# Research Report: A Theoretical and Experimental Study on the Construction of Suffix Arrays in External Memory

Andreas Crauser and Paolo Ferragina

**Author's Address**

Andreas Crauser
Max-Planck Institut für Informatik
Im Stadtwald
D-66123 Saarbrücken
crauser@mpi-sb.mpg.de




Paolo Ferragina
Dipartimento di Informatica
Università di Pisa
Corso Italia 40
56125 Pisa, Italy
ferragin@di.unipi.it

**Abstract**

The construction of full-text indexes on very large text collections is nowadays a hot problem. The suffix array [Manber-Myers, 1993] is one of the most attractive full-text indexing data structures due to its simplicity, space efficiency and powerful/fast search operations supported. In this paper we analyze, both theoretically and experimentally, the I/O-complexity and the working space of six algorithms for constructing large suffix arrays. Three of them are the state-of-the-art, the other three algorithms are our new proposals. We perform a set of experiments based on three different data sets (English texts, Amino-acid sequences and random texts) and give a precise hierarchy of these algorithms according to their working-space vs. construction-time tradeoff. Given the current trends in model design [19, 47] and disk technology [16, 40], we will pose particular attention to differentiate between "random" and "contiguous" disk accesses, in order to reasonably explain some practical I/O-phenomena which are related to the experimental behavior of these algorithms and that would be otherwise meaningless in the light of other simpler external-memory models.

At the best of our knowledge, this is the first study which provides a wide spectrum of possible approaches to the construction of suffix arrays in external memory, and thus it should be helpful to anyone who is interested in building full-text indexes on very large text collections.

Finally, we conclude our paper by addressing two other issues. The former concerns with the problem of building word-indexes; we show that our results can be successfully applied to this case too, without any loss in efficiency and without compromising the simplicity of programming so to achieve a uniform, simple and efficient approach to both the two indexing models. The latter issue is related to the intriguing and apparently counter-intuitive "contradiction" between the effective practical performance of the well-known BaezaYates-Gonnet-Snider's algorithm [24], verified in our experiments, and its unappealing (i.e., cubic) worst-case behavior. We devise a new external-memory algorithm that follows the basic philosophy underlying that algorithm but in a significantly different manner, thus resulting in a novel approach which combines good worst-case bounds with efficient practical performance.

# Contents

# 1 Introduction

In the information age, one of the fastest growing category of databases are the textual databases—like AP-news, Digital Libraries, Genome databases, book collections [22]. Their ultimate impact heavily depends on the ability to *store* and *search* efficiently the information contained into them. The continued decline in the cost of external storage devices (like disks and CD-ROMs) has made nowadays the storage issue not much of a problem, compared to the challenges posed by the fast retrieval of the user-requested informations. In order to achieve this goal, specialized indexing data structures and searching tools have been introduced so far. Their main idea is to build an *index* that allows to focus the search for a given pattern string only on a very small portion of the text collection. The improvement in the query-performance is paid by the additional space necessary to store the index. Most of the research in this field has been therefore directed to design indexing data structures which offer a good trade-off between the query time and the space usage. The two main approaches are: *word* indexes and *full-text* indexes.

Word-indexes exploit the fact that for natural linguistic texts, the universe of distinct words is small. They store all the occurrences of each word in a table that is indexed via a hashing function or a tree structure (see e.g. [42]). To reduce the size of the table, common words are often not indexed (e.g. the, at, a). The advantage is to support very fast word (or prefix-word) queries, while the two obvious weaknesses consist of the limited set of supported queries, which makes the search for phrases or complex regular expressions complicated and inefficient [24, 32], and the impossibility to index *unstructured texts*, like DNA-sequences or Chinese texts [20]. (For alternative approaches to word-indexes see [17, 51].)

Full-text indexes have been designed to overcome the limitations above by dealing with arbitrary (unstructured) texts and general queries, at the cost of an increase in the additional space occupied by the underlying indexing data structure. Examples of such indexes are: suffix trees [7, 33, 13], suffix arrays [32] (cfr. PAT-arrays [24]), PAT-trees [24], Suffix Cactus [26] and String B-trees [21]. They have been successfully applied to fundamental string-matching problems (see e.g. [7]) as well text compression [11, 31], analysis of genetic sequences [25, 34, 10, 9] and recently to the indexing of special linguistic texts [20]. General full-text indexes are therefore the natural choice to perform fast complex searches without any restrictions. The most important complexity measures for evaluating their efficiency are [8, 48]: (i) the time and the extra space required to build the index, (ii) the time required to search for a string, and (iii) the space used to store the index. Points (ii) and (iii) have been largely studied in the scientific literature [13, 21, 24, 32, 33]. In this paper, we will investigate Point (i) by addressing the efficient *construction* of full-text indexes on very large text collections. This is a hot topic nowadays [1] because the construction phase may be a bottleneck that can even prevent these indexing tools to be used in large-scale applications. In fact, known construction algorithms are very fast when employed on textual data that fit in the internal memory of computers [6, 32, 41] but their performance immediately degenerates when the text size becomes so large that the texts must be arranged on (slow) external storage devices [13, 21].

---

[1]Zobel *et al.* [48] say that: *"We have seen many papers in which the index simply 'is', without discussion of how it was created. But for an indexing scheme to be useful it must be possible for the index to be constructed in a reasonable amount of time, ....."*.

## 1.1 (Mechanical) Disk vs. (Electronic) Memories.

Although disks provide a large amount of space at low cost, their access time is from $10^5$ to $10^6$ times slower than the time to access the internal memory of computers [39]. Unfortunately, the majority of the known algorithms ignore this fact and use disks in the same manner as if all data were fitted into the internal memory and could be accessible in *unit time*. Consequently, they suffer from the so called *I/O bottleneck*: They spend most of the time in moving data to/from the disk. Since our aim is to study the efficiency of algorithms that operate on very large text collections, we need a computational model that captures this *I/O-phenomenon*. Accurate disk models are complex [40, 43, 44, 47], and it is virtually impossible to exploit all the fine points of disk characteristics systematically, either in practice or for theoretical analysis. In order to capture in an easy, yet significant, way the differences between the internal (electronic) memory and the external (mechanical) disk, we refer to the external-memory model introduced by Vitter and Shriver [44]. Here a computer is abstracted to consist of a *two-level memory*: a fast and small internal memory, of size $M$, and a slow and arbitrarily large external memory, called *disk*. Data between the internal memory and the disk are transfered in blocks of size $B$ (called *disk pages*). Since disk accesses are the dominating factor in the running time of many algorithms the usual accounting scheme [44, 47] consists of evaluating their asymptotic performance by counting the total number of disk accesses performed during the computation. Although this is a workable approximation for algorithm design, this accounting scheme does not accurately predict the running time of algorithms on real machines because it does not take into account some important specialties of new disk systems [16]. In fact, disk access costs have mainly two components: the time to fetch the first bit of requested data (seek time) and the time required to transmit the requested data (transfer rate). Transfer rates are more or less stable but seek times are highly variable [16, 40]. Hence, it is very well known [16, 40, 48] that accessing one page from the disk in most cases decreases the cost of accessing the page succeeding it, so that "bulk" I/Os are less expensive per page than "random" I/Os. This difference becomes much more prominent if we also consider the reading-ahead/buffering/caching optimizations which are common in current disks and operating systems. To deal with these specialties and avoid the introduction of new parameters, we adopt the simple, yet significant, accounting scheme introduced by Farach *et al.* [19]: Let $c$ be a constant, a *Bulk I/O* is the reading/writing of a contiguous sequence of $cM/B$ disk pages; a *random* I/O is any single disk-page access which is not part of a bulk I/O. The performance of the external-memory algorithms is therefore evaluated by measuring: (a) the number of I/Os (bulk and random), (b) the internal running time (CPU time), and (c) the number of disk pages used during the construction process (working space).

## 1.2 Previous Work.

For simplicity of exposition, we will use $N$ to denote the size of the whole text collection and we will assume throughout the paper that the index is built on *only one* text, obtained by concatenating all the available texts separated by proper special characters (i.e., endmarkers).

The most famous full-text indexing data structure is the *suffix tree* [33, 49]. Its numerous applications have been cataloged in [7]. In internal memory, a suffix tree can be constructed in $O(N)$ time [33, 18]; in external memory, Farach *et al.* [19] showed that a suffix tree can be optimally constructed within the same I/O-bound as sorting $N$ atomic items; nonetheless, known practical construction algorithms [13] for external memory still operate in a brute-force manner requiring $\Theta(N^2)$ total I/Os in the worst-case. The main limit of these algorithms is inherent in the working space which depends on the text structure, is not predictable in advance and turns out to require between $16N$ and $26N$ bytes [2] (assuming that the $N \leq 2^{32}$ [30, 32]). This makes them impractical even for moderately large text collections (consider what happens for $N \approx 100$ Mbytes, the suffix tree would occupy 1.7Gbytes !). Searching for an arbitrary

---

[2][30] was able to reduce the working space to $10.1N - 20N$. However, this was achieved by assuming that $N < 2^{27}$.

string of length $p$ takes $O(p)$ time in internal memory (which is optimal for bounded alphabets), but it does not gain any speed up from the block-transfer when the suffix tree is stored on the disk [21].

To circumvent these drawbacks, Ferragina and Grossi introduced the *String B-tree* data structure [21]. Searching for an arbitrary pattern string of length $p$ takes $O(p/B + \log_B N)$ random I/Os (which is optimal for unbounded alphabets). The total occupied space is asymptotically optimal, and needs $12.3N$ bytes. The String B-tree is a *dynamic* data structure which supports efficient update operations, but its construction from scratch on a text collection of size $N$ takes $O(N \log_B N)$ random I/Os. Hence, space and construction time may still be a bottleneck in large-scale applications.

Since the space occupancy is a crucial issue when building and using full-text indexes on very-large text collections, Manber and Myers [32] proposed the *suffix array* data structure, which consists of an array of pointers to text positions and thus occupies overall $4N$ bytes [3] (thus being 4 times smaller than a suffix tree, and 3 times smaller than a String B-tree). Suffix arrays can be efficiently constructed in $O(N \log_2 N)$ time [32] and $O((N/B)(\log_2 N) \log_{M/B}(N/B))$ random I/Os [1]. In external memory, searching is not as fast as in String B-trees but it still achieves good I/O-performances. Suffix arrays have been recently the subject of experimental investigations in internal memory [32, 41], external memory [24] and distributed memory systems [28, 36]. The motivation has to be probably found in their simplicity, reduced space occupancy, and in the small constants hidden in the big-Oh notation which make them suitable for indexing very-large text collections in practical applications (without any surprise in the final performances !). Suffix arrays also present some natural advantages over the other indexing data structures for what concerns the construction phase. Indeed, their simple topology (i.e., an array of pointers) avoids at construction time all the problems related to the efficient management of tree-based data structures (like suffix trees and String B-trees) on external storage devices [29]. Additionally and more importantly, efficient practical procedures for building suffix arrays are definitively useful for efficiently constructing suffix trees, String B-trees and other powerful indexing data structures [35], so that they can allow to overcome their main bottleneck (i.e., expensive construction phase).

## 1.3  Our Contribution.

With the exception of some preliminary and partial experimental works [32, 24, 36], to the best of our knowledge, no full-range comparison exists among the known algorithms for building large suffix arrays. This will be the main goal of our paper, where we will theoretically analyze and experimentally study various suffix-array construction algorithms. Some of them are the state-of-the-art in the practical setting [24], others are the most efficient theoretical ones [32, 1], whereas *three* other algorithms are our *new proposals* obtained either as slight variations of the previous ones or as a careful combination of known techniques which were previously employed only in the theoretical setting. We will study these algorithms by evaluating their working space and their construction complexity both in terms of number of (random and bulk) I/Os and CPU-time. In the design of the new algorithms we will address mainly two issues: (i) simple algorithmic structure, and (ii) reduced working space. The first issue has clearly an impact on the predictability and practical efficiency of the proposed algorithms, which are also *flexible* enough to be used in distributed memory systems. In fact, since our new algorithms will be based on two basic routines—*sorting and scanning of a set of items*—they will immediately provide us with very fast suffix-array construction algorithms also for $D$-disk arrays systems [12] (thus achieving a speedup factor of approximately $D$ [38, 44]) and clusters of $P$ workstations (thus achieving a speedup factor of approximately $P$ [23]). Additionally, our algorithms will be not faced with the problems of carefully setting some system parameters, as it happens in the results of [28, 36]. The second issue (i.e., space usage) will be also carefully taken into account because the real disk size is limited and thus a large working space could prevent the use of a construction algorithm even for moderately large text

---

[3]We assume that $N \leq 2^{32}$.

collections. Additionally, as Knuth [29][Sect. 6.5] says: *"space optimization is closely related to time optimization in a disk memory"*. Our aim will be therefore to keep the working space of our algorithms as small as possible without worsening their total running time.

We will discuss all the algorithms according to these two resources and we will pose particular attention to differentiate between random and bulk I/Os. This will allow us to take into account in the asymptotic analysis of the most significant disk characteristics, thus making reasonable predictions on the practical behavior of these algorithms. To validate our conjectures, we will perform an extensive set of experiments based on three text collections — English texts, Amino-acid sequences and random data. The various parameters accounted for in our theoretical analysis, will allow us to reasonably explain some *interesting I/O-phenomena* which arise during these experiments and which could not be explained in the light of other simpler external-memory models (see e.g. [47] for a survey). As a result of the theoretical and experimental analysis, we will give a precise hierarchy of suffix-array construction algorithms according to their working-space vs. construction-time tradeoff; thus providing a wide spectrum of possible approaches for anyone who is interested in building large full-text indexes.

This analysis leads us to additionally address two further issues: construction of word-indexes and worst-case performance of the BaezaYates-Gonnet-Snider's algorithm [24], one of the most effective algorithms in our proposed hierarchy. As far as the former issue is concerned, we show that our new algorithms can be successfully applied to construct word-indexes without any loss in efficiency and without compromising the ease of programming so to achieve a uniform, simple and efficient approach to both the two indexing models. The latter issue deserves much attention and is related to the intriguing, and apparently counterintuitive, "contradiction" between the effective practical performance of the BaezaYates-Gonnet-Snider's (BGS) algorithm [24], verified in our experiments, and its unappealing (i.e., cubic) worst-case behavior. This fact motivated us to deeply study its algorithmic structure and exploit more fine-grained external-memory models [19] for explaining its experimental performances. This has finally lead us to devise a *new* external-memory construction algorithm that follows the BGS's basic philosophy but in a significantly different manner, thus resulting in a novel approach which combines good practical qualities with efficient worst-case performances.

## 1.4   Map of the Paper.

In the next section, we give some basic definitions and fix our notation. In Section 2.1, we review three well-known algorithms for constructing suffix arrays and analyze their space requirements. The description may appear in some sense much detailed, but our aim is to make this paper as a *self-contained reference* to anyone who is interested in building large suffix arrays. In Section 2.2, we describe three new external-memory algorithms for constructing suffix arrays on large text collections. In Section 3, we first overview the main features of the external-memory library LEDA-SM [15], used to implement all tested algorithms; and then we introduce our benchmark suite and discuss our experimental settings. In Section 4, we present and discuss the experimental results on the three data sets: Reuters corpus, Amino-acid collection and random texts. In Section 5.1, we address the problem of constructing word-indexes and show how our results can be easily extended to this indexing model in a natural way. In Section 5.2, we describe a new algorithm for suffix-array construction which follows the basic philosophy of BaezaYates-Gonnet-Snider's algorithm [24], one of the most interesting algorithms in our experiments, but in a significantly different manner thus resulting effective both in theory and in practice. We conclude the paper describing some open problems and discussing future directions of research.

# 2 The Suffix Array data structure

Given a string $T[1, N]$, let $T[1, i]$ denote the $i$-th prefix of $T$ and $T[i, N]$ denote the $i$-th suffix of $T$. The symbol $\leq_L$ denotes the *lexicographic order* between any two strings and the symbol $\leq_i$ denotes the lexicographic order between their length-$i$ prefixes: $S \leq_i T$ if and only if $S[1, i] \leq_L T[1, i]$.

