

A Linear Lower Bound on Index Size for Text Retrieval

Erik D. Demaine* Alejandro López-Ortiz†

Abstract

Most information-retrieval systems preprocess the data to produce an auxiliary index structure. Empirically, it has been observed that there is a tradeoff between query response time and the size of the index. When indexing a large corpus, such as the web, the size of the index is an important consideration. In this case it would be ideal to produce an index that is substantially smaller than the text.

In this work we prove a linear worst-case lower bound on the size of any index that reports the location (if any) of a substring in the text in time proportional to the length of the pattern. In other words, an index supporting linear-time substring searches requires about as much space as the original text. Here “time” is measured in the number of bit probes to the text; an arbitrary amount of computation may be done on an arbitrary amount of the index. Our lower bound applies to inverted word indices as well.

1 Introduction

Text retrieval is crucial in such contexts as searching the web, news, and medical databases. The most basic problem, used as a subroutine in most search engines, is the string-matching problem of searching for a given substring (keyword or phrase) in a corpus of text. Because the text database changes infrequently relative to the frequency and abundance of queries, fundamental to any search technique is a preprocessing step to prepare an *index* for fast searches.

It has been observed empirically that the larger an index the better the query time. To illustrate with an extreme example, in the absence of an index, it is necessary to examine the entire text to see if the query string is present. For example, the Knuth-Morris-Pratt algorithm [KMP77] requires no index and runs in time proportional to the length of the text plus the pattern. In practice such a search is done often with the UNIX utility `grep`.

On the other end of the spectrum, a query can be answered in time that is linear in the length of the pattern and output size. This bound is obtained by the popular *suffix tree* data structure [Wei73, McC76, Ukk95, GK97]. The index consists of $O(N)$ words, where N is the

*MIT Laboratory for Computer Science, 200 Technology Square, Cambridge, MA 02139, USA, edemaine@mit.edu.

†Department of Computer Science, University of Waterloo, Waterloo, Ontario N2L 3G1, Canada, alopez-o@uwaterloo.ca. Supported by NSERC.

number of bits in the text. Recently, Grossi and Vitter [GV00] have shown that the space can be improved to $O(N)$ bits (proportional to the size of the text), and furthermore that the search time can be improved by a factor of $\lg N$ by exploiting a machine word at least $\lg N$ bits long. Unfortunately, this structure also requires $O(\lg^\epsilon N)$ time for each matching position that should be returned. Thus, the total cost of searching for a pattern P in a text of length N using their structure is $O(|P|/\lg N + |\text{output}|\lg^\epsilon N)$ time.

Related structures. Several other text-indexing structures have been developed. Ferragina and Manzini [FM00] proposed a structure with a similar search time, $O(|P| + |\text{output}|\lg^\epsilon N)$, but whose space requirement is related to the entropy of the text, being smaller for compressible text. Their structure also supports counting the number of occurrences of a pattern in $O(|P|)$ time, but to find even a single occurrence of the pattern still requires $\Theta(\lg^\epsilon N)$ time. Recently, Ferragina and Manzini [FM02] have also shown how to remove the $\lg^\epsilon N$ factor from the search time at the cost of adding a factor of $\lg^\epsilon N$ to the space. Sadakane [Sad00] designed a structure that applies also to large alphabets, but increases the search cost to $O(|P|\lg N + |\text{output}|\lg^\epsilon N)$. All of these structures essentially encode the entire text, and hence in the worst case also require $\Omega(N)$ bits of space. Sadakane [Sad02] showed how to extend the last four structures [GV00, FM00, FM02, Sad00] to support additional queries—finding the longest common prefix—using $O(N)$ extra bits of space.

For weaker operations, better bounds are known. In particular, static dictionaries support just membership queries: does a given word (substring) appear in the text? Fredmen, Komlós, and Szemerédi [FKS84] present a structure supporting constant-time queries and using $w \lg N + o(w \lg N)$ bits of space for w query words of interest. Pagh [Pag01] improved this space bound to $w \lg(N/w)$ bits. For $w = o(N)$ query words of interest, this index is smaller than the text of N bits.

