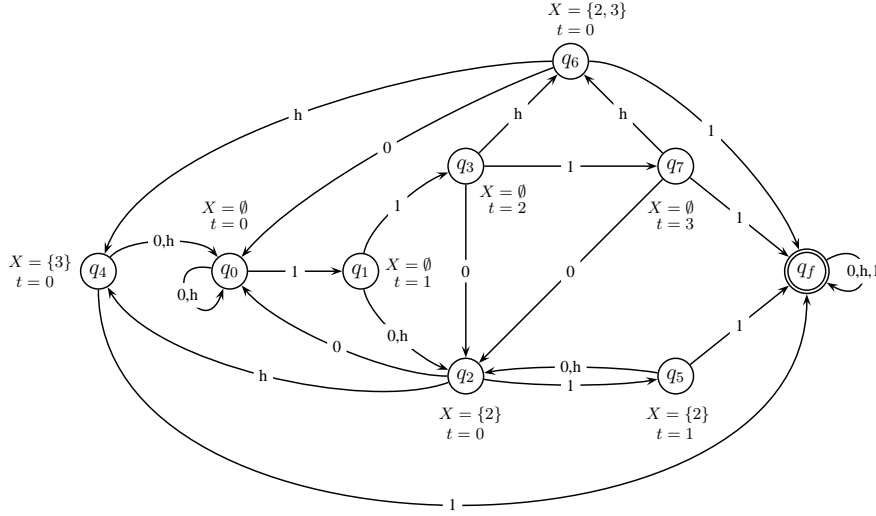


Fig. 2. Illustration to Example 3

Example 3. Still in the framework of Example 1, consider seed $\pi = \# _ \textcircled{0} \#$. Then the set R_π is $\{2, 3\}$. Possible non-final states $\langle X, t \rangle$ of S_π are states $\langle \emptyset, 0 \rangle$, $\langle \emptyset, 1 \rangle$, $\langle \emptyset, 2 \rangle$, $\langle \emptyset, 3 \rangle$, $\langle \{2\}, 0 \rangle$, $\langle \{2\}, 1 \rangle$, $\langle \{3\}, 0 \rangle$, $\langle \{2, 3\}, 0 \rangle$. All these states are reachable in S_π . Figure 2 shows the resulting automaton.

We now study main properties of automaton S_π .

Lemma 1. *The automaton S_π accepts all alignments $A \in \mathcal{A}^*$ matched by π .*

Proof. It can be verified by induction that the invariant condition on the states $\langle X, t \rangle \in Q$ is preserved by the transition function ψ . The final state verifies $\max\{X\} + t = s$ which implies that at the first time S_π gets into the final state, π matches a suffix of $a_1 \dots a_p$. \square

Lemma 2. *The number of states of the automaton S_π is no more than $(w+1)2^r$, where w is the $\#$ -weight of π .*

Proof. Assume that $R_\pi = \{z_1, z_2, \dots, z_r\}$ and $z_1 < z_2 < \dots < z_r$. Let Q_i be the set of non-final states $\langle X, t \rangle$ with $\max\{X\} = z_i$. For states $q = \langle X, t \rangle \in Q_i$ there are 2^{i-1} possible values of X and $s - z_i$ possible values of t between 0 and $s - z_i - 1$, as $\max\{X\} + t \leq s - 1$.

Thus, $|Q_i| \leq 2^{i-1}(s - z_i) \leq 2^{i-1}(s - i)$, and (1)

$$\sum_{i=1}^r |Q_i| \leq \sum_{i=1}^r 2^{i-1}(s - i) = (s - r + 1)2^r - s - 1. \quad (2)$$

Besides states Q_i , Q contains s states $\langle \emptyset, t \rangle$ ($t \in [0..s - 1]$) and one final state. Thus, $|Q| \leq (s - r + 1)2^r = (w + 1)2^r$. \square

Note that if π starts with $\#$, which is always the case for spaced seeds, then $X_i \geq i + 1$, $i \in [1..r]$, and the bound of (1) rewrites to $2^{i-1}(s - i - 1)$. This results in the same $w2^r$ bound on number of states as the one for the Aho-Corasick automaton proposed in [16] for spaced seeds (see also Lemma 4 below).

The next Lemma shows that the construction of automaton S_π is optimal in the sense that no two states can be merged in general.

