

# Full Compressed Affix Tree Representations

Rodrigo Cánovas, Eric Rivals

## ▶ To cite this version:

Rodrigo Cánovas, Eric Rivals. Full Compressed Affix Tree Representations. DCC: Data Compression Conference, IEEE, Apr 2017, Snowbird, UT, United States. pp.102-111, 10.1109/DCC.2017.39. lirmm-02093302

# HAL Id: lirmm-02093302 https://hal-lirmm.ccsd.cnrs.fr/lirmm-02093302

Submitted on 8 Apr 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Full Compressed Affix Tree Representations\*

Rodrigo Cánovas and Eric Rivals<sup>†</sup> L.I.R.M.M. - UMR 5506, CNRS et Université de Montpellier Montpellier, France

Institut Biologie Computationnelle, Montpellier, France emails: canovas-ba@lirmm.fr, rivals@lirmm.fr

4th April 2017

#### Abstract

The Suffix Tree, a crucial and versatile data structure for string analysis of large texts, is often used in pattern matching and in bioinformatics applications. The Affix Tree generalizes the Suffix Tree in that it supports full tree functionalities in both search directions. The bottleneck of Affix Trees is their space requirement for storing the data structure. Here, we discuss existing representations and classify them into two categories: *Synchronous* and *Asynchronous*. We design Compressed Affix Tree indexes in both categories and explored how to support all tree operations bidirectionally. This work compares alternative approaches for compress the Affix Tree, measuring their space and time trade-offs for different operations. Moreover, to our knowledge, this is the first work that compares all Compressed Affix Tree implementations, *i.e.*, four different approaches: the Affix Array, the Bidirectional Wavelet Index, and our two new structures, measures their space and time trade-offs, offering a practical benchmark for this structure.

#### **1** Introduction

Stringology offers many indexes for unidirectionally searching a pattern in a text, but recently the need for bidirectional search tools arose for instance in genomics. An example is the search of hairpin like structures that consist of a *loop motif* (a short core sequence) and a *tail* (such that on each side of the loop motif is one half of a palindromic string) [6]. In this search, when a candidate loop is found, one can extend the search with each branch of the hairpin. The search of such motifs has triggered the development of bidirectional indexes: first the Affix Tree [14, 23], then the Affix Array [24], and more recently the Bidirectional Wavelet Tree [22]. In all these approaches the procedure can switch the search direction at any time during the search, which is needed when the whole pattern is not a single word given at start time, but can shrink or extend in piecewise during the search. Also, recently Affix Trees have been mentioned as a way to represent complex graphs encoding string overlaps, for example, the Hierarchical Overlap Graph [7].

Here, we introduce a classification of bidirectional indexes based on how they are traversed, *Asynchronous* and *Synchronous*. We propose new Compressed Affix Tree indexes in both categories: the *Asynchronous* Compressed Affix Tree (ACAT) and the *Synchronous* Compressed Affix Tree (SCAT), and explain their practical implementations. Moreover, we exhibit both a reduced version of the index of [10], and a sampled version of ACAT. To evaluate the potential of these approaches and to compare them with different indexes for bidirectional search, we have implemented the Affix Array [24] (given that the original was not found), and compared seven indexing structures in practice. This yields the first large and practical overview of bidirectional search indexes in term of memory usage, and of operation times. Last, all the implementations designed for this article will be made available in near future.

<sup>\*</sup>Work supported by the Institut de Biologie Computationnelle (ANR-11-BINF-0002).

<sup>&</sup>lt;sup>†</sup>to whom correspondance should be addressed.

#### **2** Basic Concepts

Let  $\Sigma$  denote a finite alphabet of cardinality  $\sigma$  and T be a sequence of characters from  $\Sigma$ . The length of T is denoted |T| and let n := |T|. The characters of T are indexed from 1 to n, so that the *i*-th character of T is  $T_i$  or T[i]. A substring of T is denoted  $T_{i,j} = T_i T_{i+1} \dots T_j$ . A prefix of T is a substring of the form  $T_{1,j}$ , and a suffix is a substring of the form  $T_{i,n}$ .

**Suffix Tree** Let  $P = \{S^1, \dots, S^N\}$  be a set of *N* strings. A *Trie* [9, 12] is a tree where each edge is associated with an element of  $\Sigma$  and each node represents a distinct prefix in the set *P*. This structure supports searching for a string in the set in time proportional to the length of the sought string (independently of the set size).