The suffix array $SA$ built on the text $T[1, N]$ is an array containing the lexicographically ordered sequence of suffixes of $T$, represented via pointers to their starting positions (i.e., *integers*). For instance, if $T = ababc$ then $SA = [1, 3, 2, 4, 5]$. This way, $SA$ occupies $4N$ bytes if $N \leq 2^{32}$. Manber and Myers [32] introduced this data structure in the early 90s (cfr. PAT-array [24]) and proposed an interesting algorithm to efficiently search for an arbitrary string $P[1, p]$ in $T$ by taking advantage of the information coded into $SA$. They proved that all the $Pocc$ occurrences of $P$ in $T$ can be retrieved in $O(p \log_2 N + Pocc)$ time in the worst-case using the plain $SA$, and that this bound can be improved to $O(p + \log_2 N + pocc)$ time if another array of size $4N$ is provided. If $SA$ is stored on the disk, divided into blocks of size $B$, the search for $P$ takes $O(\lceil p/B \rceil \log_2 N + Pocc/B)$ random I/Os. In practical cases it is $p << B$, so that $\lceil p/B \rceil = 1$. The simplicity of the search procedure, the small constants (hidden in the big-Oh notation), and the reduced space occupancy are the most important characteristics that make this data structure very appealing for practical applications.

## 2.1 Constructing a suffix array

We now review three known algorithms for constructing the suffix array data structure on a string $T[1, N]$. We analyze their construction time in terms of CPU-time and number of I/Os (both random and bulk) in the external-memory model. We also address the issues related to their space requirements, by assuming $N \leq 2^{32}$ so that 4 bytes are sufficient to encode a (suffix) pointer. We remark that the working space of all algorithms is *linear* in the length $N$ of the indexed text, and thus *asymptotically optimal*. However, since the constants hidden in the big-Oh notation differ a lot and the available disk space is not unlimited, we will carefully evaluate the space usage of these algorithms in order to study their practical applicability (see Table 2.1 for a summary).

### 2.1.1 The algorithm of Manber and Myers

It is the fastest theoretically known algorithm for constructing a suffix array in internal memory [32]. It requires $O(N \log_2 N)$ worst-case time and consists of $\lceil \log_2 (N + 1) \rceil$ stages, each taking $O(N)$ time. In the first stage, the suffixes are put into buckets according to their first symbol (via radix sort). Before the generic $h$-th stage starts, the algorithm has inductively identified a sequence of buckets, each containing a set of suffixes. Any two suffixes in the same bucket share the first $2^{h-1}$ characters, whereas any two suffixes in two different buckets are $\leq_L$-sorted according to the bucket-ordering (initially, we have just one bucket containing all of $T$'s suffixes). In stage $h$, the algorithm sorts lexicographically the suffixes of each bucket according to their first $2^h$ characters, thus forming new smaller buckets which preserve the inductive hypothesis. After the last stage, all the buckets will contain one suffix, thus giving the final

| Algorithm | Working space | CPU–time | total number of I/Os |
|---|---|---|---|
| Manber–Myers (Sect. 2.1.1) | $8N$ | $N \log_2 N$ | $N \log_2 N$ |
| BaezaYates–Gonnet–Snider (Sect. 2.1.2) | $8N$ | $(N^3 \log_2 M)/M$ | $(N^3 \log_2 M)/(MB)$ |
| Doubling (Sect. 2.1.3) | $24N$ | $N(\log_2 N)^2$ | $(N/B)\,(\log_{M/B} N/B)\,\log_2 N$ |
| Doubling+Discard (Sect. 2.2.1) | $24N$ | $N(\log_2 N)^2$ | $(N/B)\,(\log_{M/B} N/B)\,\log_2 N$ |
| Doubling+Discard+Radix (Sect. 2.2.1) | $12N$ | $N(\log_2 N)^2$ | $(N/B)\,(\log_{M/(B \log N)} N)\,\log_2 N$ |
| Construction in $L$ pieces (Sect. 2.2.3) | $\max\{\frac{24N}{L},\frac{2NL+8N}{L}\}$ | $N(\log_2 N)^2$ | $(N/B)\,(\log_{M/B}(N/B))\,\log_2 N$ |
| New BGS (Sect. 5.2) | $8N$ | $N^2(\log_2 M)/M$ | $(N^2/MB)$ |

Table 2.1: *The CPU-time and the number of I/Os are expressed in big-Oh notation; the working space is evaluated exactly; L is an integer constant greater than 1. BaezaYates-Gonnet-Snider algorithm (BGS), and its new variant (called* new BGS*), operate via only disk scans, whereas all the other algorithms mainly execute random I/Os. Notice that with a tricky implementation, the working space of the BGS-algorithm can be reduced to 4N.*

$SA$. The efficiency of this algorithm depends on how fast is the sorting step in a generic stage. Manber and Myers [32] showed how to perform it in linear time by using only two integer arrays, for a total of $8N$ bytes. If this algorithm is used in a virtual memory setting, it performs $O(N \log_2 N)$ random I/Os.

### 2.1.2 The algorithm of BaezaYates-Gonnet-Snider

The algorithm [24] computes incrementally the suffix array $SA$ of the text string $T[1, N]$ in $\Theta(N/M)$ stages. Let $\ell < 1$ be a positive constant fixed below, and assume to set $m = \ell M$. The latter parameter will denote the size of the text pieces loaded in memory at each stage. We also assume for the sake of presentation that $m$ divides $N$.

The algorithm maintains at each stage the following invariant: *At the beginning of stage $h$, where $h = 1, 2, \ldots, N/m$, the algorithm has stored on the disk an array $SA_{ext}$ containing the sequence of the first $(h-1)m$ text suffixes ordered lexicographically and represented via their starting positions (i.e., integers).*

During the $h$th stage, the algorithm *incrementally updates* $SA_{ext}$ by properly inserting into it the text suffixes which start in the substring $T[(h-1)m + 1, hm]$, and by maintaining their lexicographic order. This preserves the invariant above. Hence, after all the $N/m$ stages are executed, it is $SA_{ext} = SA$. We are therefore left with showing how the generic $h$-th stage works.

The text substring $T[(h-1)m + 1, hm]$ is loaded into internal memory, and the suffix array $SA_{int}$ containing only the suffixes starting in that text substring is built by possibly accessing the disk, if needed. [1] Then, $SA_{int}$ is merged with the current $SA_{ext}$ to produce the new array and preserve the invariant. This merging process is executed in two steps with the help of a counter array $C[1, m+1]$. In the former step, the text $T$ is scanned rightwards and the lexicographic position $j$ of each text suffix $T[i, N]$, with $1 \leq i \leq (h-1)m$, is determined in $SA_{int}$ via a binary search. The entry $C[j]$ is then incremented by one unit in order to record the fact that $T[i, N]$ lexicographically lies between the $SA_{int}[j-1]$-th and the $SA_{int}[j]$-th suffix of $T$. In the latter step, the information kept in the array $C$ is employed to quickly merge $SA_{int}$ with $SA_{ext}$: entry $C[j]$ indicates how many consecutive suffixes in $SA_{ext}$ follow the $SA_{int}[j-1]$-th text suffix and precede the $SA_{int}[j]$-th text suffix. This implies that a simple disk scan of $SA_{ext}$ is sufficient to perform such a merging process. At the end of these two steps, the invariant on $SA_{ext}$ has been properly preserved so that $h$ can be incremented and the next stage can correctly start.

Some comments are in order at this point. It is clear that a phase proceeds by mainly executing two disk scans: one is performed to load $T[(h-1)m + 1, hm]$ in internal memory, and another is performed to merge $SA_{int}$ and $SA_{ext}$ via the counter-array $C$. However, some random disk accesses may be necessary in

---

[1] The comparison between any two suffixes can require to access the substring $T[hm + 1, N]$, which is still on disk, thus inducing some random I/Os.

two distinct situations: either when $SA_{int}$ is built or when the lexicographic positions of each text suffix $T[i, N]$ is determined in $SA_{int}$. In both these two cases, we may need to compare a pair of text suffixes which share a long prefix not entirely available in internal memory (i.e., out of $T[(h-1)m+1, hm]$). In the worst case, this comparison will require two sequential disk scans (initiated at the starting positions of these two suffixes) taking $O(N/M)$ bulk I/Os.

As far as the worst-case I/O-complexity is concerned, let us consider the pathological case in which we have $T = a^N$. Here, we need $O((m \log_2 m)N/m)$ bulk I/Os to build $SA_{int}$; $O((h-1)m(\log_2 m)(N/m))$ bulk I/Os to compute the array $C$; and $O(hm/M) = O(h)$ bulk I/Os to merge $SA_{int}$ with $SA_{ext}$. No random I/Os are executed, so that the global number of bulk I/Os is $O(\sum_{h=1}^{N/m} hm(\log_2 m)(N/m)) = O((N^3 \log_2 M)/M^2)$. Since the algorithm processes each loaded block, we may assume that it takes $\Theta(M)$ CPU-time to operate on each of them, thus requiring $O((N^3 \log_2 M)/M)$ overall CPU-time. The total space required is $4N$ bytes for $SA_{ext}$ and $8m$ bytes for both $C$ and $SA_{int}$; plus $m$ bytes [2] to keep $T[(h-1)m+1, hm]$ in internal memory ($\ell$'s value is derived consequently). The merging step can be easily implemented using some extra space (indeed additional $4N$ bytes), or by employing just the space allocated for $SA_{int}$ and $SA_{ext}$ via a more tricky implementation. We adopt for simplicity the former strategy.

Since the worst-case number of total I/Os is cubic, a purely theoretical analysis would classify this algorithm much less interesting than the others discussed in the following sections (see Table 2.1). But there are some considerations that are crucial to shed new light on it, and look at this algorithm from a different perspective. First of all, we must observe that in practical situations, it is very reasonable to assume that each suffix comparison finds in internal memory all the (usually, constant number of) characters needed to compare the two involved suffixes. Consequently, the practical behavior is more reasonably described by the formula: $O(N^2/M^2)$ bulk I/Os and $O((N^2 \log_2 M)/M)$ CPU time. Additionally, the accounting scheme adopted in this paper allows us to evidence some positive features that would undergo unobserved without this accounting model. Indeed, the analysis above has pointed out that all the I/Os are *sequential* and that the actual number of random seeks is $O(N/M)$ (i.e., at most a constant number per stage). Consequently, the algorithm takes fully advantage of the large bandwidth of current disks and of the high computation-power of current processors [16, 40]. Moreover, the reduced working space facilitates the prefetching and caching policies of the underlying operating system (remember Knuth's quote cited in the Introduction) and finally, a careful look to the algebraic calculations shows that the constants hidden in the big-Oh notation are very small. As a result, the adopted accounting scheme does not label the BGS-algorithm as "worse" but drives us to conjecture good practical performances, leaving the final judgment to depend on disk and system specialties. These aspects will be addressed in Section 4.

**Implementation Issues.** The implementation of this algorithm including automatic accounting is tricky if one wants to take into account various pathological cases occurring in the merging step of $SA_{int}$ and $SA_{ext}$, and in the construction of $SA_{int}$. These cases could lead the algorithm to execute many I/Os that would *apparently* be classified as random (according to the accounting scheme) but which are likely to be buffered by the system and thus can be executed much faster. This is not so harmful in itself but would destroy the accounting results. Our main idea is twofold:

- The array $SA_{int}$ is built only on the first $m - B$ suffixes of $T[(h-1)m+1, hm]$, thus we discard the last $B$ suffixes (one page) of that text piece. These discarded suffixes will be (re)considered in the next stage. In such a way, during the construction of $SA_{int}$, each suffix is guaranteed to have a *prefix of at least B* characters available in internal memory, hence *significantly* reducing the probability of a page fault for a suffix comparison.

- To merge $SA_{int}$ and $SA_{ext}$ we need to load the suffixes starting in $T[1, (h-1)m]$ and search them into $SA_{int}$. These suffixes are loaded via a rightward scan of $T$, which brings into internal memory

---

[2] We assume that the alphabet size is smaller than 256.

text pieces whose size is *twice* the size of a bulk I/O. Then, only the suffixes starting in the *first half* of this piece are searched in $SA_{int}$, thus guaranteeing to have in internal memory at least $cM$ characters for their comparison.

These two simple tricks avoid some *"border cases"* which are very likely to induce random I/Os and that instead can be easily canceled via a careful programming. We experimented a dramatic reduction in the total number of random I/Os and a consequent significant speedup in the final performance of the implemented BGS-algorithm (see Section 4).

### 2.1.3   The doubling algorithm

This algorithm was introduced in [1] as a variant of the *labeling technique* of [27], properly adapted to work in external memory. The main idea is first to *logically* pad $T$ at the end with $N$ *special* blank characters which are smaller than any other of $T$'s characters; and then to assign names (i.e. small integers) to the power-of-two length substrings of $T[1, 2N]$ in order to satisfy the so called *lexicographic naming property*: Given two substrings $\alpha$ and $\beta$ of length $2^h$ occurring in $T$, it is $\alpha \leq_L \beta$ if and only if the name of $\alpha$ is smaller than the name of $\beta$. These names are computed inductively by exploiting the following observation: the lexicographic order between any two substrings $\alpha, \beta$ of length $2^h$ can be obtained by splitting them into two equal-length parts $\alpha = \alpha_1 \alpha_2$ and $\beta = \beta_1 \beta_2$, and by using the order between the two pairs of names inductively assigned to $(\alpha_1, \alpha_2)$ and $(\beta_1, \beta_2)$. After $q = \lceil \log_2 (N + 1) \rceil$ stages, the algorithm has computed the lexicographic names for the $2^q$-length substrings of $T[1, 2N]$ starting at positions $1, \ldots, N$ (where $2^q \geq N$). Consequently, the order between any two suffixes of $T$, say $T[i, N]$ and $T[j, N]$, can be obtained in constant time by comparing the lexicographic names of $T[i, i + 2^q - 1]$ and $T[j, j + 2^q - 1]$. This is the rationale behind the doubling algorithm in [1], whose implementation details are sketched below.

At the beginning, the algorithm scans the string $T[1, N]$ and creates a list of $N$ tuples each consisting of four components, say $\langle 0, 0, i, T[i] \rangle$ (the third component will remain fixed afterwards) [3]. During all $q = \lceil \log_2 (N + 1) \rceil$ stages, the algorithm manipulates these tuples by preserving the following invariant: *At the beginning of stage $h$ (initially $h = 1$), tuple $\langle *, *, j, n_j \rangle$ keeps some information about the substring $T[j, j + 2^{h-1} - 1]$ and indeed $n_j$ is its lexicographic name.* [4] After the last $q$-th stage, the suffix array of $T$ is obtained by executing two steps: (i) sort the tuples in output from stage $q$ according to their fourth component; (ii) construct $SA$ by reading from left-to-right the third component of the tuples in the ordered sequence.

We are therefore left with showing how stage $h$ can preserve the invariant above by computing the lexicographic names of the substrings of length $2^h$ with the help of the names inductively assigned to the $2^{h-1}$-length substrings. This is done in four steps as follows.

1. The $N$ tuples (in input to stage $h$) are sorted according to their third component, thus producing a list such that its $i$-th tuple has the form $\langle *, *, i, n_i \rangle$ and thus keeps the information regarding the substring $T[i, i + 2^{h-1} - 1]$.

2. This list is scanned rightwards and the $i$-th tuple is substituted with $\langle n_i, n_{i+2^{h-1}}, i, 0 \rangle$, where $n_{i+2^{h-1}}$ is the value contained in the fourth component of the $(i + 2^{h-1})$-th tuple in the list. [5]

3. The list of tuples is sorted according to their first two components (lexicographic naming property).

---

[3] The blank characters are only *logically* appended to the end of $T$.

[4] In what follows, the symbol $*$ is used to denote an arbitrary value for a component of a tuple, which is actually not important in the discussion.

[5] The rationale behind this step is to represent each substring $T[i, i + 2^h - 1]$ by means of the lexicographic names assigned to its prefix and its suffix of length $2^{h-1}$. These names are inductively available at the beginning of stage $h$.

4. The sorted sequence is scanned rightwards and different tuples are assigned (in their fourth component) with *increasing* integer numbers. This way, the lexicographic naming property is preserved for the substrings of $T$ having length $2^h$.