Our results. An important but relatively unstudied field of research is to determine the asymptotic tradeoff between fast queries and a small index. As time is currently the most important issue in most systems, a natural problem in this field is to determine the minimum amount of space required by an index supporting substring searches in time proportional to the substring length.

In this paper, we prove the first nontrivial lower bound on this problem. Specifically, we show that any such indexing structure requires space proportional to the text itself. This bound applies in a powerful bit-probe model in which the search algorithm has free access to the *entire* index and can perform unlimited computation—only probes to the text are charged. The bound holds whenever the number of allowed probes is $O(|P| + \lg |T|)$, where $|P|$ denotes the length of the search substring and $|T|$ denotes the length of the text (measured in bits). We also show that the same lower bound holds in the practical case of *inverted word indices*, where the text is broken into $|T|/w$ allowed query “words” each of length w . (Here, $w = \lg |T| + o(\lg |T|)$.)

1.1 Model of Computation

More formally, we consider the problem of determining one of the locations, if any, of a given query string P (the *pattern*) in the text string T , in time bounded in terms of the length of the pattern, $|P|$, using an index I . The index is a static structure precomputed before query time, and for the purposes of lower bounds we do not consider the preprocessing time needed to build I . The search algorithm has access to P , T , and I , and the standard model is that the search algorithm can perform at most $\mathcal{S} = \mathcal{S}(|P|, |T|)$ total internal computations and accesses (probes) to P , T , and I . We impose one restriction to this model, that each probe to the text T retrieves only a single bit from T .

Our lower bound applies to a much stronger model. The search algorithm is allowed an unlimited amount of computation, and can read the entire index I and pattern P ; we only count the number of bit probes made to the text T , which must be at most \mathcal{S} . This model is stronger than the powerful bit-probe model [BMRS00, Mil93, Yao81], which counts all bit probes to I , P , and T . The model is also stronger than a RAM in which probes to the text are restricted to bit probes. Our lower bounds thus apply to these models as well.

2 Main Result

In this section we prove the following theorem under the preceding model of computation:

Theorem 2.1 *If I is an index supporting a search for the substring P in $\mathcal{S}(|P|, |T|) = O(|P|)$ bit probes to the text T , then $|I| = \Omega(|T|)$.*

This bound applies even when all query substrings P have length $\lg |T| + o(\lg |T|)$, and do not overlap. Thus, the bound also holds when $\mathcal{S}(|P|, |T|) = O(|P| + \lg |T|)$, and for inverted word indices.

See Theorem 2.2 in Section 2.6 for a precise statement of the constant factor.

2.1 Construction

Recall that the Kolmogorov complexity of an average element from a class of k objects is within an additive constant of the information theoretic lower bound, $\lg k$. For details on Kolmogorov complexity, see [LV97].

Now consider a permutation of the integers $0, \dots, n-1$ that is random in this Kolmogorov sense, i.e., the Kolmogorov complexity of the permutation is $\lg n! + O(1)$.¹ From such a permutation we construct the text T by writing the numbers in binary, using exactly $\lg n$ bits each, in the order given by the permutation, and terminating each binary representation with a special character $\#$. For example, for $n = 16$, the permutation $5, 12, 1, \dots, 15, 10$ would generate the string

$$T = 0101\#1100\#0001\#\dots\#1111\#1010\#.$$

Define $N = |T| = n(\lg n + 1)$.

¹We use Knuth's notation for O which bounds only the absolute value of f in $f = O(g)$.

We will be searching for the *patterns* consisting of the binary representation of a number between 0 and $n - 1$ followed by a #, such as 0101#. Thus, the size $|P|$ of such a pattern P is $\lg n + 1$.

2.2 Reducing Ternary Alphabet to Binary

While we use a ternary alphabet $\{0, 1, \#\}$ for convenience, this can easily be avoided by a standard trick. Our only requirement is that each pattern appears exactly once in the text. Thus, we can replace the # symbol with the special string 01^k0 , and prevent the special string from occurring in the rest of the string by introducing an additional 0 every $k - 1$ symbols in each maximal binary substring (e.g., 0101) of the original string. The original string T is easily recovered from this new binary string. In addition, this encoding has the advantage that each search string is expanded by the same factor.