Lemma 3. *Let $\mathcal{A} = \{0, 1\}$ and $\mathcal{B} = \{\#, _ \}$, where $\# = \{1\}$ and $_ = \{0, 1\}$. Consider a seed $\pi = \# _ \dots _ \#$ with r letters $_$ between two $\#$'s. Then the automaton S_π is reduced, that is*

- (i) *each of its states q is reachable, and*
- (ii) *any two non-final states q', q'' are not equivalent.*

Proof. (i) Let $q = \langle X, t \rangle$ be a non-final state of the automaton S_π , and let $X = \{x_1, \dots, x_k\}$ with $x_1 < \dots < x_k$. Let $A = a_1 \dots a_{x_k} \in \{0, 1\}^*$ be an alignment of length x_k defined as follows: $a_p = 1$ if, for some $i \in [1..k]$, $p = x_k - x_i + 1$, and $a_p = 0$ otherwise. Note that $1 \notin X$ and thus $a_{x_k} = 0$. Thus $\psi(\langle \emptyset, 0 \rangle, A) = \langle X, 0 \rangle$ and finally $\psi(\langle \emptyset, 0 \rangle, A \cdot 1^t) = q$.

(ii) For a set $X = \{x_1, \dots, x_k\}$ and an integer t , denote $X \oplus t = \{x_1 + t, \dots, x_k + t\}$. Let $q' = \langle X', t' \rangle$ and $q'' = \langle X'', t'' \rangle$ be non-final states of S_π . If $\max\{X'\} + t' > \max\{X''\} + t''$, then let $d = (r + 2) - (\max\{X'\} + t')$. Obviously, $\psi(q', 1^d)$ is a final state, and $\psi(q'', 1^d)$ is not.

Now assume that $\max\{X'\} + t' = \max\{X''\} + t''$. Let $g = \max\{v | (v \in X' \oplus t' \text{ and } v \notin X'' \oplus t'') \text{ or } (v \in X'' \oplus t'' \text{ and } v \notin X' \oplus t')\}$. By symmetry,

assume that the maximum is reached on the first condition, i.e. $g = x'_i + t'$ for some $x'_i \in X'$. Let $d = (r + 1) - g$ and consider word $0^d \mathbf{1}$. It is easy to see that $\psi(q', 0^d \mathbf{1})$ is a final state. We claim that $\psi(q'', 0^d \mathbf{1})$ is not. To see this, observe that none of the seed prefixes corresponding to $x \in X''$ with $x + t'' > x'_i + t'$ can lead to the final state on $0^d \mathbf{1}$, due to the last $\#$ symbol of π . The details are left to the reader. \square

Another interesting property of S_π is the existence of a surjective mapping from the states of the Aho-Corasick automaton onto reachable states of S_π . This mapping proves that even if S_π is not always minimized, it has *always* a smaller number of states than the Aho-Corasick automaton. Here, by the Aho-Corasick (AC) automaton, we mean the automaton with the states corresponding to nodes of the trie built according to the classical Aho-Corasick construction [22] from the set of all instances of the seed π . More precisely, given a seed π of span s , the set of states of the AC-automaton is $Q_{AC} = \{A \in \mathcal{A}^* \mid |A| \leq s \text{ and } A \text{ is matched by prefix } \pi_{1..|A|}\}$. The transition $\psi(A, a)$ for $A \in Q_{AC}$, $a \in \mathcal{A}$ yields the *longest* $A' \in Q_{AC}$ which is a suffix of Aa . We assume that all final states are merged into a single sink state.

Lemma 4. *Consider an alignment alphabet \mathcal{A} , a seed alphabet \mathcal{B} and a seed $\pi \in \mathcal{B}^s$ of span s . There exists a surjective mapping $f : Q_{AC} \rightarrow Q$ from the set of states of the Aho-Corasick automaton to the set of reachable states of the subset seed automaton S_π .*

Proof. We first define the mapping f . Consider a state $A \in Q_{AC}$, $|A| = p < s$, where A is matched by $\pi_{1..p}$. Decompose $A = A'1^t$, where the last letter of A' is not 1. If A' is empty, define $f(A) = \langle \emptyset, t \rangle$. Otherwise, $\pi_{1..p-t}$ matches A' and $\pi[p-t] \neq \#$. Let X be a set of positions that contains $p-t$ together with all positions $i < p-t$ such that $\pi_{1..i}$ matches a suffix of A' . Define $f(A) = \langle X, t \rangle$. It is easy to see that $\langle X, t \rangle \in Q$, that $\langle X, t \rangle$ exists in S_π and is reachable by string A .