The *Suffix Trie* of a text T is a *Trie* built over the set of all the suffixes of T (assuming that T is terminated by a special element "\$"). The *Suffix Tree* (*ST*) [16, 25] of a text T is a *Suffix Trie* where each unary path is converted into a single edge. The leaf of the Suffix Tree indicates the text position where the corresponding suffix starts. Table 1 lists some of the operations used to navigate over the nodes and leaves of the Suffix Tree and introduces a notation for these.

Operation	Description
Root()	Root of the Suffix Tree.
Ancestor(v,w)	True if <i>v</i> is an ancestor of <i>w</i> .
Depth(v) / TDepth(v)	String-depth / Tree-depth of v.
Parent(v)	Parent node of <i>v</i> .
SLink(v)	Suffix-link of <i>v</i> ; the node <i>w</i> s.th. $\pi(w) = \beta$ if $\pi(v) = a\beta$ for $a \in \Sigma$ .
WLink(v, a)	Weiner-link of <i>v</i> given <i>a</i> ; the node <i>w</i> s.th. $\pi(w) = a\beta$ if $\pi(v) = \beta$ .
LCA(v,w)	Lowest common ancestor of <i>v</i> and <i>w</i> .
Child(v,a)	The node <i>w</i> s.th. the first letter on edge $(v, w)$ is $a \in \Sigma$ .
Children(v)	The list of nodes <i>w</i> such that <i>w</i> is a child of <i>v</i> .
Degree(v)	Number of children of <i>v</i> .

Table 1: Navigational operations over the nodes and leaves of a suffix tree.

Various Compressed Suffix Tree (CST) [1, 18, 20, 21] implementations offering different space vs time trade-off are available.

**Suffix Array** A Suffix Array (*SA*) [15] is a permutation of the interval [1,n] such that  $T_{SA[i],n} <_L T_{SA[i+1],n}$  for all  $1 \le i < n$ , where  $<_L$  denotes lexicographically smaller than. Several implementations of *Compressed Suffix Arrays* have been designed [8, 17]. Each provides three main operations: 1/ given *i* extract the value of SA[i], 2/ given a position *j* in the text *T*, fetch its position in the Suffix Array ( $SA^{-1}[j]$ ), and 3/ find in SA(T) the interval of the suffixes starting with a given pattern.

The Suffix Array and the Suffix Tree of a text are connected. When the children of each node of ST(T) are ordered lexicographically by their string label, SA(T) contains the positions of the Suffix Tree leaves visited in left-to-right order. Moreover, every internal node v of the ST(T) can be represented as an interval [i, j] of SA(T) whose suffixes share a prefix corresponding to node v.

Generally the Suffix Array is equipped with the Longest Common Prefix (*LCP*) array, which allows solving many string processing problems in optimal time and space. Conceptually, the *LCP* array defines the shape of the Suffix Tree and thus allows any traversal of ST(T) to be simulated using SA(T).

**Burrows-Wheeler Transform (BWT)** For a text *T* of length *n*, the *BWT* ( $T^{bwt}$ ) [4] is the last column of an  $n \times n$  matrix, denoted  $M_T$ , whose rows are the rotations (or cyclic shifts) of *T* sorted in lexicographic order. Equivalently, the *BWT* collects the character preceding each suffix of *T* in the order of SA(T). Hence the *BWT* can be identified with *SA*. Assuming that the last symbol of *T* is lexicographically smaller than all the others and unique, it is possible to reverse the transformation: T[n-1] is located at  $T^{bwt}[1]$ , since the first element of  $M_T$  starts with T[n].



Figure 1: Example of an Affix Tree, with the Suffix Tree for FoS and text T on the left, the Suffix Tree for BAS and the reverse of T on the right. Dotted arrows represent the affix links.

The *LF*, or *Last-to-First* mapping, function allows us to move from the character  $T^{bwt}[i]$ , in the last column of  $M_T$  to its position in the first column of  $M_T$ . Thus, *LF* allows us to navigate *T* backwards,  $T[n-2] = T^{bwt}[LF(1)]$ , and for any position k,  $T[n-k] = T^{bwt}[LF^{k-1}(1)]$ . In general,  $T^{bwt}$  is stored using a *wavelet tree* [11, 17]. The  $T^{bwt}$  is mainly used for pattern search using a method called *backward search* (making use of *LF* and a sampling of the *SA*), which is implemented by the family of indexes known as FM-indexes [17].