The correctness of this approach immediately derives from the lexicographic naming property and from the invariant preserved at every stage $h$ which actually guarantees that:

**Property 1** *Each tuple* $\langle *, *, i, * \rangle$ *contains some compact* lexicographic information *about the $2^h$-length prefix of the suffix $T[i, N]$.*

As far as the I/O complexity is concerned, each stage applies twice a sorting routine (steps 1 and 3) and twice a scanning routine (steps 2 and 4) on a sequence of $N$ tuples. The total number of random I/Os is therefore $O(sort(N) \log_2 N)$, and the total number of bulk I/Os is $O((N/M) \log_2 N)$, where $sort(N) = (N/B) \log_{M/B}(N/B)$ is the (random) I/O-complexity of an optimal external-memory sorting algorithm [44]. The CPU-time is $O(N \log_2^2 N)$ since we perform $O(\log_2 N)$ sorting steps on $N$ items. As far as the space complexity is concerned, this algorithm sorts tuples of four components, each consisting of an integer (i.e., 4 bytes). Hence, it seems that 16 bytes per tuple are necessary. Instead, by carefully redesigning the code it is possible to save one entry per tuple, thus using only 12 bytes (overall $12N$ bytes). In summary, the total space complexity is $24N$ bytes because the implementation of the multiway mergesort routine [29], used in our experiments to sort the tuples, needs $2Xb$ bytes for sorting $X$ items of $b$ bytes each (see Section 3).

Two practical improvements are still possible. The first improvement can be obtained by coding four consecutive characters of $T$ into one integer before that the first stage is executed. This allows the saving of the first two stages and hence overall four sorting and four scanning steps. This improvement is not negligible in practice due to the time required by the sorting routine (see Figure 3.1). The second improvement comes from the observation that it is not necessary to perform $\Theta(\log_2 N)$ iterations, but the doubling process can be stopped as soon as all the $2^h$-length substrings of $T$ are different (i.e. all tuples get different names in step 4). This modification does not change the worst-case complexity of the algorithm, but it ensures that only six stages [6] are usually sufficient for natural texts [13].

## 2.2  Our three new proposals

In this section we introduce three new algorithms which asymptotically improve upon the previously known ones by offering better trade-offs between total number of I/Os and working space. Their algorithmic structure is simple because they are based only upon *sorting* and *scanning* routines. This feature has two immediate advantages: The algorithms are expected to be fast in practice because they can benefit from the prefetching of the disk [39]; they can be easily adapted to work efficiently on $D$-disk arrays and clusters of $P$ workstations. It suffices to plug-in proper sorting/scanning routines to obtain a speed-up of a factor $D$ [38] or $P$ [23], approximately (cfr. [28, 36]).

### 2.2.1  Doubling combined with a discarding stage

Our first new algorithm is based on the following observation: *In each stage of the doubling approach, all tuples are considered although the final position in SA of some of them might be already known.* Therefore all those tuples could be discarded from the succeeding sorting steps, thus reducing the overall number of operations (hence I/Os) executed in the next stages. Although this discarding strategy does not give an asymptotic speed up on the overall performance of the algorithm, it is nonetheless expected to induce

---

[6]In natural texts, suffixes share the first 64 characters on average. This accounts for six stages. Together with the first improvement, four stages are sufficient on average.

a significant improvement on experimental data sets because it tends to reduce the number of items on which the sorting/scanning routines are required to work on. [7]

The main idea is therefore to identify in the step 4 of the doubling algorithm (see Section 2.1.3), all the tuples whose final lexicographic position can be inferred using the available information, and then discard them from the next sorting stages. However, these tuples cannot be completely excluded because they might be necessary in the step 1 of the succeeding stages in order to compute the names of longer prefixes of suffixes whose position has been not yet established. In what follows, we describe how to cope with this problem (see Property 2).

As in the original doubling algorithm, we assume that a tuple has three entries (the fourth one has been dropped, see Section 2.1.3), and we call a tuple *finished* if its second component is set to $-1$. The new algorithm inductively keeps two lists of tuples: UT (finished tuples) and UT (unfinished tuples). The former is a sorted list of tuples corresponding to suffixes whose final position in $SA$ is known; they have the form $\langle pos, -1, i \rangle$, where $pos$ is the final position of suffix $T[i, N]$ in $SA$ (i.e., $SA[pos] = i$). UT is a list of tuples $\langle x, y, i \rangle$, corresponding to suffixes whose final position is not yet known, and are such that $x, y \geq 0$ denote *lexicographic names* and $T[i, N]$ is the suffix to which this (unfinished) tuple refers.

At the beginning, the algorithm creates the list UT with tuples having the form $\langle 0, T[i], i \rangle$, for $1 \leq i \leq N$, sets FT to the empty list and initializes the counter $j = 0$. Then, the algorithm proceeds into stages each consisting of the following steps:

1. **Sort** the tuples in UT according to their first two components. If UT is empty goto step 6.

2. **Scan** UT, mark the "finished" tuples and assign new names to all tuples in UT. Formally, a tuple is "finished" if it is preceded and followed in UT by two tuples which are different in at least one of their first two components; in this case, the algorithm marks "finished" the current tuple by setting its second component to $-1$. The new names for all tuples of UT are computed differently from what was done in step 4 of the doubling algorithm (see Section 2.1.3). Indeed, the first component of a tuple $t = \langle x, y, * \rangle$ is now set equal to $(x + c)$, where $c$ is the number of tuples that precede $t$ in UT and have the form $\langle x, p, * \rangle$ with $p \neq y$.

3. **Sort** UT according to the third component of its tuples (i.e., according to the starting position of the corresponding suffix).

4. **Merge** the lists UT and FT according to the third component of their tuples. UT will keep the final merged sequence, whereas FT will be emptied.

5. **Scan** UT and for each *not-finished* tuple $t = \langle x, y, i \rangle$ (with $y \neq -1$), take the next tuple at distance $2^j$ (say $\langle z, *, i + 2^j \rangle$) and change $t$ to $\langle x, z, i \rangle$. If a tuple is marked "finished" (i.e., $y = -1$), then it is discarded from UT and put into FT. Finally set $j = j + 1$ and go to step 2.

6. (UT is empty) FT is sorted according to the first component of its tuples. The third component of the sorted tuples read rightwards, give the final suffix array $SA$.

The correctness can be proved by the following invariant:

**Property 2** *At a generic stage $j$ ($j \geq 0$), the execution of step 2 ensures that a tuple $t = \langle x, y, i \rangle$ satisfies the property that $x$ is the number of $T$'s suffixes whose prefix of length $2^j$ is strictly smaller than $T[i, i + 2^j - 1]$.*

**Proof:** Before step 2 is executed, $x$ inductively accounts for the number of suffixes in $T$ whose $2^{j-1}$-length prefix is lexicographically smaller than the corresponding one of $T[i, N]$. At the first stage ($j = 0$), the

---

[7]Knuth [29][Sect. 6.5] says: *"space optimization is closely related to time optimization in a disk memory"*.

algorithm has indeed safely set the first component of each tuple to 0. When step 2 is executed and tuple $t$ is processed, the variable $c$ accounts for the number of suffixes whose $2^{j-1}$-length prefix is equal to $T[i, i + 2^{j-1} - 1]$ but their $2^j$-length prefix is smaller than $T[i, i + 2^j - 1]$. From the inductive hypothesis on the value of $x$, it then follows that the new value $(x + c)$ correctly accounts for the number of suffixes of $T$ whose $2^j$-length prefix is lexicographically smaller than the corresponding one of $T[i, N]$. ∎

The logic underlying the algorithm above is similar to the one behind the original doubling algorithm (see Section 2.1.3). However, here the new names are assigned in step 2 by following a completely different approach which guarantees not only the lexicographic naming property but also a proper coding of some useful informations (Property 2). This way when a tuple is marked "finished", its first component correctly gives the final position in $SA$ of the corresponding suffix (denoted by its third component). Therefore, the tuple can be safely discarded from UT and put into FT (step 5).

Doubling combined with the discarding strategy performs $O((N/B)(\log_{M/B}(N/B))\log(N))$ random I/Os, $O((N/M)\log_2 N)$ bulk I/Os, and occupies $24N$ bytes (see Section 2.1.3), exactly the same I/O-complexity as the Doubling algorithm. In our implementation, we will also use the compression scheme discussed at the end of Section 2.1.3, to save the first two stages and thus four sorting and four scanning steps. As Section 4 shows, the discarding strategy induces a speed-up in the practical performance of the Doubling approach.

### 2.2.2 Doubling+Discard and Radix Heaps

Although the doubling technique gives the two most I/O-efficient algorithms for constructing large suffix arrays, it has the major drawback that its working space is large (i.e. $24N$ bytes) compared to the other known approaches (see also Table 2.1). This is due to the fact that it uses an external mergesort [29] for ordering the list of tuples and this requires an auxiliary array to store the intermediate results. Our new idea here is to reduce the overall space requirements by making use of an external version of the radix heap data structure introduced in [3]. This data structure [14] uses an exact number $N/B$ of disk pages to store $N$ items, and it is a *monotone integer priority queue* that supports the insertion of a new item in amortized $O(1/B)$ I/Os and the deletion of the item having minimum priority in amortized $O(1/B \log_{M/(B\log_2 V)} V)$ I/Os, where $V$ is the maximum priority value. Hence radix heaps are space efficient but their I/O-performance degenerates when $V$ is large. Our new construction algorithm replaces the mergesort in the Doubling+Discard algorithm (see steps 1 and 3 in Section 2.2.1) with a sorting routine based on that external radix heap. This reduces the overall required space to $12N$ bytes, but at the cost of increasing the I/O–complexity to $O((N/B)(\log_{M/(B\log_2 N)} N)\log_2 N)$ random I/Os (and $O((N/M)\log_2 N)$ bulk I/Os), because $V = N$ in the step 3 of Section 2.2.1 (indeed, the third component of the sorted tuples ranges in $[1, N]$). We observe that this algorithm should outperform the doubling approach during the first stages because the range of assigned names, and thus the value of $V$, is sufficiently small to take advantage from the radix-heap structure. On the other hand, the algorithm performance degenerates as more stages are executed because $V$ becomes larger and larger. It is therefore interesting to experimentally investigate this solution since it significantly saves space and is expected to behave well in practice even in the light of the reduction in the number of tuples to be sorted in each stage.

### 2.2.3 Construction in $L$ pieces

The approaches described before are I/O-efficient but they use at least $8N$ bytes of working space. If the space issue is a primary concern, and we still wish to keep the total number of I/Os small (unlike [24]), different approaches must be devised that require much less space but still guarantee good

I/O-performances [8]. In this section, we describe one such approach which improves over all previous algorithms in terms of both I/O complexity, CPU time and space occupancy. It constructs the suffix array into *pieces of equal sizes* and thus it turns out useful either when the underlying application *does not* need the suffix array as a unique data structure, but allows to keep it in a distributed fashion [10]; or when we operate in a distributed-memory environment [28, 36].

The most obvious way to achieve this goal might be to partition the original text string $T[1, N]$ into equal-length substrings and then apply on these pieces any known suffix-array construction algorithm. However, this approach should cope with the problem of correctly handling the suffixes which *start close to the border* of any two text pieces. To circumvent this problem, some authors [28, 36] reduce the suffix array construction process to a string sorting process associating to each suffix of $T$, its prefix of length $X$, and then sorting $N$ strings of length $X$ each. Clearly, the correctness of this approach heavily depends on the value of $X$ which also influences the space occupancy (it is actually $NX$ bytes). Statistical considerations and structural informations about the underlying text might help, but anyway the choice of the parameter $X$ strongly influences the final performance (see e.g. [28, 36]).

The approach we introduce below is very simple and applies in a different way, useful for practical purposes, a basic idea known so far only in the theoretical setting (see e.g. [18]). Let us denote by $\mathcal{A}_{sa}$ any external-memory algorithm for building a suffix array, $\mathcal{A}_{string}$ any external-memory algorithm for sorting a set of strings, and let $L$ be a constant integer parameter to be fixed later. For simplicity of exposition, we assume that $N$ is a multiple of $L$, and that $T$ is *logically* padded with $L$ blank characters.

The new approach constructs $L$ suffix arrays, say $SA_1, SA_2, \ldots, SA_L$ each of size $\Theta(N/L)$. Array $SA_i$ stores the lexicographically ordered sequence of suffixes $\{T[i, N], T[i + L, N], T[i + 2L, N], \ldots, \}$. The logic underlying our algorithm is to first construct $SA_L$, by using $\mathcal{A}_{sa}$ and $\mathcal{A}_{string}$, and then derive all the others arrays $SA_{L-1}, SA_{L-2}, \ldots, SA_1$ by means of a simple external-memory algorithm for sorting *triples of integers* (e.g., multiway mergesort [29]).

The suffix array $SA_L$ is built in two main stages: In the first stage, the string set $S = \{T[L, 2L - 1], T[2L, 3L - 1] \ldots, T[N - L, N - 1], T[N, N + L - 1]\}$ is formed and lexicographically sorted by means of algorithm $\mathcal{A}_{string}$. In the second stage, a compressed text $T'$ of length $N/L$ is derived from $T[L, N + L - 1]$ by replacing each string having the form $T[iL, (i + 1)L - 1]$, for $i \geq 1$, with its *rank* in the *sorted* $S$. Then, algorithm $\mathcal{A}_{sa}$ builds the suffix array $SA'$ of $T'$, and finally derives $SA_L$ by setting $SA_L[j] = SA'[j] \cdot L$, for $j = 1, 2, \ldots, N/L$.

Subsequently, the other $L - 1$ suffix arrays are constructed by exploiting the following observation: *Any suffix $T[i + kL, N]$ in $SA_i$ can be seen as the concatenation of the character $T[i + kL]$ and the suffix $T[i + 1 + kL, N]$ occurring to $SA_{i+1}$.* So that, if $SA_{i+1}$ is known, the order between $T[i + kL, N]$ and $T[i + hL, N]$ can be obtained by comparing the two pairs of integers $\langle T[i + kL], pos_{i+1+kL} \rangle$ and $\langle T[i + hL], pos_{i+1+hL} \rangle$, where $pos_s$ is the position in $SA_{i+1}$ of suffix $T[s, N]$. This immediately means that the construction of $SA_i$ can be reduced to sorting $\Theta(N/L)$ tuples, once $SA_{i+1}$ is known.

Sorting via $\mathcal{A}_{string}$ the short strings of length $L$ takes $O(Sort(N))$ random I/Os and $2N + 8N/L$ bytes, where $Sort(N) = (N/B) \log_{M/B}(N/B)$ [1]. Building the $L$ suffix arrays $SA_i$ takes $O(Sort(N/L) \log_2(N/L) + L\, Sort(N/L)) = O(Sort(N) \log_2 N)$ random I/Os, $O(N/M \log_2(N/M))$ bulk I/Os and $24N/L$ bytes. Of course, the larger is the constant $L$, the larger is the number of suffix arrays that will be constructed, but the smaller is the working space required. By setting $L = 4$, we get an interesting algorithm for constructing large suffix arrays: it needs $6N$ bytes working space, $O(Sort(N) \log_2 N)$ random I/Os and $O(N/M \log_2(N/M))$ bulk I/Os. Its practical performance will be evaluated in Section 4. Notice that this approach builds four suffix arrays, thus its query performance is slowed down by a constant factor four, but this is practically negligible in the light of suffix-array search performance [32].