By choosing k appropriately, the size of this binary string can be made $N + o(N)$, so the expansion factor is $1 + o(1)$. The binary portion of the string expands by a $1 + \frac{1}{k-1}$ factor, and each # symbol expands to k bits (adding $k - 1$ symbols). Thus the size of the new binary string is $(k - 1)n + \left(1 + \frac{1}{k-1}\right)N$. This size is minimized by choosing $k = 1 + \sqrt{1 + \lg n}$, for which the size is $N + \Theta(N/\sqrt{\lg n})$.

In the remainder of the proof, we assume that the alphabet is ternary for simplicity of presentation. We can use the preceding encoding to extend our bounds to a binary alphabet, because each search string remains uniquely occurring, and all search strings expand by the same easy-to-compute subconstant factor. (The latter property is necessary e.g. in Algorithm Encoding in Section 2.5.)

2.3 Proof Outline

Because the permutation and hence the text T has Kolmogorov complexity $\lg n! + O(1) = n \lg n + O(n)$ by Stirling's approximation, any encoding of T requires at least $N + O(N/\lg N)$ bits.

The heart of the proof is to show that the text T can be reconstructed from the index I plus two small auxiliary structures A and B . More precisely, based on the hypothetical query algorithm and the bit probes it makes, we will construct structures A and B so that, together, I , A , and B form an encoding of T . We can combine I , A , and B into one string by first encoding the length of I , using $\Theta(\lg |I|)$ bits, so that I and AB can later be partitioned. (Separating A and B will be easy.) By the Kolmogorov argument above, this proves that

$$\Theta(\lg |I|) + |I| + |A| + |B| \geq N + O(N/\lg N). \quad (1)$$

Thus, if we can build A and B sufficiently small, say $o(N)$, we obtain a lower bound of $N + o(N)$ on the size of I . But note that A and B do not even have to be sublinear in the size N of the text. Kolmogorov complexity gives us a tight bound on the encoding size, specifying that the leading constant on N is 1. Thus, provided $|A| + |B| \leq cN$ for some constant $c < 1$, we obtain a linear lower bound on $|I|$, namely $|I| \geq (1 - c)N + O(N/\lg N)$.

2.4 Ideas for Encoding

If the search algorithm made no probes to the text T , it follows trivially that the index I encodes all positional information for each pattern. Specifically, suppose there were a search algorithm **Search** that computes the location of a pattern P using only the index I . Then the original text T can be reconstructed by sequentially querying all patterns, as follows:

Algorithm Decoding

1. For each pattern P from $00 \cdots 0\#$ to $11 \cdots 1\#$:
 - (a) Let b denote the beginning of the occurrence of the pattern P found by **Search** (P).
 - (b) Write the pattern P into positions $b, b + 1, \dots, b + |P| - 1$ of the reconstructed text T .

Of course, our task is not this easy: the **Search** algorithm can perform as many as \mathcal{S} probes to the text T . To understand the structure of these probes, we can think of the **Search** algorithm as having the following outline:

Generic Algorithm Search (P, I, T)

1. Initialize our “knowledge” to the pattern P and index I .
2. For i from 1 up to at most \mathcal{S} ,
 - (a) Using current knowledge, compute which bit to next probe the text T ; call the location $\ell_{P,i}$.
 - (b) Probe $T[\ell_{P,i}]$, and add it to our knowledge.

Here we have isolated the probes to the text. Our goal is to encode in auxiliary structures the probed bits $T[\ell_{P,1}], \dots, T[\ell_{P,\mathcal{S}}]$: if these were somehow known for every pattern P , the **Search** algorithm could run without access to T , so we could recover the text by applying **Algorithm Decoding** above. We do not have to encode where the **Search** algorithm probes (the $\ell_{P,i}$'s); these positions are given by the **Search** algorithm itself.