**Affix Tree** The *Affix Tree* was introduced and tested by Stoye [23], then further studied by Maaß [14]. To support bidirectional applications, it combines the Suffix Tree of T (denoted ST) with the Suffix Tree of the reversed text of T (denoted RT). The Affix Tree distinguishes four kinds of nodes in a Suffix Tree: *explicit* nodes, which have at least two children; *implicit* nodes, which belong to an edge between two explicit nodes; *explicit leaves*, which cannot be extend any further in the tree; and *implicit leaves*, which belong to the edge linking an explicit leaf to its parent. In the Affix Tree, each internal node (explicit and implicit) of ST is associated to a node in RT, and vice versa, by an *affix link*. The construction of an Affix Tree theoretically takes linear time and space, but in practice it "required large amount of memory (approximately 45n bytes) and its complex data structure decreased the algorithmic performance" [24].

### **3** A Classification of Bi-directional Indexes

A bidirectional index in general includes a data structure to search from left to right in T, and another to search from right to left. We call them respectively the Forward Structure (FoS) and the Backward Structure (BAS) and say that each is the *companion structure* of the other. In this work we consider two classifications for storing and maintaining the current state: *Asynchronous* and *Synchronous*.

In the Asynchronous mode, the index maintains a single current state (or current node), say the current state for FoS, and when a switch of direction is needed, it computes the corresponding state in BAS, before carrying on the search in the opposite direction. To ease updating the current state, the index stores in an extra table the corresponding node in BAS for each explicit node in FoS. A symmetrical table is needed to switch from BAS to FoS. Notice that under this definition, it is only possible to which between the structures from a explicit node to an internal node. In this mode an state is defined as a quadruplet  $\langle d, l, node, k \rangle$ , where d is the direction of the data structure to which the (explicit) node with path length l belong, and k is the context length. The context is the prefix of the path from the root to node that is ignored.

In the Synchronous mode, one maintains a pair of current states: one for FoS, one for BAS. Hence, at each

modification of the matched word, both states are updated. These indexes store less information then *Asynchronous* indexes, but pay an extra time at each extension of the search. The pair of states must always be synchronized. In this case an state is defined by the triplet  $\langle node_{FOS}, node_{BAS}, l \rangle$ . Figure 2 shows an example of states for both categories.



Figure 2: Example of possible states in the *Asynchronous* and *Synchronous* modes. Note that the nodes that are associated are not necessarily the same in both cases.

Table 2 lists all approaches explored in this work.

Reference	[24]	New	New	[10]	New	[13, 22]	New
Space	$2 \cdot (SA + LCP + child\ tables + ALink)$	$2 \cdot (CST + ALink)$	$2\cdot (CST + Alink_{ ext{sampled}})$	$2 \cdot (CST + rminq + rmaxq)$	$2 \cdot (CST + rminq)$	2 · (FM-Index)	$2 \cdot (CST)$
Description	Strothmann's Affix Array	Asynchronous Affix Tree implementation	Asynchronous Affix Tree implementation	Gog et al. Affix Tree	Reduced version of ACATN	Bidirectional BWT	Synchronous Affix Tree implementation
Full Tree Operations	No	Yes	Yes	Yes	Yes	No	Yes
Category	Asynchronous	Asynchronous	Asynchronous	Asynchronous	Asynchronous	Synchronous	Synchronous
Approach	AFA	ACAT	ACATS	ACATN	RACATN	BIDWT	SCAT

Table 2: Compressed Affix Tree approaches studied in this work.

#### 4 Asynchronous Approaches

The first attempt to reduce the space required by Affix Trees was the *Affix Array* (AFA) proposed by Strothmann [24]. He replaced each Suffix Tree by a Suffix Array plus an *LCP* array. In addition, he included two *child tables* [2] to ease locating the children of any node/interval. Finally, the AFA also stores for each explicit node (of FoS and BAS) the pointer to its corresponding node in the companion structure. The arrays storing these pointers are called the *Affix Link* arrays (*ALink*). The AFA allows both fast access to the child of a node, and fast switching between FoS and BAS. Also it supports the operation *suffix-children* (*prefix-children*), which returns all possible nodes that are an extension of a node by one character to the right (left).