---

[8]See the footnote 7 and refer to [48] where Zobel *et al.* say that: "*A space-economical index is not cheap if large amounts of working storage are required to create it.*"

# 3 Our experimental settings

## 3.1 An external-memory library

We implemented all algorithms discussed so far by using a recently developed external-memory library of algorithms and data structures called `LEDA-SM` [15] (an acronym for "LEDA for Secondary Memory"). [1] This library is an extension of the internal-memory library *LEDA* [37] and therefore follows its main underlying ideas. One of them is *portability* to a variety of platforms and another one is *high level specification* of data structures. Library `LEDA-SM` consists of a sizeable collection of efficient data structures and algorithms explicitly designed to work in an external memory setting. The system underlying `LEDA-SM` reflects a real view of the external memory model [44]: The internal memory is directly provided by the internal memory of the computer; whereas the $D$ abstract disks are modeled by the file system, which also provides proper tools for implementing the low-level I/O via block transfers. Each disk is modeled with a *single file* and it is divided into logical blocks of a fixed size $B$ (disk pages). The size of this file is fixed, thus modeling the fact that real disk space is bounded. Since `LEDA-SM` uses the file system, it explicitly takes advantage of the underlying *I/O-buffering* and *read-ahead* strategy at no additional implementation effort.[2] The part of the system which controls the I/O is the so called *external memory manager*. It is implemented as an interface to which are connected all high level data structures and algorithms, thus guaranteeing that the I/O and the management of the data on the disks is hidden to the programmer. Nevertheless, the programmer can still keep track of the number of disk accesses (bulk, random and total) performed by his/her algorithms, because the system allows the explicit counting of the number of both writes and reads to the $D$ disks. At the moment, `LEDA-SM` provides the possibility to choose among five different file system access methods, namely standard I/O (stdio), system call I/O (syscall), asynchronous I/O (aio), serial file I/O (sfio) and memory mapped I/O (mmap_io) (see [45] for a deeper discussion). Both the external memory manager and the data structures provided by `LEDA-SM` are implemented in C++ as a set of template classes and functions. The specialty of `LEDA-SM`'s data structures is that during their constructions it is possible to specify (and therefore control) the maximum amount of internal memory that they are allowed to use. Combining this feature with the counting of I/Os, library `LEDA-SM` allows the programmer to experimentally investigate how the model parameters $M$ and $B$ influence the performance of an external-memory algorithm.

For what concerns the implementation of our suffix-array construction algorithms, we used the external array data structure and the external sorting/scanning algorithms provided by `LEDA-SM` library. In particular, we used an implementation of multiway mergesort that needs $2Xb$ bytes for sorting $X$ items of $b$ bytes each (see Figure 3.1). The other in-core algorithms and data structures used in our experiments are taken from the `LEDA`-library. To avoid that the internal-memory size prevents the use of Manber-Myers' algorithm on large text collections, we run it in a virtual memory setting by using swap space. All other algorithms are not faced with this problem because they are directly designed to work in

---

[1] There exist other external memory libraries, the most notable one is TPIE [46].

[2] As the file system is allowed to buffer disk pages, some of the disk requests can immediately be satisfied by the I/O–buffer. This notably speeds up the real performances of the external-memory algorithms, as we discussed in Section 1.
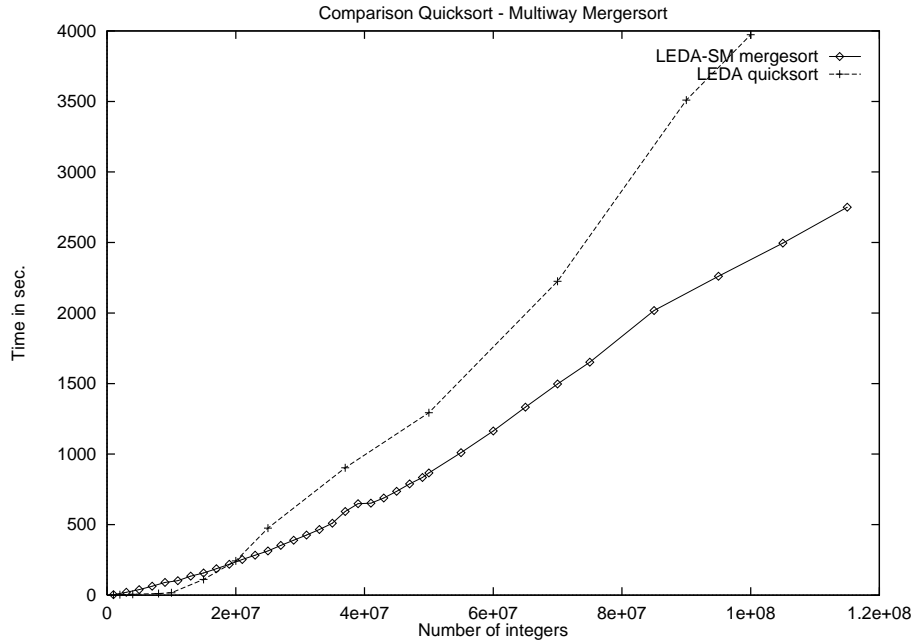
Figure 3.1: *Comparison between LEDA-quicksort and LEDA-SM multiway mergesort. Quicksort is allowed to use at most 64 Mbytes of internal memory, the further needs of memory are satisfied by using swap space. Multiway mergesort is only allowed to use 12 Mbytes of internal memory.*

external memory. All the construction algorithms for suffix arrays discussed in this paper are available and necessitate both the **LEDA** and **LEDA-SM** libraries, please contact one of the authors if you wish to play with them.

### 3.1.1 System parameters

The computer used in our experiments is a SUN ULTRA-SPARC 1/143 with 64 Mbytes of internal memory running the SUN Solaris 2.5.1 operating system. It is connected to one single Seagate Elite-9 SCSI disk via a fast-wide differential SCSI controller. The external memory library **LEDA-SM** provides logical disks of maximum size of 2 Gbytes (this is the file size limit of Solaris 2.5.1), the disks are divided into blocks of $B = 8$ Kbytes [3]. We have chosen the standard I/O (stdio) for the file system accesses because it is available on almost every system. We have designed our external-memory algorithms so that they use approximately 36 Mbytes of internal memory. This is obtained by properly choosing the internal memory size dedicated to the external memory data structures (e.g. external arrays and mergesort buckets). Thus, we guarantee that the studied algorithms do not incur in the *paging phenomenon* even for accessing their internal data structures and furthermore, that there is left room in internal memory to keep the I/O buffers of the operating system.

### 3.1.2 Choosing the bulk-load size

Modern SCSI disk drives are equipped with a large cache, which is used to hide the speed gap between the SCSI bus and the disk drive. The disk drive *prefetches* some data into the cache before it initiates the transfer of the requested data from cache to main memory, in order to hide the difference between the rotational delay of the disk and the bandwidth of the bus. Moreover, during read accesses the disk

---

[3]We note that $B = 64$ Kbytes is "optimal" for random disk access. We have chosen $B = 8$ Kbytes because we also compare to virtual memory algorithms that use the machine's page size which is exactly 8 Kbytes.

prefetches further data into the cache in order to hopefully serve the next read requests directly from the disk cache, thus avoiding the costly (mechanical) *seek operation* or suffer losses due to *rotational delay*.

According to the accounting scheme we adopted in this paper (see Section 1), we need to choose a reasonable value for the *bulk load size* (notice that the parameters $B$ and $M$ are fixed by default according to the computer features). Since the transfer rates are more or less stable (and currently large) while seek times are highly variable and costly (because of their mechanical nature), our idea is to choose a value for the bulk size which allows to hide the *extra cost* induced by the seek step when data are fetched. This would allow us to access "uniformly" to every datum stored on the disk, thus working at the highest speed allowed by its bandwidth. In some sense, we would like to *hide* the mechanical nature of the disk system [16].

Let the transfer time be described by the formula `t_seek + bulk_size/disk_bandwidth` [40]. We wish that the first term is much smaller than the second one. Consequently, from one side we should increase `bulk_size` as much as possible (at maximum $M/B$), but from the other side a large `bulk_size` might reduce the significance of our accounting scheme because many *long* sequential disk scans could result shorter than `bulk_size` and thus counted as 'random', whereas they nicely exploit the disk caching and prefetching strategies [19]. Hence a proper "tuning" of this parameter is needed according to the mechanical features of the underlying disk system.

In the disk used for our experiments, the average `t_seek` is 11 msecs, the `disk_bandwidth` is 7 Mbytes/sec. We have therefore chosen `bulk_size` = 64 disk pages, for a total of 512 Kbytes. It follows that `t_seek` is 15% of the total transfer time needed for a bulk I/O. Additionally, the bulk size of 512 Kbytes allows us to achieve 81% of the maximum data transfer rate of our disk while keeping the service time of the requests still low.

According to our considerations above (see also [16]), we think that this is a reasonable choice even in the light of the disk cache size of 1 Mbytes. Surprisingly, we also noticed in our experiments that this value allows us to catch the execution of numerous bulk I/Os by subroutines (e.g. multiway mergesort) which were defined "mainly random" at a theoretical investigation. Clearly, other values for `bulk_size` might be chosen and experimented, thus achieving different trade-offs between random/bulk disk accesses. However, the *qualitative* considerations on the algorithmic performance drawn at the end of the next section will remain mostly unchanged, thus fitting our experimental desires.

## 3.2   Textual data collections

For our experiments we collected over various WEB sites three textual data sets. They consist of:

- The Reuters corpus[4] together with other natural English texts whose size sum up to 26 Mbytes. This collection has the nice feature of presenting long repeated substrings.

- A set of amino-acid sequences taken from a SWISSPROT database[5] summing up to around 26 Mbytes. This collection has the nice feature of being an unstructured text so that full-text indexing is the obvious choice to process these data.

- A set of *randomly generated* texts consisting of three collections: one formed by texts randomly drawn from an alphabet of size 4, another formed by texts randomly drawn from an alphabet of size 16, and the last one formed by texts randomly drawn from an alphabet of size 64. These collections have two nice features: they are formed by unstructured texts, and they constitute a good test-bed to investigate the influence of the *length of the repeated substrings* on the performance

---

[4]We used the text collection "Reuters-21578, Distribution 1.0" available from David D. Lewis' professional home page, currently: `http://www.research.att.com/~lewis`

[5]See the site: `http://www.bic.nus.edu.sg/swprot.html`

of some studied algorithms. For each alphabet size we consider texts of 25 Mbytes and 50 Mbytes, thus further enlarging the spectrum of text sizes on which the algorithms are tested.

A comment is in order at this point. The reader might observe that, although this paper is on constructing suffix arrays on *large* text collections, our experimental data sets seem indeed *"not very large"*! If we look just to their sizes, the involved numbers are actually not very big (at most 50 Mbytes); but, as it will soon appear clear, our data sets are sufficiently large to evaluate/compare in a fair way the I/O-performance of the analyzed algorithms, and investigate their *scalability* in an external-memory setting. In fact, the suffix array $SA$ needs $4N$ bytes to index a text of length $N$. Hence, the text plus $SA$ globally occupy 200 Mbytes, when $N = 50$ Mbytes. Additionally, each of the algorithms discussed in our paper requires at least $8N$ bytes of working space; this means 400 Mbytes for the largest text size. In summary, *more than* 600 Mbytes will be used during the overall construction of $SA$ in each of the experimented algorithms, when $N = 50$ Mbytes ! Now, since 64 Mbytes is the size of the available internal memory of our computer, all the experiments will run on disk, and therefore the performance of the studied algorithms will properly reflect their I/O-behavior.

# 4 Experimental Results

## 4.1 Experiments on natural data

We now comment the results obtained for our six different construction algorithms: Manber and Myers' algorithm (MM), BaezaYates-Gonnet-Snider's algorithm (BGS), original doubling (Doubl), doubling with discarding (Doubl+Disc), doubling with discarding and external radix heaps (Doubl+Disc+Radix), and the 'construction into pieces' approach (L-pieces). The overall results are reported in Figure 4.1, and they are detailed in Tables 4.1 and 4.2 below.

| The Reuters corpus | | | | | | |
|---|---|---|---|---|---|---|
| N | MM | BGS | Doubl | Doubl+Disc | Doubl+Disc+Radix | L-pieces |
| 1324350 | 67 | 125 | 828 | 982 | 1965 | 331 |
| 2578790 | 141 | 346 | 1597 | 1894 | 3739 | 582 |
| 5199134 | 293 | 1058 | 3665 | 4070 | 7833 | 1119 |
| 10509432 | 223200 | 4808 | 8456 | 8812 | 16257 | 2701 |
| 20680547 | – | 16670 | 23171 | 20434 | 37412 | 5937 |
| 26419271 | – | 27178 | 42192 | 28937 | 50324 | 7729 |
| The Amino-Acid Test | | | | | | |
| 26358530 | – | 20703 | 37963 | 24817 | 41595 | 6918 |

Table 4.1: *Construction time (in seconds) of all experimented algorithms on two text collections: the Reuters corpus and the Amino-acid data set. N is the size of the text collection in bytes. The symbol '–' indicates that the test was stopped after 63 hours.*

### 4.1.1 Results for the Manber-Myers' algorithm.

It is not astonishing to observe that the construction time of MM-algorithm is outperformed by every other algorithm studied in this paper as soon as the working space exceeds the memory size (i.e., 64 Mbytes). This worse behavior is due to the fact that the algorithm accesses the suffix array in an unstructured and unpredictable way. In fact its paging activity almost crashes the system, as we monitored by using the Solaris-tool *vmstat*. The vmstat value *sr*, which monitors the number of page scans per second, was constantly higher than 200. According to the Solaris system guide, this indicates that the system is constantly out of memory. Looking at Table 2.1, we infer that MM-algorithm should be chosen only when the text is shorter than $M/8$. In this case, every data structure fits in internal memory and thus the disk is never used. In our experimental setting, this actually happens for $N \leq 8$ Mbytes. When $N > 8$ Mbytes, the time complexity of MM-algorithm is still *quasi-linear* but the constant hidden in the big-Oh notation is very large due to the paging activity, thus making the algorithmic behavior unacceptable.

| The Reuters corpus | | | | | |
|---|---|---|---|---|---|
| N | BGS | Doubl | Doubl+Disc | Doubl+Disc+Radix | L-pieces |
| 1324350 | 120/7865 | 2349/256389 | 2242/199535 | 4872/377549 | 837/57282 |
| 2578790 | 317/20708 | 4517/500151 | 4383/395018 | 10075/787916 | 1693/177003 |
| 5199134 | 929/60419 | 9095/1009359 | 8916/809603 | 22466/1761273 | 3386/360210 |
| 10509432 | 4347/282320 | 18284/2041285 | 18126/1655751 | 47571/3728159 | 6849/730651 |
| 20680547 | 14377/933064 | 35935/4017664 | 35904/3293234 | 96292/7550794 | 14243/1530995 |
| 26419271 | 24185/1568947 | 45911/5132822 | 45842/4202902 | 129071/10001152 | 18178/1956557 |
| The Amino-Acid Test | | | | | |
| 26358530 | 24181/1568773 | 41709/4656578 | 39499/3539148 | 105956/8222236 | 16118/1719883 |

Table 4.2: *Number of I/Os (bulk/total) of all experimented algorithms on two text collections: the Reuters corpus and the Amino-acid data set. N is the size of the text collection in bytes; 64 disk pages form a bulk-I/O.*

### 4.1.2 Results for the BaezaYates–Gonnet–Snider's algorithm.

As observed in Section 2.1.2, the main *theoretical* drawback of this algorithm is the high (i.e., cubic) *worst-case* asymptotic complexity, but its small working space, its regular pattern of disk accesses and the small constants hidden in its big-Oh complexity has lead us in the previous sections to think favorably of BGS for practical uses. Moreover, we conjectured a quadratic I/O-behavior in practice because of the *short* repeated substrings which usually occur in real texts. Our experiments show that we were right in all these suppositions. Indeed, if we double the text size, the running time increases by nearly a factor of four (see Table 4.1), and the number of bulk and random I/Os increase quadratically (see Table 4.2). The number of total and bulk I/Os is nearly identical for all data sets (Reuters, Amino-Acid and Random, see table 4.4), so that the *practical behavior* is actually quadratic. Furthermore, it is not astonishing to verify experimentally that BGS is *faster* than any Doubling variant on the Reuters corpus and the Amino-Acid data set (see Figure 4.1). Consequently, it turns out to be the fastest algorithm for building a (unique) suffix array when $N \leq 25$ Mbytes. This scenario probably remains unchanged for text collections which are slightly larger than the ones we experimented in this paper; after that, the quadratic behavior of BGS will be probably no longer *"hidden"* by its nice algorithmic properties (largely discussed in Section 2.1.2).

In Table 4.2, we notice that (i) only the 1% of all disk accesses are random I/Os (hence, most of them are bulk I/Os !); (ii) the algorithm performs the least number of random I/Os on both the data sets; (iii) BGS is the fastest algorithm to construct one unique suffix array, and it is the second fastest algorithm in general. Additionally, we observe that the quadratic CPU-time complexity of the BGS-algorithm heavily influences (i.e., slows down) its efficiency and thus I/O is not the only bottleneck.