One trivial option would be to simply write down the bits in the order they would be probed when executing **Decoding**, i.e., store the string

$$T[\ell_{00\dots 0,1}]T[\ell_{00\dots 0,2}] \cdots T[\ell_{00\dots 0,\mathcal{S}}] \cdots \cdots T[\ell_{11\dots 1,1}]T[\ell_{11\dots 1,2}] \cdots T[\ell_{11\dots 1,\mathcal{S}}]$$

as an auxiliary structure. (Define $\ell_{P,i}$ arbitrarily if fewer than i probes were made during the call to **Search** (P, I, T).) This auxiliary structure uses $\mathcal{S}n = \frac{\mathcal{S}}{\lg n + 1} N$ bits. The hard case is thus when $\mathcal{S} \geq \lg n + o(\lg n)$: if \mathcal{S} were $\leq c \lg n + o(\lg n)$ for $c < 1$, this encoding would be sufficiently small.

Slightly more sophisticated is to encode just the *new* probed bits. In other words, the encoding algorithm for T simulates the execution of **Decoding**, and whenever its subroutine **Search** probes location ℓ , the encoder first checks whether it has already encoded the bit $T[\ell]$. If so, the encoder simply omits the bit, with no explicit code; otherwise, the bit is simply

written to the encoded string. The decoder, on the other hand, has perfect knowledge of what bits it has probed so far, so whenever it comes across a new bit to probe, it plays the same game: if it has never probed that bit before, it reads the next unread bit in the encoded string, taking that as the probe result; and otherwise, it uses its knowledge of the past to determine the result of the probe. This scheme is somewhat more efficient than the trivial encoding above, but is still insufficient; however, this idea will be used heavily in our encoding.

2.5 Our Encoding Scheme

We use two auxiliary structures, A and B , to efficiently encode the probed bits during the simulated execution of **Decoding**. The basic strategy is that when searching for a pattern P , we decide whether to record in A the new bits probed in T . After all queries have been made and A has been filled, any remaining probed bits in T that remain unrecorded will be recorded in B .

More precisely, the auxiliary structure A is a table with n rows, one per pattern from $00 \cdots 0$ to $11 \cdots 1$. Algorithm **Encoding** (detailed below) loops over the patterns, following a simulated execution of **Decoding**. Each row of table A contains a bit specifying whether the row is “present.” If it is present, the row effectively records the bits probed by **Search** (P, I, T) that were not already probed (as in the second method above). The idea is that we will record the probed bits *if* they will give us a significant bonus in implicitly revealing other bits. The advantage of recording the probe values for a **Search** is that the decoder can later simulate the execution of **Search** and actually determine the location of the pattern P : this reveals $|P| = \lg n + 1$ bits of the text to be the pattern P (including the $\#$ symbol).² Some of these bits may already be known to the decoder (such as the $\#$ symbol), but some may not. If a constant fraction is not known, say at least $r \lg n$ out of the $\lg n + 1$ bits are newly revealed, we make the row of A present. We will optimize the constant r , $0 < r < 1$, later.

If a row is absent, it consists only of the presence bit, 0. If the row is present, it continues with up to four additional fields. Again, the goal of these fields is just to encode the bits probed by **Search** (P, I, T) in order to reveal the bits where the pattern P occurs. However, **Search** might probe a bit (or several) that is discovered, at the end of the search, to occur within the found instance of P . We need to avoid encoding these bits, for otherwise our “bonus” will be lost.

Thus, if **Search** probes a bit that occurs within the to-be-found instance of P , the first field of the row encodes a 1 bit to state this fact, the second field encodes the first probe with this property (in binary using exactly $\lceil \lg \mathcal{S} \rceil$ bits), and the third field encodes the distance from the beginning of the instance of P to the probed bit (in binary using exactly $\lceil \lg(\lg n + 1) \rceil$ bits).³ Then the fourth field encodes any bits probed by **Search** that were not probed or revealed before and are not within the occurrence of the pattern P . On the other

²As described in Section 2.2, the binary encoding of P in fact has $\lg n + O(\sqrt{\lg n})$ bits, not $\lg n + 1$, but for simplicity, we treat the special symbol $\#$ as just another bit. This simplification has no asymptotic effect, changing only lower-order terms.