To compute the *ALink* arrays, Strothmann's algorithm precomputes for FoS the *SLink* (see Table 1) of all explicit nodes. Then, it traverses the explicit nodes in depth-first preorder, checks whether the *ALink* value of the current node is already stored, otherwise computes it. For a node v whose *ALink* value is unknown, it uses the formula:  $ALink(v) := Child_{BAS}(ALink(SLink(v)), c)$ , where c is the first character on the string path of v. Computing the *ALink* array of BAS is symmetrical.

We notice that Strothmann's algorithm for compute the *ALink* for a node v can take many recursive calls (in the worse case, a number equal to the length of the string path of v). Instead, we propose to visit the nodes in level order, such that the computation of the *ALink* for a node v takes a number of recursive steps that is bounded by the length of the edge linking v to its parent. We experimentally observed that Strothmann's algorithm exhausted the available memory for a 50 MB text. So, in this work we compute the *ALink* array using our proposed methodology.

Strothmann also proved that for every state  $\langle 0, |s|, v, 0 \rangle$ , its reverse associated state is given by  $\langle 1, |v|, ALink(v), k = |v| - |s| \rangle$ . Note that when |s| < |v| (*i.e.*, *s* is an implicit node), the path of *s* is a prefix of the path of *v*. Thus, the path of the reverse of *s* is a suffix of the path of *ALink*(*v*) starting at position *k* (see Figure 2).

While the AFA improves on the Affix Tree regarding space, it still requires a huge amount of memory and supports only operations that can be solved using the *Child* operation. Since Strothmann's approach, several compressed representations of the Suffix Tree [1, 18, 20, 21] have been reported. We propose to upgrade the AFA by keeping the *ALink* arrays but representing FoS and BAS using *CST*s. This approach, denoted as ACAT (*Asynchronous* Compressed Affix Tree), uses less memory than AFA and supports all tree operations.

Moreover, given the *ALink* formula, it is clear that we don't need to store all the *ALink* values (provided that FOS and BAS support *SLink*). According to a user defined parameter *z*, indicating the maximum distance to a sample node in number of suffix links, we store whether or not the *ALink* value of a node is sampled. To check if we need to sample the *ALink* value of a node, we just need to examine that none of the nodes traversed by doing the operation *SLink z* times is already sampled. Otherwise the node is not sampled. We name this approach the ACATS (*Asynchronous* Compressed Affix Tree Sampled). It reduces the space used by the Affix Link arrays, at the cost of longer time for switching between the structures. In the worse case, for each switch at a unsampled node, ACATS would take additional *z SLink* and *Child* operations to get the *ALink* value of that node.

The lightest scenario for ACATS is when we do not store the *ALink* arrays at all. Then, each time we want to compute the corresponding state in the companion structure, we compute it starting from the root. Albrecht and Heun [3] studied how to optimally obtain the *ALink* in this case using binary search. A faster approach was presented by Gog *et al.*[10]. As the procedure is symmetrical for a node BAS or of FOS, we explain it for the latter. They noted that if an internal explicit node *v* is represented as a *SA* interval [l, r], then all the text positions *SA*(*k*) (with  $l \le k \le r$ ) where the string path of the current node occurs are known. This also implies that it is possible to compute all the positions where the reverse string path occurs in the reverse of *T* as  $w_k = n - (SA(k) + Depth(v))$ . With this on mind, Gog *et al.*proposed first to compute the array  $A[i] := SA_{BAS}^{-1}(n - SA(i))$  for  $1 \le i \le n$ , which converts an entry of the *SA* of FOS into the matching entry in the *SA* of BAS. Then to find the reverse representation of *v*, it is enough to have an structure to compute the minimum and another to compute maximum values of A[k] for  $l \le k \le r$ , which points to the leafs in BAS associated with these values, and to return their lower common ancestor (*LCA* in Table 1). We refer to this approach as ACATN. For more details of how the array *A* is stored and queried, we refer the reader to [10].

Finally we notice that the space used by ACATN could be further reduced. When we switch from one structure to another, it is enough to know only one extremity of the *SA* interval in the reverse structure, given that the lengths of both intervals are equal [22, 24]. This simple fact implies that one of two structures is superfluous (say that for compute the maximum), and that it reduces the number of operations performed during a switch. We denote this reduced approach as RACATN.

#### 5 Synchronous Approaches

The *Synchronous* scheme aims at getting rid of Affix Link arrays by keeping FoS and BAS continuously synchronized. This necessarily increases the search time. For example, if in FoS one extends the current string path by one symbol to the right, then one also needs to extend the string path in BAS by the same symbol to the left.