In summary, the BGS-algorithm is amazingly fast on *medium-size* text collections, and remains reasonably fast on larger data sets. It is not absolutely the fastest on larger and larger text collections because of its quadratic CPU- and I/O-complexities. Nonetheless, the small working space (possibly $4N$ bytes via a tricky implementation) and the ease of programming make the BGS-algorithm very appealing for software developers and practitioners, especially in applications where the space issue is the primary concern.

### 4.1.3 Results for the Doubling algorithm.

The doubling algorithm performs 11 stages on the Reuters corpus, hence it performs 21 scans and 21 sorting steps on $N$ tuples, where $N$ is the text size. Consequently, we can conclude that there is a repeated substring in this text collection of length about $2^{12}$ (namely we detected a duplicated article). Figure 4.1
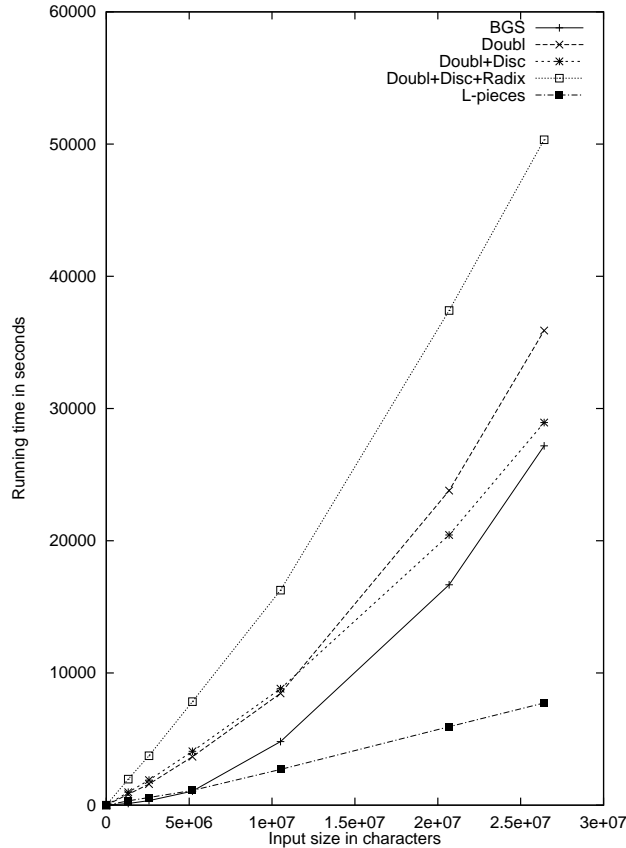
Figure 4.1: *Run-time of all construction approaches on the Reuters corpus.*

shows that the Doubling algorithm scales well in the tested input range: [1] If the size of the text doubles, the total running time doubles too. Among all disk accesses the 41% of them are random, and the number of bulk and random I/Os is larger than the ones executed by every other tested algorithm, except the Doubl+Disc+Radix variant which has higher worst-case complexity but smaller working space (see Table 4.2 and Section 2.2). Due to the high number of random I/Os and to the large working space (see Table 2.1), the Doubling algorithm is expected to surpass the performance of BGS only for *very large* values of $N$.

In summary, although theoretically interesting and almost asymptotically optimal, the Doubling algorithm is not appealing in practice; this motivated our development of the two variants discussed in Section 2.2 (namely, Doubl+Disc and Doubl+Disc+Radix algorithms).

### 4.1.4 Results for the Doubl+Disc algorithm.

If we add the discarding strategy described in Section 2.2.1 to the Doubling algorithm, we achieve better performances as conjectured. The gain induced by the discarding step is approximately the 32% of the original running time both for the Reuters corpus and the Amino-acid data set (for large $N$). If we look in detail at the number of discarded tuples (see Figure 4.3), we see that for the Reuters corpus, this

---

[1]Notice that the curve for Doubling is not always linear. There is a "jump" in the running time as soon as the input has size $21 \cdot 10^6$. This is due to the system's pager-daemon that is responsible for managing the memory-pages. During the merge-step of multiway-mergesort, the number of free pages falls below a given threshold and the pager-daemon tries to free some memory pages. This goal can't be achieved because the mergesort constantly loads new blocks into memory. Therefore, this pager process runs almost all the time and increases the CPU-time of the merge-step of a factor of two.
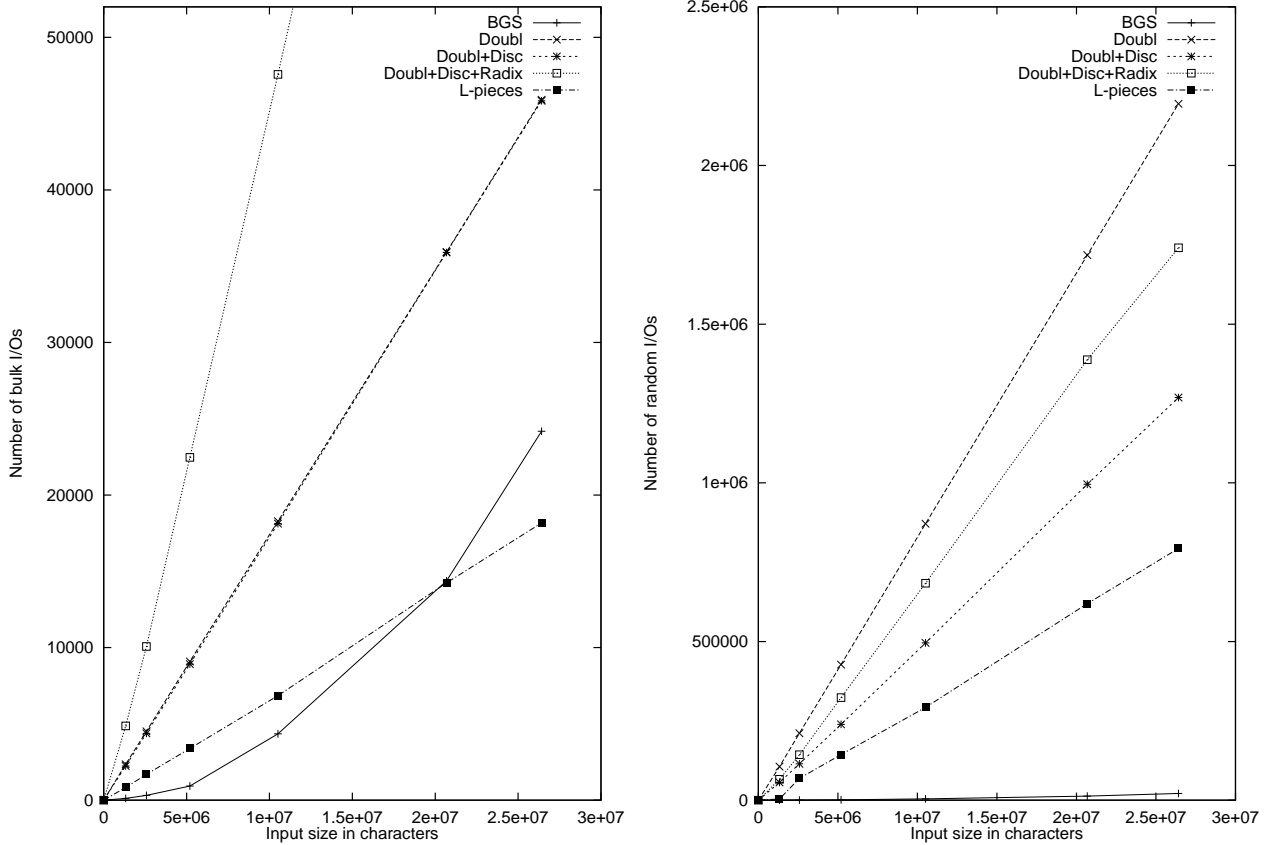
Figure 4.2: *Bulk and Random I/Os for all construction approaches. The bulk size is* 64 *disk pages.*

is small in the first two stages, while it becomes significant in stages three and four, where nearly 55% of the tuples are thrown away. Since we double the length of the substrings at each stage and we use the compression scheme of Section 2.1.3, we can infer that the Reuters corpus has a lot of substrings of length $16 - 32$ that occur once in the collection but their prefixes of length $8 - 16$ occur at least twice. We also point out that the curve indicating the number of discarded tuples is nearly the same as the size of the test set increases. This means that the number of discarded tuples is a "function" of the structure of the indexed text. For our experimental data sets, we save approximately 19% of the I/Os compared to Doubling. The percentage of random I/Os is 28%, this is much less than Doubling (42%), and drives naturally us to observe that discarding helps in reducing mainly the random I/Os (see also Table 4.2). The saving induced by the discarding strategy is expected to pay much more on larger text collections, because of the significant reduction in the number of manipulated tuples at the early stages, which should facilitate caching and prefetching operations (see footnote 7). Consequently, if the time complexity is a much more important concern than the space occupancy, the Doubl+Disc algorithm is definitively the choice for building (unique) *very large* suffix arrays.

### 4.1.5    Results for the Doubl+Disc+Radix algorithm.

This algorithm is not as fast as we conjectured in Section 2.2.2, even for the tuple ordering of Step 2. The reason we have drawn from the experiments is the following. Notice that radix heaps are integer priority queues, and thus we cannot keep the parameter $V$ small (to exploit Radix Heaps properties) by defining the first two components of a tuple as its priority. Hence, the sorting in Step 2 must be implemented via
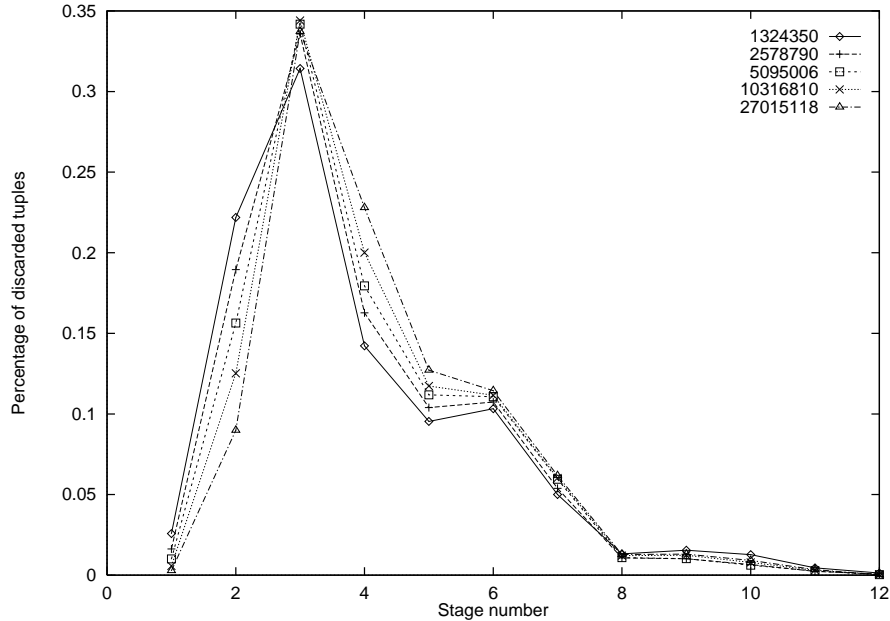
Figure 4.3: *The percentage of discarded tuples at each stage of the Doubl+Disc algorithm on the Reuters corpus.*

two sorting phases and this naturally *doubles* the work executed in Step 2. Additionally, the compression scheme of Section 2.1.3 cannot be used here because it would increase $V$ too much in the early stages. Hence, the algorithm performs two more stages than Doubl-algorithm, and furthermore it does not take advantage of the discarding strategy because the number of discarded tuples in this two initial stages is very small. This way, it is not surprising to observe in Table 4.2 that the Doubl+Disc+Radix algorithm performs *twice* the I/Os of the other Doubling variants, and it is the slowest among all the tested algorithms. This result allows us to conclude that the *compensation* conjectured in Section 2.2.2 between the number of discarded tuples and the increase in the I/O-complexity do not actually arises in practice.

In summary, the Doubl+Disc+Radix algorithm can be interesting only in the light of its space requirements. However, if we compare the space vs. time trade-off we can reasonably consider the Doubl+Disc+Radix algorithm *worse* than the BGS-algorithm because the former requires larger working space and it is expected to surpass the BGS-performance only for *very large* text collections (see Section 6 for further comments).

### 4.1.6 Results for the L-pieces algorithm.

We fixed $L = 4$, used multi-way mergesort for sorting short strings and the Doubl-algorithm for constructing $SA_L$ (see Section 2.2.3). Looking at Table 4.2 we notice that 40% of the total I/Os are random I/Os, and that the algorithm executes slightly more I/Os than the BGS-algorithm. Nonetheless, as shown in Figure 4.1 this algorithm is the fastest one. It is three to four times faster than BGS (due to its quadratic CPU-time) and four times faster than the Doubl+Disc algorithm (due to the larger number of I/Os). The running time distributes as follows: 63% of the overall time is used to build the compressed suffix array $SA_L$; the sorting of the short strings required only 4% of the overall time; the rest is used to build the other three suffix arrays. It must be said that for our test sizes, the short strings fit in internal memory at once thus making their sorting stage very fast. However, it is also clear that sorting short strings takes no more time than the one needed by *one* stage of the Doubl-algorithm. Furthermore, the construction of

the other three suffix arrays, when executed entirely on disk, would account for no more than 1.5 stages of the Doubl-algorithm. Consequently, even in the case where the short strings and the other three suffix arrays reside on the disk, it is not hazardous to *conjecture* that this algorithm is still significantly faster than all the other approaches. The only "limit" of this algorithm is that it constructs the suffix array in four distinct pieces. Clearly, if the underlying text retrieval applications does not impose to have *one unique* suffix array, then this approach turns out to be de-facto *'the'* choice for constructing such a data structure (see Section 6 for further discussions).

### 4.1.7 Comparison of all construction approaches.

We first compare all algorithms with respect to their time performances for increasing text lengths (see Figure 4.1). It is obvious from the discussions above that the MM-algorithm is the fastest one when the suffix array can be built entirely in internal memory. As soon as the working space crosses the 'memory border', there are various possible choices. The L-pieces algorithm is the natural choice whenever the splitting of the suffix array does not prevent its use in the underlying application. If instead a unique suffix array is needed, then the choice lies between the BGS-algorithm and the Doubl+Disc algorithm. For text collections which are not very big, the BGS-algorithm is preferable: it offers fast performance and very small working space. For larger text collections, the choice depends on the primary resource to be minimized: time or space? In the former case, the Doubl+Disc algorithm is the winning choice; in the latter case, the BGS-algorithm is the best.

The Doubl+Disc+Radix variant could be interesting on huge text collections but this consideration cannot be claimed for sure because when $N$ grows, the internal-memory size is also expected to grow ! Consequently, although the idea of using a space-efficient external-memory heap to implement the sorting step is valuable, a different heap structure should be devised to efficiently cope with large integer priorities and efficient I/O-performance. (This topic will be addressed in Section 6.)

With respect to the ease of programming, the BGS-algorithm seems still the best choice, unless the software developer has a library of general external sorting routines, in which case all Doubling variants turns out to be simple too.

## 4.2 Experiments on random data

The previous section left open three important questions:

1. How do the structural properties of the indexed text influence the performance of the Doubl+Disc algorithm? Do these properties decrease the number of stages, and thus significantly improve its overall performance? How big is the improvement induced by the discarding strategy on more structured texts?

2. What happens to the BGS-algorithm if we increase the text size? Is its practical behavior "independent" of the text structure?

3. What happens to the L-pieces algorithm when the four suffix arrays reside on disk? Is this algorithm still significantly faster than all the other ones?

In this section we will answer *positively* all these questions by performing an extensive set of experiments on *three* sets of textual data which are *randomly and uniformly* drawn from alphabets of size 4 (random-small), size 16 (random-middle) and size 64 (random-large), respectively. For each alphabet size, we will indeed generate *two* text collections of 25 Mbytes and 50 Mbytes.