³We do not actually have to store this distance, because it can be derived by dividing the location of the probed bit by the pattern length, $|P| = \lg n + 1$. We will encode the distance for the sake of having an encoding as explicit as possible; it will only affect the encoding length by a lower-order term.

hand, the decoder reads the answers from probes made by **Search** using the fourth field, until it reaches the probed bit specified in the second field. Then it determines the location of the pattern P using the third field, and adds those bits to its knowledge. Finally, the decoder continues reading the fourth field as usual until reaching the next pattern P and row of A .

The simpler case is when **Search** does not probe any bit in the to-be-found instance of P . Then the first field just encodes a 0 bit to state this fact, the second and third fields are empty, and the fourth field encodes all probed bits that were not already encoded in A .

After the simulated execution of **Algorithm Decoding**, some of the probed bits have been encoded in A , and other bits have been revealed by knowing the locations of patterns. Any remaining unknown bits of T are now written into B , ordered left-to-right as they occur in T . Thus, by applying essentially the same process as **Algorithm Encoding**, we can decode A and B to reconstruct the original text T . The remaining question is how efficient our encoding is.

Algorithm Encoding (P)

- Mark the predetermined locations of $\#$ symbols as known.⁴
- For P from $00 \cdots 0\#$ to $11 \cdots 1\#$:
 1. Apply **Search** (P, I, T), let b denote the beginning of the occurrence of the pattern P , and let $\ell_{P,1}, \dots, \ell_{P,k}$ denote the probed bits in order.
 2. If more than $r \lg n$ of the bits $T[b], \dots, T[b + \lg n]$ are marked known,
 - Write “0” (row of A is absent).
 3. Otherwise:
 - (a) Write “1” (row of A is present).
 - (b) If one of the probed bits is in between $T[b]$ and $T[b + \lg n]$ inclusive:
 - i. Write “1” (bit of P was probed).
 - ii. Let $\ell_{P,i}$ denote the first probed bit that is in between $T[b]$ and $T[b + \lg n]$.
 - iii. Write $i - 1$ in binary using exactly $\lceil \lg \mathcal{S} \rceil$ bits.
 - iv. Write $\ell_{P,i} - b$ in binary using exactly $\lceil \lg(\lg n + 1) \rceil$ bits.
 - (c) Otherwise, write “0” (no bit of P was probed).
 - (d) For i from 1 to k :
 - If $\ell_{P,i}$ is not in between b and $b + \lg n$ inclusive, and $T[\ell_{P,i}]$ is not marked as known:
 - i. Write $T[\ell_{P,i}]$ (as a new probed bit).
 - ii. Mark $T[\ell_{P,i}]$ as known.
 - (e) Mark the bits $T[b], \dots, T[b + \lg n]$ as known.
- For i from 1 to N ,
 4. If $T[i]$ is not marked as known, write $T[i]$ (as a bit of B).

⁴Note that these positions can be computed even with the binary encoding from Section 2.2, because all search strings expand by the same predetermined subconstant factor.

2.6 Size of Encoding

It may be surprising that our carefully worded encoding method actually saves space in the worst case over the simpler encoding algorithms, but we will show that the savings are significant. The basic idea is two-fold. First, as mentioned, if most of the bits in the instance of P are unknown, then it pays to encode the probe bits, because then we learn where P is and hence learn most of those $\lg n$ bits for free. Second, we will show that a significant portion of pattern instances will have this property, i.e. A will have many rows, inducing much savings.

Equation (1) tells us that the encoding of the triple (I, A, B) has size at least $n \lg n + O(n)$. Equivalently, we have

$$|I| + \Theta(\lg |I|) \geq n \lg n - |A| - |B| + O(n). \quad (2)$$

Now we need to estimate the sizes of A and B . Let a be the number of present rows in A , and let b be the number of absent rows in A . Thus we have

$$a + b = n. \quad (3)$$

Each row in table A consists of the present bit, sometimes (a out of n times) followed by the values of the newly probed bits optionally prefixed by some information about the location of P . Thus,

$$|A| \leq n[2 \text{ bits}] + a[\lceil \lg \mathcal{S} \rceil + \lceil \lg(\lg n + 1) \rceil \text{ bits}] + [\text{number of probed bits in } A].$$

Applying this to Equation (2) gives us

$$\begin{aligned} |I| + \Theta(\lg |I|) &\geq n \lg n - a(\lceil \lg \mathcal{S} \rceil + \lceil \lg(\lg n + 1) \rceil) \\ &\quad - \underbrace{([\text{number of probed bits in } A] + |B|)} + O(n) \end{aligned} \quad (4)$$

The work so far has simply converted the encoding method into algebra. The important part that remains is marked with a brace in Equation (4): bound the total number of bits encoded explicitly in A and B , i.e., the number of probe bits encoded in A plus the entire size of B .