To improve over the AFA in terms of space and still support bidirectional search, two groups have proposed in parallel to use an Affix Array like structure based on the BWT [13, 22]. This approach is called the *Bidirectional BWT* (BIDWT for short). The underlying idea is to store FoS and BAS using data structures that support backward search based on a wavelet tree (FM-Indexes). Then backward and forward extensions of a search are both possible. Indeed, to extend a current node to the right with character c, one performs a backward search of c from the current state in BAS. For more details, we refer the reader to [22].

While the BIDWT requires much less space than the Affix Array and offers fast search extension by one character, it does not support any tree operation which allow to shrink the current pattern. Moreover, the BIDWT only offers search extension for a single character. This means that to perform an operation like *suffix-children*, one needs to explore the possible extensions for all letters of the alphabet. A similar issue occurs for the *Child* operation, whenever the child path is longer than one symbol.

We propose to upgrade the BIDWT by storing FOS and BAS as *CST* instead of FM-Indexes. We call this approach the *Synchronous* Compressed Affix Tree (SCAT). The SCAT can support all tree operations within each structure, but the issue of maintaining both structures synchronized remains. Fortunately, we noticed that each tree operation has an equivalent operation in the reverse tree. That is, *Parent-Slink* and *Child-WLink*, and all other operations can be solved using these four operations. Let us give the algorithms for computing *Parent* as an illustration.

Let  $(node_{FOS}, node_{BAS}, l)$  be the current state. To compute the *Parent* in FOS, we return

$$\langle Parent(node_{FOS}), SLink^p(node_{BAS}), l-p \rangle$$
,

where p is the length of the edge between the *node*<sub>FOS</sub> and its parent.

### **6** Results

In order to compare the approaches presented in this work, we used four different input files obtained from the pizzachilli text collection<sup>1</sup>: DNA, PROTEINS, XML, and ENGLISH. The alphabet of size of the text tested are 16, 27, 97, and 239 respectively, and we use the 50 Megabytes version of each of them (given the high memory requirement used by AFA). The pizzachilli web-page provides a detailed descriptions of the input files. All experiments were performed on a computer with Intel(R) Xeon(R) CPU E5-2623 v3 up to 3.00 GHz, The operating system was Ubuntu 14.04.01, version 3.19.0-59-generic-generic Linux kernel. Finally all the approaches described here were implemented in C++11, using version 4.8.4 of the g++ compiler.

Each implementation uses one specific type of data structure to represent FoS and BAS (but the same for both). We compared the following methodologies: AFA (which we implemented following Strothmann's description [24], since we did not find an available implementation); BIDWT (a wrap implementation using the *sdsl library* [11], which includes the bidirectional backward and forward search presented in [22]); ACAT; ACATS (using as sample parameter z = 1); ACATN (based on the work of Gog et al. [10]); RACATN (our reduced version of ACATN); and SCAT.

In this work all the approaches that used a *CST* were implemented using the *sct3* version from the sdsl library [18], which offers a good space/time trade-off in practice. Even if other *CST* implementations could offer different trade-offs, or reduce the number of operations to run, we assume that the *sct3* implementation gives a fair methodology to compare all studied approaches.

For space limitations, we show the results of each experiment only over the DNA and ENGLISH texts, which are the ones with smallest and largest alphabet. Also in all figures of this section, the horizontal and vertical axes are in logarithmic scale.

http://pizzachili.dcc.uchile.cl/texts.html



Figure 3: Construction time versus storage space.

#### 6.1 Space usage and construction time

Figure 3 shows the space used by each structure (measured in bytes per character of the original text) vs its construction time (in milliseconds). We observe that our AFA implementation takes over 20 bpc, which is consistent with the figures in [24]. In comparison the ACAT takes 10 bpc (less than half the AFA), which makes it scalable enough to index a Human genome (whose length is  $\sim 3$  Gb). Clearly, the *Asynchronous* indexes takes longer time and more space than *Synchronous* indexes, but most of their construction time is spent for computing the Affix Link arrays (which use  $\sim 6.5$  bpc when is completely stored) or some equivalent (the array A for ACATN and RACATN). The SCAT structure uses almost twice the space of the BIDWT, auguring well of their scalability. This comparatively higher memory usage is the price to pay for supporting all tree operations, and improving the versatility.