The choice of "randomly generated data sets" is motivated by the following two observations. By varying the alphabet size we can study the impact that the *average length* of the repeated substrings

has on the performance of the discarding strategy. Indeed, the smaller is this length, the larger should be the number of tuples which are discarded at earlier stages, and thus the bigger should be the speed-up obtained by the Doubl+Disc algorithm. The experiments carried out on the Reuters corpus (see Section 4) did not allow us to complete this analysis because of the structural properties of this text collection. Unfortunately, the Reuters corpus represents a pathological case because it has many *long* repeated text-substrings and this is usually untypical for natural linguistic texts. This was the reason why we conjectured at the end of Section 4 a much larger gain from the discarding strategy on more structured texts. An early validation of this conjecture was provided by the experiments carried out on the Amino-acid data set (see Table 4.1). Now we expect that the two random collections will be a good test-bed for providing further evidence. The same thing can be said about the BGS-algorithm whose "independent behavior from the structure of the indexed text in practice", conjectured in Section 4, can now be tested on *variously structured* text collections.

Second, by varying the size of the indexed collection we can investigate the behavior of the L-pieces algorithm when the ordering of the short strings and the construction of the suffix arrays $SA_1, SA_2, SA_3$ operate directly on the disk (and not in internal memory). We can also test if the average length of the repeated substrings can influence the performance of $\mathcal{A}_{sa}$ when building $SA_4$. We notice that the construction of the other three suffix arrays is clearly not influenced by the structure of the underlying text, because it consists of just three sorting steps executed on a sequence of integer triples. In Section 4 we conjectured that a larger text collection should not significantly influence the overall performance of the L-pieces algorithm because, apart from the construction of $SA_4$, all the other algorithmic steps account for 2.5 stages of the Doubl-algorithm. Consequently, the overall work should be always much smaller than the one executed by all the Doubling variants. In what follows, we will validate this conjecture by running the L-pieces algorithm on larger data sets.

### 4.2.1  Results for the BGS-algorithm.

It may appear surprising that BGS is the *slowest* algorithm on the random texts, after its successes on real text collections claimed in Section 4. It is more than twice slower than Doubl+Disc and up to nine times slower than L-pieces. However, nothing strange is going on in these experiments on random data because if we compare Table 4.4 to Table 4.2, we notice two things: (i) the number of bulk and total I/Os executed by BGS do not depend on the alphabet size and they are almost identical to those obtained on the Reuters corpus; (ii) the execution time of BGS decreases as the alphabet size grows, and it is smaller than the time required on the Reuters corpus (when $N \approx 25$ Mbytes). The former observation implies that our conjecture on the " quadratic I/O-behavior in practice" is true, and in fact if we double the text size, the total number of I/Os increases by a factor of approximately four. The latter observation allows us to conclude that the CPU-time of BGS is affected by the length of the repeated substrings occurring in the data set. As the alphabet size increases, this length decreases on the average, and in turn decreases the running time of BGS. Such a dependence also shows that *both* I/O and CPU-time are significant bottlenecks for BGS, as pointed out already in Section 4.

### 4.2.2  Results for the Doubl+Disc algorithm.

The algorithm performs 4 stages on alphabet size 4, and 3 stages on alphabet sizes 16 and 32. Therefore, the random test with alphabet size 4 is the *worst case* for Doubl+Disc. In fact, the gain of the first two stages is negligible ($0.8 \cdot 10^{-8}\%$ discarded tuples); in the third stage, 98% of the tuples are discarded; the rest of the tuples is thrown away in the last fourth stage. This gives us the following insight: There are many *different* text substrings of 32 characters (i.e., 98%) whose 16-length prefix occurs at least twice in the collection. This validates our conjecture that the Doubl+Disc performance heavily depends on the average length of the repeated text-substrings.

| Running times on the random texts | | | |
|---|---|---|---|
| N | L-pieces | Doubl+Disc | BGS |
| Alphabet Size 4 | | | |
| 25000000 | 4133 | 14130 | 21485 |
| 50000000 | 8334 | 34599 | 72552 |
| Alphabet Size 16 | | | |
| 25000000 | 3838 | 11011 | 16162 |
| 50000000 | 7753 | 26450 | 62001 |
| Alphabet Size 64 | | | |
| 25000000 | 3400 | 10080 | 15195 |
| 50000000 | 6802 | 25606 | 58417 |

Table 4.3: *Construction time (in seconds) required for random texts built on alphabet-sizes* 4, 16 *and* 64.

| I/Os (bulk/total) on the random texts | | | |
|---|---|---|---|
| N | L-pieces | Doubl+Disc | BGS |
| Alphabet Size 4 | | | |
| 25000000 | 9466/970256 | 20661/2068791 | 22347/1449884 |
| 50000000 | 18912/1942979 | 41418/4143990 | 85481/5543929 |
| Alphabet Size 16 | | | |
| 25000000 | 8499/860120 | 15320/1518126 | 22347/1449884 |
| 50000000 | 16984/1722707 | 30738/3042700 | 85481/5543929 |
| Alphabet Size 64 | | | |
| 25000000 | 7531/749984 | 14643/1434118 | 22347/1449884 |
| 50000000 | 15056/1502435 | 30375/2996868 | 85481/5543929 |

Table 4.4: *Bulk and total I/Os required for random texts built on alphabets of size* 4, 16 *and* 64.

If we look at the results on alphabet size 64 (see Tables 4.3 and 4.4), we see that the gain of the first stage is much bigger (about 5%), whereas the second stage throws away about 94% of the tuples. Consequently, for increasing alphabet sizes (approaching natural texts), Doubl+Disc gets faster and faster. Apart from the Reuters corpus, which seems indeed a pathological case, we therefore expect a much better performance on natural (and more structured) texts. So that we suggest its use in practice in place of the (plain) Doubling algorithm.

At this point it is worth to notice that the former two stages of Doubl+Disc do not discard any significant number of tuples (like on Reuters). Looking carefully to their algorithmic structure, we observe that these stages execute more work than the one required by the corresponding stages of Doubling because of Step 4 (Section 2.2.1). We experimentally checked this fact by running the Doubling algorithm on the random data sets and verifying an improvement of a factor of two ! Clearly, in the early two stages, the algorithm compares short substrings of length $1-16$, and therefore it is unlikely that those substrings occur only once in a long text (and can be then discarded). Consequently, an insight coming from these experiments is that a *tuned* Doubl+Disc should follow an *hybrid* approach to gain the highest advantage from both the doubling and the discarding strategies: (plain) Doubling executed in the early (e.g. two) stages, Doubl+Disc for the next stages.

### 4.2.3   Results for the L-pieces algorithm.

We again fix $L = 4$ and use the Doubling algorithm to construct $SA_4$. Using bigger texts (up to 50 Mbytes), we are ensured that all the steps of this algorithm operate on disk. The running time distributes as follows: 6% is used to sort short strings, 37% is used to build $SA_4$, and 47% is used to build the other three suffix arrays (via multiway mergesort).[2] Comparing Table 4.1 to Table 4.3, we conclude that the ordering of the short strings remains fast even when it is executed on disk, and it will never be a bottleneck. Moreover, the time required to build $SA_4$ clearly depends on the length of the longest repeated substring (see comments on Doubling), but $\mathcal{A}_{sa}$ is executed on a *compressed* text (of size $N/4$) where that length is reduced by a factor 4. Consequently, the L-pieces algorithm is usually up to 2.9 times faster than Doubl+Disc (see Table 4.3). This speed-up is larger than the one we observed on the Reuters corpus (see Section 4), and thus validates our conjecture that a bigger text collection does not slow down the L-pieces algorithm.

### 4.2.4   Concluding remarks on our experiments

We first compare all experimented algorithms with respect to their time performance for increasing text lengths (see Table 4.1 for a summary). It is obvious from the previous discussions that the MM-algorithm is the fastest one when the suffix array can be built entirely in internal memory. As soon as the working space exceeds the memory size, we can choose among different algorithms depending on the resource to be minimized, either *time* or *space*. The L-pieces algorithm is the obvious choice whenever the splitting of the suffix array does not prevent its use in the underlying application. It is *more than* 3 times faster on the Reuters corpus than any other experimented algorithm; it is *more than* twice faster than the best Doubling variant on random texts. If, instead, a unique suffix array is needed, the choice depends on the structural properties of the text to be indexed. In presence of long repeated substrings, BGS is a good choice till *reasonably large* collections. For *very large* text collections, the hybrid variant of the Doubl+Disc algorithm is definitely worth to be used.

If the space resource is of primary concern, then BGS is a very good choice till reasonably large text collections. For huge sizes, Doubl+Disc+Radix is expected to be better in the light of its asymptotic I/O- and CPU-time complexity. However if one is allowed to keep the suffix array distributed into pieces, then the best construction algorithm results definitely the L-pieces algorithm: It is both the fastest and the cheapest in term of space occupancy (it only needs $6N$ bytes).

We wish to conclude this long discussion on our experimental data and tested algorithms by making a further, and we think necessary, consideration. The *running time* evaluations indicated in the previous tables and pictures are not clearly intended to be definitive. Algorithmic engineering and software tuning of the C++-code might definitely lead to improvements without anyway changing the algorithmic features of the experimented algorithms, and therefore without affecting significantly the scenario that we have depicted in these pages. Consequently, we feel not confident to give an *absolute quantitative* measure of the time performance of these algorithms in order to claim which is the "winner". There are too many system parameters ($M$, buffer size, cache size, memory bandwidth), disk parameters ($B$, seek, latency, bandwidth, cache), and structural properties of the indexed text that heavily influence those times. Nevertheless, the *qualitative* analysis developed in these sections should, in our opinion, safely route and clarify to the software developers which is the algorithm that best fits their wishes and needs.

We conclude our paper by addressing two other issues. The former concerns with the problem of building word-indexes on large text collection; we show in the next section that our results can be successfully applied to this case too without any loss in efficiency and without compromising the ease of programming so to achieve a uniform, simple and efficient approach to both the two indexing models.

---

[2]An attentive reader might observe the the distribution of the time only sums up to 90% of the total construction time. The missing 10% is required to copy back the computed suffix array.

The latter issue is related to the intriguing, and apparently counterintuitive, "contradiction" between the effective practical performance of the BGS-algorithm and its unappealing (i.e., cubic) worst-case behavior. In Section 5.2, we deeply study its algorithmic structure and propose a new approach that follows its *basic philosophy* but in a significantly different manner, thus resulting in a novel algorithm which combines good practical qualities with efficient worst-case performances.

# 5 Extensions

## 5.1 Constructing word-indexes

By using a simple and efficient preprocessing phase, we are also able to build a suffix array on a text where only the beginning of each word is indexed (hereafter called *word-based* suffix array). Our idea is based on a proposal of Andersson *et al.* [5] which was formulated in the context of suffix trees. We greatly simplify their presentation by exploiting the properties of suffix arrays. The preprocessing phase consists of the following steps:

1. **Scan** the text $T$ and define as *index points* the text positions where a non-alphabetic character is followed by an alphabetic one.

2. Form a sequence $S$ of strings which correspond to substrings of $T$ occurring between consecutive *index points*.

3. **Sort** the strings in $S$ via multiway mergesort.

4. Associate with each string its *rank* in the lexicographically ordered sequence, so that given $s_1, s_2 \in S$: if $s_1 = s_2$ then $name(s_1) = name(s_2)$, and if $s_1 <_L s_2$ then $name(s_1) < name(s_2)$.

5. **Sort** (backwards) $S$ according to the starting positions in the original text $T$ of its strings.

6. Create a compressed text $T'$ via a simultaneous **scan** of $T$ and (the sorted) $S$. Here, every substring of $T$ occurring in $S$ is replaced with its *rank*.

The symbols of $T'$ are now *integers* in the range $[1, N]$. It is easy to show that the *word-based* suffix array for $T$ is exactly the same as the suffix array of $T'$. Indeed, let us consider two suffixes $T[i, N]$ and $T[j, N]$ starting at the beginning of a word. They occur also in $T'$ in a *"compressed form"* which preserves their lexicographic order because of the naming process. Consequently, the lexicographic comparison between $T[i, N]$ and $T[j, N]$ is the same as the one among the corresponding compressed suffixes of $T'$.

The cost of the preprocessing phase is dominated by the cost of sorting the string set $S$ (step 2). As the average length of an English word is *six* characters [13], we immediately obtain from [1] (model A, strings shorter than $B$) that the I/O complexity of step 2 is $\Theta(N/B \, \log_{M/B} N/B)$, where $N$ is the total number of string characters. Multiway mergesort is therefore an optimal algorithm to solve the string sorting problem in step 3. In general, we can guarantee that each string of $S$ is always shorter than $B$ characters by introducing some *dummy index points* that split the long substrings of $T$ into *equal-length* shorter pieces. This approach does not suffer from the existence of very long repeated substrings that was reported by the authors of [36] in their quicksort/mergesort–based approaches; and therefore it is expected to work better on very large texts. With respect to [5], our approach does not use *tries* to manage $S$'s strings and thus it does not incur in the very well-known problems related to the management

of unbalanced trees in external memory [29]. Additionally, it reduces the overall working space by making use of arrays. Finally, since the preprocessing phase is based on sorting and scanning routines, we can again infer that this approach scales well on multi-disk and multi-processor machines, as we have largely discussed in the previous sections.

From the experiments executed on the L-pieces algorithm, we know that step 3 and step 5 above will be fast in practice. Furthermore, the compression in step 6 reduces the length of the repeated substrings in $T$, so that Doubl+Disc is expected to require very few stages when applied on $T'$. Consequently, we can expect that constructing word-indexes via suffix arrays is effective in real situations, and can benefit a lot from the study carried out in this paper.

## 5.2 The new BGS-algorithm

The experimental results of the previous sections have lead us to conclude that the BGS-algorithm is attractive for software developers because requires only $4N$ bytes of working space (see footnote 7), it accesses the disk in a sequential manner (thus taking fully advantage from the caching/prefetching strategies of current disks as well as from their high bandwidth [40, 16]), and finally it is very simple to be programmed. However, its worst-case performance is poor and thus its real behavior is not well predictable but heavily depends on the structure of the indexed text. This limits the broad applicability of BGS, making it questionable at theoretical eyes.

In this section, we propose a new algorithm which deploys the *basic philosophy* underlying BGS (i.e., very long disk scans) but in a completely different manner: *the text $T$ is examined from the right to the left*. The algorithm will choreograph this new approach with some additional data structures that allow to perform the suffix comparisons using only the information *available* in internal memory, thus avoiding the *random I/Os* in the worst case. The resulting algorithm still uses small working space (i.e. $8N$ bytes on disk), it is very simple to be programmed, it has small constants, and additionally it achieves effective *worst-case performance* (namely $O(N^2/M^2)$ worst-case bulk I/Os). This makes the practical behavior of the final algorithm *guaranteed on any* indexed text *independently of its structure*, thus overcoming the (theoretical) limitations of the BGS-algorithm, and still keeping its attractive practical properties.

Let us set $m = \ell M$, where $\ell < 1$ is a positive constant to be fixed later in order to guarantee that some auxiliary data structures can be fit into the internal memory. In order to simplify our discussion, let us also assume that $N$ is a multiple of $m$, say $N = km$ for a positive integer $k$.

We divide the text $T$ into $k$ non-overlapping substrings of length $m$ each, namely $T = T_k T_{k-1} \cdots T_2 T_1$ (they are numbered going from the right to the left, thus reflecting the "stage" when the algorithm will process each of them). For the sake of presentation, we also introduce an operator $\diamond$ that allows to go from *absolute positions* of $T$ to *relative positions* in the text pieces constituting $T$. Namely, $x \diamond m = ((x - 1) \bmod m) + 1$: This way, if $T[x]$ lies in the text piece $T_h$ then $T[x] = T_h[x \diamond m]$.