Now show that for every reachable state $\langle X, t \rangle \in Q$ of S_π there exists $A \in Q_{AC}$ such that $f(A) = \langle X, t \rangle$. Consider a string $C \in \mathcal{A}^*$ that gets S_π to the state $\langle X, t \rangle$. Then $C = C'1^t$ and the last letter of C' is not 1. If X is empty then define $A = 1^t$. If X is not empty, then consider the suffix A' of C' of length $x = \max\{X\}$ and define $A = A'1^t$. Since $\pi_{1..x}$ matches A' , and $x + t \leq s$, then $\pi_{1..x+t}$ matches A and therefore $A \in Q_{AC}$. It is easy to see that $f(A) = \langle X, t \rangle$. \square

Observe that the mapping of Lemma 4 is actually a morphism from the Aho-Corasick automaton to S_π .

Table 1 shows experimentally estimated average sizes of the Aho-Corasick automaton, subset seed automaton, and minimal automaton. The two tables correspond respectively to the binary alphabet (spaced seeds) and ternary alphabet (see Example 1). For Aho-Corasick and subset seed automata, the ratio to the average size of the minimal automaton is shown. Each line corresponds to a seed weight ($\#$ -weight for $|\mathcal{A}| = 3$). In each case, 10000 random seeds of different span have been generated to estimate the average.

$ \mathcal{A} = 2$		Aho-Corasick		S_π		Minimized		$ \mathcal{A} = 3$		Aho-Corasick		S_π		Minimized	
w		avg.	ratio	avg.	ratio	avg.		avg.		avg.	ratio	avg.	ratio	avg.	
9	130.98	2.46	67.03	1.260	53.18	9	1103.5	16.46	86.71	1.293	67.05				
10	140.28	2.51	70.27	1.255	55.98	10	1187.7	16.91	90.67	1.291	70.25				
11	150.16	2.55	73.99	1.254	58.99	11	1265.3	17.18	95.05	1.291	73.65				
12	159.26	2.57	77.39	1.248	62.00	12	1346.1	17.50	98.99	1.287	76.90				
13	168.19	2.59	80.92	1.246	64.92	13	1419.3	17.67	103.10	1.284	80.31				

Table 1. Average number of states of Aho-Corasick, S_π and minimal automaton

4 Subset seed automaton implementation

As in section 3, consider a subset seed π of #-weight w and span s , and let $r = s - w$ be the number of non-# positions. A straightforward generation of the transition table of the automaton S_π can be performed in time $\mathcal{O}(r \cdot w \cdot 2^r \cdot |\mathcal{A}|)$. In this section, we show that S_π can be constructed in time proportional to its size, which is bounded by $(w+1)2^r$, according to Lemma 2. In practice, however, the number of states is usually much smaller.

The algorithm generates the states of the automaton incrementally by traversing them in the breadth-first manner. Transitions $\psi(\langle X, t \rangle, a)$ are computed using previously computed transitions $\psi(\langle X', t \rangle, a)$. A tricky part of the algorithm corresponds to the case where state $\psi(\langle X, t \rangle, a)$ has already been created before and should be retrieved.

The whole construction of the automaton is given in Algorithm 1. We now describe it in more details.

Let $R_\pi = \{z_1, \dots, z_r\}$ and $z_1 < z_2 \dots < z_r$. Consider $X \subseteq R_\pi$. To retrieve the maximal element of X , the algorithm maintains a function $k(X)$ defined by

$$k(X) = \max\{i | z_i \in X\}, \quad k(\emptyset) = 0.$$

Let $q = \langle X, t \rangle$ be a non-final and reachable state of S_π , $X = \{x_1, \dots, x_i\} \subseteq R_\pi$ and $x_1 < x_2 \dots < x_i$. We define $X' = X \setminus \{z_{k(X)}\} = \{x_1, \dots, x_{i-1}\}$ and $q' = \langle X', t \rangle$. The following lemma holds.