#### 6.2 Performance of navigational operations

We tested four different operations (chosen to be supported by all indexes): *Forward-Backward* search, *Degree*, *suffix-children* and *prefix-children*. Among the operations that are supported by full Affix Trees (*i.e.*, not by AFA nor BIDWT), we only show experiments results for *SLink* due to space limitations (Note that the performace of tree operations can change depending on the *CST* structure used). For the *Forward-Backward* search, we took samples of 1000 random sequences of length 100 from the text and searched them performing forward and backward searches alternatively (thereby switching the direction after each symbol), and repeating the experiment with 1000 samples. For the operations *suffix-children*, *prefix-children*, and *SLink* we performed each on random samples of 10000 internal nodes. In all experiments, we report the average time taken to perform one operation.

Figure 4 shows the space used by the indexes (in bpc) versus the time (in microseconds) obtained for *Forward-Backward* search. The only competitive solutions against BIDWT, are the AFA and SCAT approaches. While AFA offers a faster performance but required huge amounts of space, SCAT, which is almost 10 times slower, supports all tree operations. Note that the solutions containing a representation of the *CST* could instead incorporate a FM-index (like the one used by the BIDWT). In this case, it would be possible to re-implement the operation *Forward-Backward* search to obtain the same performance than BIDWT.

Figure 5 and 6 show the results obtained for *suffix-children* and *prefix-children* respectively. As for the *Forward-Backward* search, AFA approach displayed the best time performance, since almost all required information by these queries is explicitly stored. We also notice that with BIDWT the query time for these two operations depends on the alphabet size of the input text, while SCAT depends on the *CST* structure chosen. The differences in time for the these operations between the *Asynchronous* approaches are due to the way that each index switches between FoS and BAS. So comparatively, a direct computation of *Children* is an advantage only with larger alphabets, when checking all possible symbols becomes expensive for the BIDWT.



Figure 4: Forward-Backward search time vs space usage.



Figure 5: Time of operation *suffix-children* vs space usage.

Figure 7 shows the results obtained for the operation *Degree*. To solve *Degree*, BIDWT takes similar times than for *suffix-children*, while the time of other indexes depends mostly on how the Compressed Suffix Trees are stored. This gives a clear advantage ( $\approx 1$  order of magnitude) to all approaches over the BIDWT, where *Degree* must be computed by brute force.

Figure 8 shows the average time to perform the *SLink* operation (which is not supported by BIDWT and AFA). While the *Asynchronous* approaches needs only one operation to return the *SLink*, SCAT always had to perform two operations, which explains the longer average times obtained with SCAT than with the other methodologies. We found this behavior repeated when tested other tree operations.

Summarizing, while BIDWT uses less space and offers faster backward/forward search of extension of one character, it doesn't supports all the tree operations and is sensitive to the alphabet size when computing operations related with child. The AFA approach offers generally fastest operation times, but requires too much space, and also miss some tree operations. When full affix tree operations are required, SCAT uses less space, while being competitive with the *Asynchronous* approaches. We prove that RACATN improve both in space and time over ACATN. We show that ACATS with a sample parameter of 1 (z := 1) uses  $\approx 70\%$  of the space of ACAT, while offering similar time performances. Moreover, although increasing z would reduce the space of ACATS, SCAT gives a lower bound for space and achieves almost equivalent running times (compared to ACATS with z = 1).



Figure 6: Time of operation *prefix-children* vs space usage.



Figure 7: Time of operation Degree vs space usage.

### 7 Conclusions

We investigated the design of compressed versions of Affix Trees proposing several new structures and introduced a way to classify these approaches (*Synchronous* or *Asynchronous*), offering a first benchmark of Affix Tree compressed data structures.

As future work, it will be interesting to see how the structures presented in this paper can be used for pattern search with errors, where it is necessary to be able to extend and shrink the pattern during the search. Also, new approaches for bidirectional search will keep emerging (for example [19]), being of great interest to explore how these new solutions would compare with the ones presented in this work.

#### **Remark and acknowledgements**

This report is an extended version of the article published in the *Data Compression Conference* in 2017 [5]. We thank the reviewers for insightful comments and suggestions. This work is supported by the Institut de Biologie Computationnelle (ANR-11-BINF-0002).



Figure 8: Time of operation Slink vs space usage.