The algorithm executes $\Theta(k) = \Theta(N/M)$ stages (like BGS) and processes the text *from the right to the left* (unlike BGS). The following invariant is kept inductively before stage $h$ starts: $S = T_{h-1} T_{h-2} \cdots T_1$ *is the text part processed in the previous $(h - 1)$ stages. The algorithm has computed and stored on disk the following two data structures: The suffix array $SA_{ext}$ of the string $S$ and its "inverse" array $Pos_{ext}$, which keeps at each entry $Pos_{ext}[j]$ the position in $SA_{ext}$ of the suffix $S[j, |S|]$. (See Figure 5.1.) After all $k$ stages are executed, we have $S = T$ and thus $SA = SA_{ext}$.*

The main idea underlying the *leftward*-scanning of the text is that when the $h$-th stage processes the text suffixes starting in $T_h$, it has already accumulated into $SA_{ext}$ and $Pos_{ext}$ some informations about the text suffixes starting to the right of $T_h$. This way, the comparison of the former text suffixes will eventually exploit these two arrays, and thus use only *localized* information which eliminates the need of *random I/Os*. The next Lemma formalizes this intuition:

**Lemma 1** *A text suffix $T[i, N]$ starting into the substring $T_h$ can be represented succinctly via the pair $(T[i, i+m-1], Pos_{ext}[(i+m) \diamond m])$. Consequently, all the text suffixes starting in $T_h$ can be represented using overall $O(m)$ space.*

**Proof:** Text suffix $T[i, N]$ can be seen as the concatenation of two strings $T[i, i+m-1]$ and $T[i+m, N]$, where the second string is an arbitrarily long text suffix. The position $i+m$ occurs in $T_{h-1} \cdots T_1 (= S)$ and in particular $T[i+m, n] = S[(i+m) \diamond m, |S|]$. This string can be succinctly represented with a number which denotes its (lexicographic) position among $S$'s suffixes (i.e., its position in $SA_{ext}$). This number is $Pos_{ext}[(i+m) \diamond m]$. The space-bound easily holds by storing in internal memory the text substring $T_h \cdot T_{h-1}$ and the array $Pos_{ext}[2, m+1]$. ∎

Stage $h$ preserves the invariant above and thus updates $SA_{ext}$ and $Pos_{ext}$ by *properly* inserting into them the "information" regarding the text suffixes of $T$ which start in $T_h$. This way, the new $SA_{ext}$ and $Pos_{ext}$ will correctly refer to the "extended" string $T_h \cdot S$ ($T_h$ concatenated with $S$), thus preserving the invariant for the next $(h+1)$th stage (where $S = T_h \cdot S = T_h T_{h-1} \cdots T_2 T_1$). Details on the stage $h$ follow (see Figure 5.1 below).



$$SA_{ext} = [\,8, 2, 4, 1, 3, 5, 7, 6\,]$$
$$Pos_{ext} = [\,4, 2, 5, 3, 6, 8, 7, 1\,]$$
$\left.\right\}$ refer to string $S$

$$SA^* = [\,2, 4, 3, 1\,]$$
$$Pos^* = [\,4, 1, 3, 2\,]$$
$\left.\right\}$ refer to $T^*[1,m]$

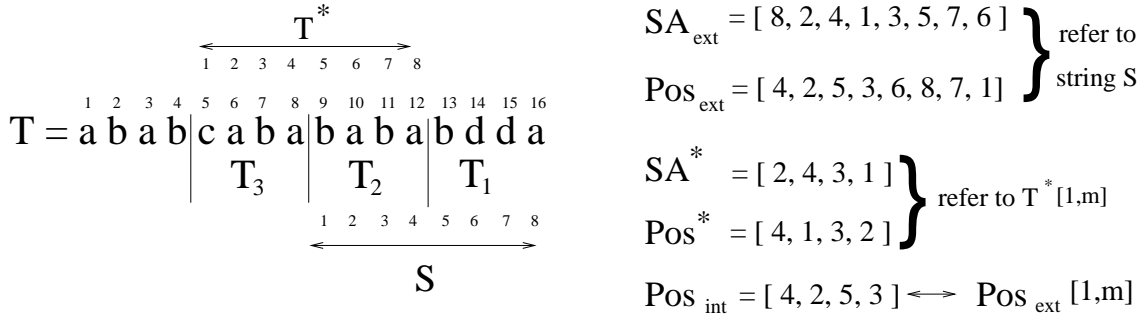$$Pos_{int} = [\,4, 2, 5, 3\,] \longleftrightarrow Pos_{ext}[1,m]$$

Figure 5.1: *The figure depicts the data structures used during stage $h = 3$, where $T^* = T_3 \cdot T_2$ and $S = T_2 \cdot T_1$. $SA^*$ contains only the first $m$ suffixes of $T^*$. Notice the order in $SA^*$ of the two text suffixes starting at positions 2 and 4 of $T^*$ (i.e., 6 and 8 of $T$). Restricted to their prefixes lying into $T^*$, these two suffixes satisfy the relation $4 \leq_L 2$, but considering them in their entirety (till the end of $T$), it is $2 \leq_L 4$. We can represent compactly $T[6, 16]$ via the pair $\langle$ 'abab', 2$\rangle$ and $T[8, 16]$ via the pair $\langle$ 'abab', 3$\rangle$; hence the comparison of those pairs gives the correct order and can be executed in internal memory (Lemma 1).*

1. Load $T_h$ and $T_{h-1}$ in the internal memory and set $T^* = T_h \cdot T_{h-1}$. ($O(1)$ bulk I/Os.)

2. Load the first $m$ entries of $Pos_{ext}$ into the array $Pos_{int}$. ($O(1)$ bulk I/Os.)

3. Construct the suffix array $SA^*$ which contains the lexicographically ordered sequence of the text suffixes starting into $T_h$ (recall that they extend till the end of $T$). This is done by means of three substeps which exploit the information kept in internal memory:

   (a) Build the suffix array of the string $T^*$ using any internal memory algorithm (e.g. [32]). Then, throw away all suffixes of $T^*$ which start in its second half. (No I/Os are required.)

   (b) Store the remaining $m$ entries into $SA^*$.[1] (No I/Os are required.)

---

[1]Notice that this is not yet the correct $SA^*$ because there might exist two text suffixes which start in $T^*$ but have a common prefix which extends outside $T^*$. The next Step 3c does this computation without accessing the disk !

(c) Refine the order in $SA^*$ taking into account the suffixes in their entirety (i.e., considering also their characters outside $T^*$). This is done as follows. Let $T^*[x, 2m]$ and $T^*[y, 2m]$ be two suffixes which lie adjacent into the current array $SA^*$, namely $SA^*[j] = x$ and $SA^*[j+1] = y$ for some value $j$, and such that one of them is the *prefix* of the other. Their order in $SA^*$ may possibly be not correct (see Figure 5.1). Hence, the correct order is computed by comparing the two pairs $\langle T^*[x, x+m-1], Pos_{int}[(x+m) \diamond m] \rangle$ and $T^*[y, y+m-1], Pos_{int}[(y+m) \diamond m]$ (see Lemma 1). [2] This comparison is done for all the suffixes of $T^*$ for which the ambiguity above arises. (No I/Os are required.)

4. Scan simultaneously the string $S$ and the array $Pos_{ext}$ ($O(h)$ bulk I/Os in total). For each suffix $S[j, |S|]$ do the following substeps:

   (a) Via a binary search, find the lexicographic position $pos(j)$ of that suffix into $SA^*$ so that $S[j, |S|]$ follows the suffix of $T$ starting at position $SA^*[pos(j) - 1]$ of $T_h$ and precedes the suffix of $T$ starting at position $SA^*[pos(j)]$ of $T_h$. Lemma 1 allows to perform the suffix comparisons of the binary search using only the internal memory.

   (b) Increment $C[pos(j)]$ by one unit.

   (c) Update the entry $Pos_{ext}[j] = Pos_{ext}[j] + pos(j) - 1$.

5. Build the new array $SA_{ext}[1, hm]$ by merging the old (external and shorter) array $SA_{ext}$ with the (internal) array $SA^*$ by means of the information available into $C[1, m+1]$. This requires a single disk scan (like BGS) during which the algorithm also executes the computation $SA_{ext}[j] = SA_{ext}[j] + m$, in order to take into account the fact that in the next $(h+1)$th stage the new string $S$ has appended in front the text piece $T_h$. (Globally $O(h)$ bulk I/Os.)

6. Process in internal memory the array $C$ as follows (no I/Os are executed):

   (a) Compute the Prefix-Sum of $C$.

   (b) Set $C[j] = C[j] + j$, for $j = 1, \ldots, m+1$.

   (c) Compute $Pos^*[1, m]$ as the inverse of $SA^*$, and then permute $C[1, m]$ according to $Pos^*[1, m]$. Namely, $C[i] = C[Pos^*[i]]$, simultaneously for all $i = 1, 2, \ldots, m$.

7. Build the new array $Pos_{ext}[1, hm]$ by appending $C[1, m]$ to the front of the current $Pos_{ext}$. Namely, $Pos_{ext} = C[1, m] \cdot Pos_{ext}$. ($O(h)$ bulk I/Os.)

The correctness of the algorithm immediately comes from Lemma 1 and from the following observations. As far as Step 6 is concerned, we observe that:

- Step 6a determines for each entry $C[j]$ the number of text suffixes which start in $S$ and are *lexicographically smaller* than the text suffix starting at position $SA^*[j]$ of $T_h$.

- Step 6b keeps into account also the number of suffixes which start in $T_h$ and are *lexicographically smaller* than the text suffix starting at position $SA^*[j]$ of $T_h$. These suffixes are $(j-1)$, so that the algorithm sums $j$ to compute the final rank of that suffix (i.e., the one starting at $T_h[SA^*[j]]$) among all the suffixes of the string $T_h \cdot S$.

- Step 6c permutes the array $C$ so that $C[j]$ gives the rank of the text suffix starting at $T_h[j]$ among all suffixes of the string $T_h \cdot S = T_h \cdots T_1$.

---

[2]This step is executed by employing the extra-space available into $C$ and not yet used.

Since the string $T_h \cdots T_1$ corresponds to the string $S$ of the next $(h+1)$th stage, we can conclude that Step 6 correctly stores in $C$ the first $m$ entries of the new array $Pos_{ext}$ (Step 7). Finally, we point out that Step 4c correctly updates the entries of $Pos_{ext}[1, (h-1)m]$ by taking into account the insertion in $SA_{ext}$ of the $m$ text suffixes starting in $T_h$; and Step 5 correctly updates the entries of $SA_{ext}[1, (h-1)m]$ by taking into account the insertion in front of $S$ of the $m$ suffixes starting in $T_h$. In summary we can state the following result,

**Theorem 1** *The suffix array of a text $T[1, N]$ can be constructed in $O(N^2/M^2)$ bulk-I/Os, no random-I/Os, and $8N$ disk space in the worst case. The overall CPU time is $O(\frac{N^2}{M} \log_2 M)$.*

The above parameter $\ell$ is set to fit $SA^*$, $Pos^*$, $Pos_{int}$, $T^*$ and $C$ into internal memory (notice that some space can be reused). We remark that the algorithm requires $4N$ bytes *more than* the space-optimized variant of the BGS-algorithm (see Section 2.1.2). Nonetheless, we can still get the $4N$ space-bound if we accept to compute only the array $Pos_{ext} = SA^{-1}$ (implicit $SA$). In any case, the overall working space is *much less* than the one required by all Doubling variants, and it is *exactly equal* to the one required by our implementation of the BGS-algorithm. The new BGS-algorithm is also very simple to be programmed and has small constants (hidden in the big-Oh notation). Therefore, it preserves the good algorithmic properties of BGS, but it now guarantees good worst-case performances.

# 6 Conclusions

It is often observed that practitioners use algorithms which tend to be different from what is claimed as optimal by theoreticians. This is doubtless because theoretical models tend to be simplifications of reality, and theoretical analysis needs to use conservative assumptions. Our paper provided to some extent an example of this phenomenon—apparently "bad" theoretical algorithms are good in practice (see BGS and new-BGS). In the present paper we actually tried to "bridge" this difference by analyzing more deeply some suffix-array construction algorithms via the new accounting scheme of [19], taking more into account the specialties of current disk systems. This has lead us to a reasonable and significant taxonomy of all these algorithms. As it appears clear from the experiments, the final choice of the "best" algorithm depends on the available disk space, on the disk characteristics (which induce different trade-offs between random and bulk I/Os), on the structural features of the indexed text, and also on the patience of the user to wait for the completion of the suffix-array construction. The moral we draw from our experiments is that the design of an external-memory algorithm must take more and more into account the current technological trends [16], which boost interest toward the development of algorithms which *prefer* bulk rather than random I/Os because this paradigm can take advantage of the large disk-bandwidth and the high computational power of current computer systems.

The description of all algorithms discussed in this paper may have been much detailed, but our aim has been to make this paper as much *self-contained and readable* as possible, without leaving "dangling pointers" to the literature in order to additionally offer a unique reference for anyone who is interested in building large suffix arrays. The study and experiments so far conducted leads us to indicate some other directions of research that in our opinion deserve further investigation:

- Radix Heaps were chosen to reduce the space requirements of the doubling algorithm. Recently, in [14] has been presented a novel external-memory heap mainly based on an array topology, which requires strictly linear space and an optimal number of I/Os. The performance of this heap does not depend on the value of the priorities stored into it, so that it has the space-advantages of the radix-heap but it overcomes its drawbacks. Additionally, since the priority data type is not restricted to integers, this heap may allow us to perform the Step 1 of Doubling (see Section 2.2.1) *in one shot* without increasing the overall work (as instead happens with radix-heaps). Therefore, it would be interesting to plug this new data structure into the Doubl+Disc algorithm and evaluate its time vs. space tradeoff.

- In [50], it is discussed a variant of the multiway mergesort which uses $NX + (N/B)\epsilon$ bytes to sort $N$ items of size $X$ bytes each ($\epsilon > 0$). This approach is I/O-optimal but the pattern of disk accesses is distributed randomly so that it is unclear how it might behave in real situations. It is therefore interesting to evaluate how this approach compares with the sorting routine above based on array-heaps [14].

- The algorithm presented in Section 2.2.3 is very efficient both in terms of I/Os and working space, but it produces $L = 4$ distinct suffix arrays. It would be interesting to establish how much it costs

in practice to *merge* these four arrays into one unique suffix array. In this respect, it would be worth to use the *lazy-trie* data structure proposed in [1]. Although this approach would be not I/O-efficient from a theoretical point of view (it may take more than $N$ I/Os in the worst case), the properties of the lazy trie and the limited (i.e., four) number of sequences to be merged, seem to offer an hope to achieve good performance in practice.

- In Section 5.1 we proposed a method for constructing word-based suffix arrays and conjectured its good practical performance based on the promising experimental results shown in Section 4. In the near future, we would like to validate this conjecture by performing some experimental tests.

- Recently, Farach *et al.* [19] devised the first I/O-optimal algorithm for building *suffix trees*. Although asymptotically optimal (both in time and space), that algorithm uses more space than the algorithms discussed in this paper because it operates on a tree topology. It is not yet clear how this approach could be used to *directly* construct suffix arrays. This is an important problem both theoretically and experimentally.

- Our new algorithms scale theoretically well in a distributed environment where $D$ disks or $P$ workstations are available. However, there are only few practical experiments in this setting [36]. We hope in the future to extend our results to this high-performance environment.

# Bibliography

[1] L. Arge, P. Ferragina, R. Grossi and J. S. Vitter. On sorting Strings in External Memory. *ACM Symposium on Theory of Computing*, pp. 540-548, 1997.

[2] A. Aggarwal, A. Chandra, and M. Snir. Hierarchical memory with block transfer. *IEEE Symposium on Foundations of Computer Science*, pp. 204–216, 1987.

[3] R. Ahuja, K. Mehlhorn, J. B. Orlin and R. E. Tarjan. Faster Algorithms for the Shortest Path Problem. *Journal of the ACM(2)*, pp. 213-223, 1990.

[4] B. Alpern, L. Carter, and E. Feig. Uniform memory hierarchies. *IEEE Symposium on Foundations of Computer Science*, pp. 600–609, 1990.

[5] A. Andersson, N. J. Larson and K. Swanson. Suffix Trees on Words. *Combinatorial Pattern Matching Conference*, LNCS 1075, pp. 102-115, 1996.

[6] A. Andersson and S. Nilsson. Efficient implementation of Suffix Trees. *Software Practice and Experience*, 2(25): 129-141, 1995.