Lemma 5. *If $q = \langle X, t \rangle$ is reachable, then $q' = \langle X', t \rangle$ is reachable and has been processed before in a breadth-first computation of S_π .*

Proof. First prove that $\langle X', t \rangle$ is reachable. If $\langle X, t \rangle$ is reachable, then $\langle X, 0 \rangle$ is reachable due to the definition of transition function for $t > 0$. Thus, there is a word A of length $x_i = z_{k(X)}$ such that $\forall j \in [1..r]$, $z_j \in X$ iff the seed suffix $\pi_{1..z_j}$ matches the word suffix $A_{x_i - z_j + 1} \dots A_{x_i}$. Define A' to be the suffix of A of length $x_{i-1} = z_{k(X')}$ and observe that reading A' gets the automaton to the state $\langle X', 0 \rangle$, and then reading $A' \cdot 1^t$ leads to the state $\langle X', t \rangle$. Finally, as $|A' \cdot 1^t| < |A \cdot 1^t|$, then the breadth-first traversal of states of A_π always processes state $\langle X', t \rangle$ before $\langle X, t \rangle$. \square

To retrieve X' from X , the algorithm maintains a function $\text{FAIL}(q)$, similar to the *failure* function of the Aho-Corasick automaton, such that $\text{FAIL}(\langle X, t \rangle) = \langle X', t \rangle$ for $X \neq \emptyset$, and $\text{FAIL}(\langle \emptyset, t \rangle) = \langle \emptyset, \max\{t - 1, 0\} \rangle$.

We now explain how values $\psi(q, a)$ are computed by Algorithm 1. Note first that if $a = 1$, state $\psi(q, a) = \langle X, t + 1 \rangle$ can be computed in constant time (part a. of Algorithm 1). Moreover, since this is the only way to reach state $\langle X, t + 1 \rangle$, it is created and added once to the set of states.

Assume now that $a \neq 1$. To compute $\psi(q, a) = \langle Y, 0 \rangle$, we retrieve state $q' = \text{FAIL}(q) = \langle X', t \rangle$ and then retrieve $\psi(q', a) = \langle Y', 0 \rangle$. Note that this is well-defined as by Lemma 5, q' has been processed before q .

Observe now that since X' and X differ by only one seed prefix $\pi_{1..z_k(X)}$ the only possible difference between Y and Y' can be the prefix $\pi_{1..z_k(X)+t+1}$ depending on whether $\pi_{z_k(X)+t+1}$ matches a or not. As $a \neq 1$, this is equivalent to testing whether $(z_k(X) + t + 1) \in R_\pi$ and $\pi_{z_k(X)+t+1}$ matches a . This information can be precomputed for different values $k(X)$ and t .

For every $a \neq 1$, we define

$$V(k, t, a) = \begin{cases} \{z_k + t + 1\} & \text{if } z_k + t + 1 \in R_\pi \text{ and } \pi_{z_k+t+1} \text{ matches } a, \\ \emptyset & \text{otherwise.} \end{cases}$$

Thus, $Y = Y' \cup V(k(X), t, a)$ (part c. of Algorithm 1). Function $V(k, t, a)$ can be precomputed in time and space $\mathcal{O}(|\mathcal{A}| \cdot r \cdot s)$.

Note that if $V(k, t, a)$ is empty, then $\langle Y, 0 \rangle$ is equal to an already created state $\langle Y', 0 \rangle$ and no new state needs to be created in this case (part e. of Algorithm 1).

If $V(k, t, a)$ is not empty, we need to find out if $\langle Y, 0 \rangle$ has already been created or not and if it has, we need to retrieve it. To do that, we need an additional construction. For each state $q' = \langle X', t \rangle$, we maintain another function $\text{REVMAXFAIL}(q')$, that gives the *last created* state $q = \langle X, t \rangle$ such that $X \setminus z_k(X) = X'$ (part d. of Algorithm 1). Since the state generation is breadth-first, new states $\langle X, t \rangle$ are created in a non-decreasing order of the quantity $(z_k(X) + t)$. Therefore, among all states $\langle X, t \rangle$ such that $\text{FAIL}(\langle X, t \rangle) = \langle X', t \rangle$, $\text{REVMAXFAIL}(\langle X', t \rangle)$ returns the one with the largest $z_k(X)$.