## References

- [1] A. Abeliuk, R. Cánovas, and G. Navarro. Practical Compressed Suffix Trees. *Algorithms*, 6(2):319–351, 2013.
- [2] M. I. Abouelhoda, S. Kurtz, and E. Ohlebusch. Replacing suffix trees with enhanced suffix arrays. *Journal of Discrete Algorithms*, 2(1):53–86, 2004.
- [3] B. Albrecht and V. Heun. Space Efficient Modifications to Structator A Fast Index-Based Search Tool for RNA Sequence-Structure Patterns. In 11th International Symposium on Experimental Algorithms, (SEA), pages 27–38, 2012.
- [4] M. Burrows and D. Wheeler. A block-sorting lossless data compression algorithm. Technical report, Digital Equipment Corporation, 1994.
- [5] R. Cánovas and E. Rivals. Full compressed affix tree representations. In A. Bilgin, M. W. Marcellin, J. Serra-Sagristà, and J. A. Storer, editors, 2017 Data Compression Conference, DCC 2017, Snowbird, UT, USA, April 4-7, 2017, pages 102–111. IEEE, 2017.
- [6] Y. Carafa, E. Brody, and C. Thermes. Prediction of Rho-independent *Escherichia coli* Transcription Terminators. *Journal of Molecular Biology*, 216(4):835–858, 1990.
- [7] B. Cazaux, R. Cánovas, and E. Rivals. Shortest DNA cyclic cover in compressed space. In *Data Compression Conference DCC*, pages 536–545, 2016.
- [8] P. Ferragina, R. González, G. Navarro, and R. Venturini. Compressed Text Indexes: From Theory to Practice. ACM Journal of Experimental Algorithmics (JEA), 13:article 12, 2009.
- [9] E. Fredkin. Trie memory. In Communications of the ACM 3, pages 490–500, 1960.
- [10] S. Gog, K. Karhu, J. Kärkkäinen, V. Mäkinen, and N. Välimäki. Multi-pattern matching with bidirectional indexes. *Journal of Discrete Algorithms*, 24:26–39, 2014.
- [11] R. Grossi, A. Gupta, and J. Vitter. High-order entropy-compressed text indexes. In Proc. 14th Symposium on Discrete Algorithms (SODA), pages 841–850, 2003.
- [12] D. Knuth. The Art of Computer Programming, Volume 3: Sorting and Searching. Addison-Wesley, 1973.
- [13] T. W. Lam, R. Li, A. Tam, S. Wong, E. Wu, and S. M. Yiu. High Throughput Short Read Alignment via Bi-directional BWT. In *Bioinformatics and Biomedicine*, pages 31–36, 2009.

- [14] M. G. Maaß. Linear Bidirectional On-Line Construction of Affix Trees. Algorithmica, 37(1):43–74, 2003.
- [15] U. Manber and E. Myers. Suffix arrays: a new method for on-line string searches. SIAM Journal on Computing, pages 935–948, 1993.
- [16] E. McCreight. A space-economical suffix tree construction algorithm. *Journal of the ACM*, 32(2):262–272, 1976.
- [17] G. Navarro and V. Mäkinen. Compressed Full-Text Indexes. ACM Computing Surveys, 39(1):article 2, 2007.
- [18] E. Ohlebusch, J. Fischer, and S. Gog. CST++. In Proc. 17th International Symposium, SPIRE 2010, pages 322–333, 2010.
- [19] C. Pockrandt, M. Ehrhardt, and K. Reinert. Constant-time and space-efficient unidirectional and bidirectional FM-indices using EPR-dictionaries. *ArXiv e-prints*, 2016.
- [20] L. Russo, G. Navarro, and A. Oliveira. Fully-Compressed Suffix Trees. ACM Transactions on Algorithms (TALG), 7(4):article 53, 2011. 35 pages.
- [21] K. Sadakane. Compressed suffix trees with full functionality. *Theory of Computing Systems*, 41(4):589–607, 2007.
- [22] T. Schnattinger, E. Ohlebusch, and S. Gog. Bidirectional search in a string with wavelet trees and bidirectional matching statistics. *Information and Computation*, 213:13–22, 2012.
- [23] J. Stoye. Affix Trees. Technical Report, University of Bielefeld, 2000.
- [24] D. Strothmann. The affix array data structure and its applications to RNA secondary structure analysis. *Theoret-ical Computer Science*, 389(1–2):278–294, 2007.
- [25] P. Weiner. Linear pattern matching algorithms. In Proc. 14th Annual Symposium on Switching and Automata Theory, pages 1–11, 1973.