[7] A. Apostolico. The myriad virtues of subword trees. In *Combinatorial Algorithms on Words*, A. Apostolico anf Z. Galil, Eds., NATO ASI Series F: Computer and System Sciences, Springer Verlag, pp. 85-96, 1985.

[8] R. Baeza-Yates, E. F. Barbosa and N. Ziviani. Hierarchies of Indices for Text Searching. *Information Systems*, 21(6):497-514, 1996.

[9] P. Bieganski, J. Riedl, J. V. Carlis and E. F. Retzel. Generalized suffix trees for biological sequence data: applications and implementation. *Hawaii International Conference on System Sciences*, pp. 35–44, IEEE press, 1994.

[10] S. Burkhard, A. Crauser, P. Ferragina, H. Lenhof, E. Rivals and M. Vingron. *q*-gram based database searching unsing a suffix array (QUASAR). *International Conference on Computational Molecular Biology*, 1999. See also *Tech. Rep. MPI-I-98-1-024*, Oct. 1998, MPI für Informatik, Saarbruücken, Germany.

[11] M. Burrows and D. J. Wheeler. A block-sorting lossless data compression algorithm. *TR 124, Digital SRC Research Report*, 1994.

[12] P.M. Chen, E.K. Lee, G.A. Gibson, R.H. Katz, and D.A. Patterson. RAID: high-performance, reliable secondary storage. *ACM Computing Surveys*, 26(2):145–185, 1994.

[13] D. R. Clark and J. I. Munro. Efficient Suffix Trees on Secondary Storage. *ACM-SIAM Symposium on Discrete Algorithms*, pp. 383-391, 1996.

[14] A. Crauser, P. Ferragina and U. Meyer. Practical and Efficient Priority Queues for External Memory. Technical Report MPI, see WEB pages of the authors.

[15] A. Crauser and K. Mehlhorn. LEDA-SM: A Library Prototype for Computing in Secondary Memory. Technical Report MPI, http://www.mpi-sb.mpg.de/c̆rauser/leda-sm.html.

[16] M. Dahlin. Trends in disk technology. http://www.cs.utexas.edu/users/dahlin/techTrends/, 1998.

[17] C. Faloutsos. Access Methods for text. *ACM Computing Surveys*, 17, pp.49-74, 1985.

[18] M. Farach. Optimal suffix tree construction with large alphabets. *IEEE Symposium on Foundations of Computer Science*, pp. 137–143, 1997.

[19] M. Farach, P. Ferragina and S. Muthukrishnan. Overcoming the memory bottleneck in Suffix Tree construction. *IEEE Symposium on Foundations of Computer Science*, 1998.

[20] C.L. Feng. PAT-Tree-Based Keyword Extraction for Chinese Information Retrieval. *ACM SIGIR*, pp. 50-58, 1997.

[21] P. Ferragina and R. Grossi. A Fully-Dynamic Data Structure for External Substring Search. *ACM Symposium on Theory of Computing*, pp. 693-702, 1995. Also *Journal of the ACM* (to appear).

[22] E. A. Fox, R. M. Akcyn, R. K. Furuta and J. J. Leggett. Special Issue: Digital libraries. *Communications of the ACM*, 38(4), 1995.

[23] T. Goodrich. Communication-efficient parallel sorting. *ACM Symposium on Theory of Computing*, pp. 247–256, 1996.

[24] G. H. Gonnet, R. A. Baeza-Yates and T. Snider. New indices for text:PAT trees and PAT arrays. In *Information Retrieval – Data Structures and Algorithms*, W.B. Frakes and R. BaezaYates Editors, pp. 66-82, Prentice-Hall, 1992.

[25] D. Gusfield. *Algorithms on strings, trees and sequences: computer science and computational biology*, Cambridge University Press, 1997.

[26] J. Kärkkäinen. Suffix Cactus: A Cross between Suffix Tree and Suffix Array. In *Combinatorial Pattern Matching*, LNCS 937, 191-204, 1995.

[27] R. M. Karp, R. E. Miller and A. L. Rosenberg. Rapid identification of repeated patterns in strings, arrays and trees. *ACM Symposium on Theory of Computing*, pp. 125-136, 1972.

[28] J. P. Kitajima, G. Navarro, B. Ribeiro-Neto and N. Ziviani, Distributed Generation of Suffix Arrays: A Quicksort Based Approach. *South American Workshop on String Processing*, Carleton University Press, 53–69, 1997.

[29] D. E. Knuth. *The Art of Computer Programming: Sorting and Searching.* Vol. 3, Addison-Wesley Publishing Co. 1969.

[30] S. Kurtz. Reducing the Space Requirement of Suffix Trees. *Technical Report 98-03*, University of Bielefeld, 1998.

[31] N. J. Larsson. Extended application of suffix trees to data compression. *IEEE Data Compression Conference*, 1996.

[32] U. Manber and G. Myers. Suffix arrays: a new method for on-line string searches. *SIAM Journal of Computing 22*, 5,pp. 935-948, 1993.

[33] E. M. McCreight. A space-economical suffix tree construction algorithm. *Journal of the ACM* 23(2):262-272, 1976.

[34] H. W. Mewes and K. Heumann. Genome Analysis: Pattern Search in Biological Macromolecules. *Combinatorial Pattern Matching Conference*, LNCS 937, pp. 261-285, 1995.

[35] I. Munro, V. Raman and S. S. Rao. Space efficient suffix trees. *Foundations of Software Technology and Theoretical Computer Science*, 1998.

[36] G. Navarro, J.P. Kitajima, B.A. Ribeiro-Neto and N. Ziviani. Distributed Generation of Suffix Arrays. *Combinatorial Pattern Matching Conference*, pp. 103–115, 1997.

[37] S. Näher and K. Mehlhorn. LEDA: A Platform for Combiantorial and Geometric Computing. *Communications of the ACM* 38, 1995.

[38] M. H. Nodine and J. S. Vitter. Greed sort: An optimal sorting algorithm for multiple disks. *Journal of the ACM*, pp. 919-933, 1995.

[39] Y. N. Patt. The I/O subsystem: a candidate for improvement. *Guest Editor's Introduction in IEEE Computer*, 27(3), pp. 15-16, 1994.

[40] C. Ruemmler and J. Wilkes. An introduction to disk drive modeling. *IEEE Computer*, 27(3):17–29, 1994.

[41] K. Sadakane. A Fast Algorithm for Making Suffix Arrays and for Burrows-Wheeler Transformation. *IEEE Data Compression Conference*, 1998.

[42] G. Salton and M. J. McGill. *Introduction to Modern Information Retrieval.* McGraw-Hill, New York, 1983.

[43] E. Shriver. *Performance modeling for realistic storage devices.* PhD thesis, Department of Computer Science, New York University, New York, NY, 1997.

[44] E. A. Shriver and J. S. Vitter. Algorithms for parallel memory I: two-level memories. *Algorithmica*, 12(2-3):110-147, 1994.

[45] W. R. Stevens. *Advanced Programming in the UNIX Environment.* Addison-Wesley Professional Computing Series, 1992.

[46] D. E. Vengroff and J. S. Vitter. I/O-efficient scientific computing using TPIE. *IEEE Symposium on Parallel and Distributed Computing*, 1995.

[47] J. Vitter. External memory algorithms. Invited Tutorial, *ACM Symposium on Principles of Database Systems (PODS '98)*, Seattle, WA, 1998. Also Invited Paper in *European Symposium on Algorithms (ESA '98)*, Venice, Italy, 1998.

[48] J. Zobel, A. Moffat and K. Ramamohanarao. Guidelines for presentation and comparison of indexing techniques. *SIGMAD Record 25*, 3:10–15, 1996.

[49] P. Weiner. Linear pattern matching algorithm. *IEEE Symposium on Switching and Automata Theory*, pp. 1-11, 1973.

[50] I. H. Witten, A. Moffat, and T. C. Bell. *Managing Gigabytes.* Van Nostrand Reinhold, 115th Avenue, New York, 1994.

[51] S. Wu and U. Manber. GLIMPSE: A Tool to Search Through Entire File Systems. TR 93-34, Dept. of Computer Science, University of Arizona, 1993.

Below you find a list of the most recent technical reports of the Max-Planck-Institut für Informatik. They are available by anonymous ftp from `ftp.mpi-sb.mpg.de` under the directory `pub/papers/reports`. Most of the reports are also accessible via WWW using the URL `http://www.mpi-sb.mpg.de`. If you have any questions concerning ftp or WWW access, please contact `reports@mpi-sb.mpg.de`. Paper copies (which are not necessarily free of charge) can be ordered either by regular mail or by e-mail at the address below.

Max-Planck-Institut für Informatik
Library
attn. Birgit Hofmann
Im Stadtwald
D-66123 Saarbrücken
GERMANY
e-mail: `library@mpi-sb.mpg.de`

| | | |
|---|---|---|
| MPI-I-1999-2-003 | U. Waldmann | Cancellative Superposition Decides the Theory of Divisible Torsion-Free Abelian Groups |
| MPI-I-1999-2-001 | W. Charatonik | Automata on DAG Representations of Finite Trees |
| MPI-I-98-2-018 | F. Eisenbrand | A Note on the Membership Problem for the First Elementary Closure of a Polyhedron |
| MPI-I-98-2-017 | M. Tzakova, P. Blackburn | Hybridizing Concept Languages |
| MPI-I-98-2-014 | Y. Gurevich, M. Veanes | Partisan Corroboration, and Shifted Pairing |
| MPI-I-98-2-013 | H. Ganzinger, F. Jacquemard, M. Veanes | Rigid Reachability |
| MPI-I-98-2-012 | G. Delzanno, A. Podelski | Model Checking Infinite-state Systems in CLP |
| MPI-I-98-2-011 | A. Degtyarev, A. Voronkov | Equality Reasoning in Sequent-Based Calculi |
| MPI-I-98-2-010 | S. Ramangalahy | Strategies for Conformance Testing |
| MPI-I-98-2-009 | S. Vorobyov | The Undecidability of the First-Order Theories of One Step Rewriting in Linear Canonical Systems |
| MPI-I-98-2-008 | S. Vorobyov | AE-Equational theory of context unification is Co-RE-Hard |
| MPI-I-98-2-007 | S. Vorobyov | The Most Nonelementary Theory (A Direct Lower Bound Proof) |
| MPI-I-98-2-006 | P. Blackburn, M. Tzakova | Hybrid Languages and Temporal Logic |
| MPI-I-98-2-005 | M. Veanes | The Relation Between Second-Order Unification and Simultaneous Rigid $E$-Unification |
| MPI-I-98-2-004 | S. Vorobyov | Satisfiability of Functional+Record Subtype Constraints is NP-Hard |
| MPI-I-98-2-003 | R.A. Schmidt | E-Unification for Subsystems of S4 |
| MPI-I-98-2-002 | F. Jaquemard, C. Meyer, C. Weidenbach | Unification in Extensions of Shallow Equational Theories |
| MPI-I-98-1-031 | G.W. Klau, P. Mutzel | Optimal Compaction of Orthogonal Grid Drawings |
| MPI-I-98-1-030 | H. Brönniman, L. Kettner, S. Schirra, R. Veltkamp | Applications of the Generic Programming Paradigm in the Design of CGAL |
| MPI-I-98-1-029 | P. Mutzel, R. Weiskircher | Optimizing Over All Combinatorial Embeddings of a Planar Graph |
| MPI-I-98-1-028 | A. Crauser, K. Mehlhorn, E. Althaus, K. Brengel, T. Buchheit, J. Keller, H. Krone, O. Lambert, R. Schulte, S. Thiel, M. Westphal, R. Wirth | On the performance of LEDA-SM |
| MPI-I-98-1-027 | C. Burnikel | Delaunay Graphs by Divide and Conquer |

| | | |
|---|---|---|
| MPI-I-98-1-026 | K. Jansen, L. Porkolab | Improved Approximation Schemes for Scheduling Unrelated Parallel Machines |
| MPI-I-98-1-025 | K. Jansen, L. Porkolab | Linear-time Approximation Schemes for Scheduling Malleable Parallel Tasks |
| MPI-I-98-1-024 | S. Burkhardt, A. Crauser, P. Ferragina, H. Lenhof, E. Rivals, M. Vingron | $q$-gram Based Database Searching Using a Suffix Array (QUASAR) |
| MPI-I-98-1-023 | C. Burnikel | Rational Points on Circles |
| MPI-I-98-1-022 | C. Burnikel, J. Ziegler | Fast Recursive Division |
| MPI-I-98-1-021 | S. Albers, G. Schmidt | Scheduling with Unexpected Machine Breakdowns |
| MPI-I-98-1-020 | C. Rüb | On Wallace's Method for the Generation of Normal Variates |
| MPI-I-98-1-019 | | 2nd Workshop on Algorithm Engineering WAE '98 - Proceedings |
| MPI-I-98-1-018 | D. Dubhashi, D. Ranjan | On Positive Influence and Negative Dependence |
| MPI-I-98-1-017 | A. Crauser, P. Ferragina, K. Mehlhorn, U. Meyer, E. Ramos | Randomized External-Memory Algorithms for Some Geometric Problems |
| MPI-I-98-1-016 | P. Krysta, K. Loryś | New Approximation Algorithms for the Achromatic Number |
| MPI-I-98-1-015 | M.R. Henzinger, S. Leonardi | Scheduling Multicasts on Unit-Capacity Trees and Meshes |
| MPI-I-98-1-014 | U. Meyer, J.F. Sibeyn | Time-Independent Gossiping on Full-Port Tori |
| MPI-I-98-1-013 | G.W. Klau, P. Mutzel | Quasi-Orthogonal Drawing of Planar Graphs |
| MPI-I-98-1-012 | S. Mahajan, E.A. Ramos, K.V. Subrahmanyam | Solving some discrepancy problems in NC* |
| MPI-I-98-1-011 | G.N. Frederickson, R. Solis-Oba | Robustness analysis in combinatorial optimization |
| MPI-I-98-1-010 | R. Solis-Oba | 2-Approximation algorithm for finding a spanning tree with maximum number of leaves |
| MPI-I-98-1-009 | D. Frigioni, A. Marchetti-Spaccamela, U. Nanni | Fully dynamic shortest paths and negative cycle detection on diagraphs with Arbitrary Arc Weights |
| MPI-I-98-1-008 | M. Jünger, S. Leipert, P. Mutzel | A Note on Computing a Maximal Planar Subgraph using PQ-Trees |
| MPI-I-98-1-007 | A. Fabri, G. Giezeman, L. Kettner, S. Schirra, S. Schönherr | On the Design of CGAL, the Computational Geometry Algorithms Library |
| MPI-I-98-1-006 | K. Jansen | A new characterization for parity graphs and a coloring problem with costs |
| MPI-I-98-1-005 | K. Jansen | The mutual exclusion scheduling problem for permutation and comparability graphs |
| MPI-I-98-1-004 | S. Schirra | Robustness and Precision Issues in Geometric Computation |
| MPI-I-98-1-003 | S. Schirra | Parameterized Implementations of Classical Planar Convex Hull Algorithms and Extreme Point Compuations |
| MPI-I-98-1-002 | G.S. Brodal, M.C. Pinotti | Comparator Networks for Binary Heap Construction |
| MPI-I-98-1-001 | T. Hagerup | Simpler and Faster Static $AC^0$ Dictionaries |
| MPI-I-97-2-012 | L. Bachmair, H. Ganzinger, A. Voronkov | Elimination of Equality via Transformation with Ordering Constraints |
| MPI-I-97-2-011 | L. Bachmair, H. Ganzinger | Strict Basic Superposition and Chaining |
| MPI-I-97-2-010 | S. Vorobyov, A. Voronkov | Complexity of Nonrecursive Logic Programs with Complex Values |
| MPI-I-97-2-009 | A. Bockmayr, F. Eisenbrand | On the Chvátal Rank of Polytopes in the 0/1 Cube |
| MPI-I-97-2-008 | A. Bockmayr, T. Kasper | A Unifying Framework for Integer and Finite Domain Constraint Programming |
| MPI-I-97-2-007 | P. Blackburn, M. Tzakova | Two Hybrid Logics |
| MPI-I-97-2-006 | S. Vorobyov | Third-order matching in $\lambda \rightarrow$-Curry is undecidable |