Now, observe that if $V(k, t, a)$ is not empty, i.e. $Y = Y' \cup \{z_k(X) + t + 1\}$, then $\text{FAIL}(\langle Y, 0 \rangle) = \langle Y', 0 \rangle$. Since state $\langle Y, 0 \rangle$ has the maximal possible current value $z_k(Y) + 0 = z_k(X) + t + 1$, by the above remark, we conclude that if $\langle Y, 0 \rangle$ has already been created, then $\text{REVMAXFAIL}(\langle Y', 0 \rangle) = \langle Y, 0 \rangle$. This allows us to check if this is indeed the case and to retrieve the state $\langle Y, 0 \rangle$ if it exists (part d. of Algorithm 1).

The generation of states $\langle X, t \rangle$ with $X = \emptyset$ represents a special case (part b. of Algorithm 1). Here another precomputed function is used:

$$U(t, a) = \cup \{x \mid x \leq t + 1 \text{ and } a \text{ matches } \pi_x\}$$

$U(t, a)$ gives the set of seed prefixes that match the word $1^t \cdot a$. In this case, checking if resulting states have been already added is done in a similar way to $V(k, t, a)$. Details are left out.

Algorithm 1: computation of S_π

Data: a seed $\pi = \pi_1\pi_2 \dots \pi_s$
Result: an automaton $S_\pi = \langle Q, q_0, q_F, \mathcal{A}, \psi \rangle$
 $q_F \leftarrow \text{createstate}(\langle \rangle)$; $q_0 \leftarrow \text{createstate}(\langle \emptyset, 0 \rangle)$;
/ process the first level of states to set FAIL and REVMAXFAIL */*
for $a \in \mathcal{A}$ **do**
 if $a \in \pi_1$ **then**
 if $a = 1$ **then**
 $\langle Y, t_y \rangle \leftarrow \langle \emptyset, 1 \rangle$;
 else
 $\langle Y, t_y \rangle \leftarrow \langle \{1\}, 0 \rangle$;
 if $z_k(Y) + t_y \geq s$ **then**
 $q_y \leftarrow q_F$;
 else
 $q_y \leftarrow \text{createstate}(\langle Y, t_y \rangle)$;
 $\text{FAIL}(q_y) \leftarrow q_0$; $\text{REVMAXFAIL}(q_0) \leftarrow q_y$;
 $\text{push}(\text{Queue}, q_y)$;
 else
 $q_y \leftarrow q_0$;
 $\psi(q_0, a) \leftarrow q_y$;
/ breadth-first processing */*
while $\text{Queue} \neq \emptyset$ **do**
 $q : \langle X, t_X \rangle \leftarrow \text{pop}(\text{Queue})$;
 $q' \leftarrow \text{FAIL}(q)$;
 for $a \in \mathcal{A}$ **do**
 / compute $\psi(\langle X, t_X \rangle, a) = \langle Y, t_y \rangle$ */*
 $q'_y : \langle Y', t'_y \rangle \leftarrow \psi(q', a)$;
 if $a = 1$ **then**
 $Y \leftarrow X$;
 $t_y \leftarrow t_X + 1$;
 else
 if $X = \emptyset$ **then**
 $Y \leftarrow U(t_X, a)$;
 else
 $Y \leftarrow Y' \cup V(k(X), t_X, a)$;
 $t_y \leftarrow 0$;
 / create a new state unless it already exists or it is final */*
 $q_{rev} : \langle Y_{rev}, t_{rev} \rangle \leftarrow \text{REVMAXFAIL}(q'_y)$;
 if $\text{defined}(q_{rev})$ **and** $t_y = t_{rev}$ **and** $Y = Y_{rev}$ **then**
 $q_y \leftarrow q_{rev}$;
 else if $t_y = t'_y$ **and** $Y = Y'$ **then**
 $q_y \leftarrow q'_y$;
 else
 if $z_k(Y) + t_y \geq s$ **then**
 $q_y \leftarrow q_F$;
 else
 $q_y \leftarrow \text{createstate}(\langle Y, t_y \rangle)$;
 $\text{FAIL}(q_y) \leftarrow q'_y$; $\text{REVMAXFAIL}(q'_y) \leftarrow q_y$;
 $\text{push}(\text{Queue}, q_y)$;
 $\psi(q, a) \leftarrow q_y$;

We summarize the results of this section with the following Lemma.

Lemma 6. *After a preprocessing of seed π within time $\mathcal{O}(|\mathcal{A}| \cdot s^2)$, the automaton S_π can be constructed by incrementally generating all reachable states so that every transition $\psi(q, a)$ is computed in constant time.*

5 Possible extensions

An important remark is that the automaton defined in this paper can be easily generalized to the case of multiple seeds. For seeds π^1, \dots, π^k , a state of the automaton recognizing the alignments matched by one of the seeds would be a tuple $\langle X_1, \dots, X_k, t \rangle$, where X_1, \dots, X_k contain the set of respective prefixes, similarly to the construction of the paper. Interestingly, Lemma 4 still holds for the case of multiple seeds. This means that although the size of the union of individual seed automata could potentially grow as the product of sizes, it actually does not, as it is bounded by the size of the Aho-Corasick automaton which grows additively with respect to subsets of underlying words. In practice, our automaton is still substantially smaller than the Aho-Corasick automaton, as illustrated by Table 2. Similar to Table 1, 10000 random seed pairs have been generated here in each case to estimate the average size.

$ \mathcal{A} = 2$		Aho-Corasick		S_π		Minimized	$ \mathcal{A} = 3$		Aho-Corasick		S_π		Minimized
w		avg.	ratio	avg.	ratio	avg.	w		avg.	ratio	avg.	ratio	avg.
9	224.49	2.01	122.82	1.10	111.43		9	2130.6	12.09	201.69	1.15	176.27	
10	243.32	2.07	129.68	1.10	117.71		10	2297.8	12.53	209.75	1.14	183.40	
11	264.04	2.11	137.78	1.10	125.02		11	2456.5	12.86	218.27	1.14	191.04	
12	282.51	2.15	144.97	1.10	131.68		12	2600.6	13.14	226.14	1.14	198.00	
13	300.59	2.18	151.59	1.10	137.74		13	2778.0	13.39	236.62	1.14	207.51	

Table 2. Average number of states of Aho-Corasick, S_π and minimized automata for the case of two seeds

Another interesting observation is that the construction of a matching automaton where each state is associated with a set of “compatible” prefixes of the pattern is a general one and can be applied to the general problem of subset matching [18,23,19,20]. Recall that in subset matching, a pattern is composed of subsets of alphabet letters. This is the case, for example, with IUPAC genomic motifs, such as motif ANDGR representing the subset motif $A[ACGT][AGT]G[AG]$. Note that the text can also be composed of subset letters, with two possible matching interpretations [20]: a seed letter b matches a text letter a either if $a \subseteq b$ or if $a \cap b \neq \emptyset$.

Interestingly, the automaton construction of this paper still applies to these cases with minor modifications due to the absence of text letter 1 matched by any seed letter. With this modification, the automaton construction algorithm of Section 4 still applies. As a test case, we applied it to subset motif $[GA][GA]GGGNNNNAN[CT]ATGNN[AT]NNNNN[CTG]$ mentioned in [20] as a motif describing the translation initiation site in the *E.coli* genome. For a regular 4-letters genomic text, the automaton obtained with our approach has only 138 states, while the minimal automaton has 126 states. For a text composed of 15 subsets of 4 letters and the inclusion matching relation, our automaton contains

139 states, compared to 127 states of the minimal automaton. However, in the case of intersection matching relation, the automaton size increases drastically: it contains 87617 states compared to the 10482 states of the minimal automaton.

References

1. Kucherov, G., Noé, L., Roytberg, M.: A unifying framework for seed sensitivity and its application to subset seeds. *JBCB* **4**(2) (2006) 553–569

19. Holub, J., Smyth, W.F., Wang, S.: Fast pattern-matching on indeterminate strings. *Journal of Discrete Algorithms* (2006)
20. Rahman, S., Iliopoulos, C., Mouchard, L.: Pattern matching in degenerate DNA/RNA sequences. In: *Proceedings of the Workshop on Algorithms and Computation (WALCOM)*. (2007) 109–120
21. Noé, L., Kucherov, G.: Improved hit criteria for DNA local alignment. *BMC Bioinformatics* **5**(149) (2004)
22. Aho, A.V., Corasick, M.J.: Efficient string matching: An aid to bibliographic search. *Communications of the ACM* **18**(6) (1975) 333–340
23. Amir, A., Porat, E., Lewenstein, M.: Approximate subset matching with don't cares. In: *Proceedings of 12th Symposium on Discrete Algorithms (SODA)*. (2001) 305–306