Consider any bit in the text T . The bit is within exactly one of the query patterns, say P . We distinguish two cases:

Case 1: The row of A corresponding to P is absent. There are $b(\lg n + 1)$ such bits.

In this case, the bit is either encoded as a probe bit in another row of A , or it will be encoded in B . Indeed, this accounts for all bits encoded in B , and some of the probe bits encoded in A .

Case 2: The row of A corresponding to P is present.

In this case, the bit is implicitly encoded by the discovered location of P . This means that the bit will not be explicitly encoded in any future rows of A or in B . However, it may have already been encoded in an earlier row of A .

Fortunately, by the definition of A , we have a bound on the number of such bits. Namely, this row of A was chosen to be present because the number of already known bits of P was at most $r \lg n$. Thus, the number of bits of this type that are encoded as probe bits in A (over all present rows of A) is at most $ar \lg n$.

In total, we have the estimate

$$[\text{number of probed bits in } A] + |B| \leq b \lg n + b + ar \lg n.$$

Combining this equation with (3) and (4), we obtain

$$\begin{aligned} |I| + \Theta(\lg |I|) &\geq n \lg n - a(\lceil \lg \mathcal{S} \rceil + \lceil \lg(\lg n + 1) \rceil) - b \lg n - b - ar \lg n + O(n) \\ &= (n - b - ar) \lg n - a(\lg \mathcal{S} + \lg \lg n) + O(n) \\ &= (1 - r)a \lg n - a(\lg \mathcal{S} + \lg \lg n) + O(n). \end{aligned} \quad (5)$$

Lastly, for this bound to be useful, we must bound a from below. To do this, we recount which of the $n \lg n$ bits of T are encoded where (excluding $\#$ symbols, which are known a priori). On the one hand, each row of A encodes at most \mathcal{S} probe bits, plus implicitly encoding at most $\lg n$ pattern bits. Thus,

$$\text{number of bits encoded by } A \leq a(\lg n + \mathcal{S}). \quad (6)$$

On the other hand, whenever a portion of a pattern P is encoded in B , at least $r \lg n$ of the bits were already explicitly or implicitly encoded in A , leaving at most $(1 - r) \lg n$ bits of the pattern to be encoded in B . Thus,

$$\text{number of bits encoded by } B = |B| \leq b(1 - r) \lg n. \quad (7)$$

Combining these two equations,

$$n \lg n = \text{total number of bits encoded} \leq a(\lg n + \mathcal{S}) + b(1 - r) \lg n.$$

Dividing both sides by $\lg n$ and simplifying, we obtain

$$\begin{aligned} n &\leq a(1 + \mathcal{S}/\lg n) + b(1 - r) \\ &= a(1 + \mathcal{S}/\lg n) + (n - a)(1 - r) \\ &= a(r + \mathcal{S}/\lg n) + n(1 - r) \end{aligned}$$

Solving for a , we obtain the desired bound:

$$a \geq \frac{r}{r + \mathcal{S}/\lg n} n.$$

Substituting this into Equation (5), we obtain

$$\begin{aligned} |I| + \Theta(\lg |I|) &\geq (1 - r) \left(\frac{r}{r + \mathcal{S}/\lg n} n \right) \lg n - \left(\frac{r}{r + \mathcal{S}/\lg n} n \right) (\lg \mathcal{S} + \lg \lg n) + O(n) \\ &= \frac{r(1 - r)}{r + \mathcal{S}/\lg n} n \lg n - \frac{r}{r + \mathcal{S}/\lg n} n (\lg \mathcal{S} + \lg \lg n) + O(n) \end{aligned} \quad (8)$$

Focusing on just the lead $n \lg n$ term, this expression is maximized when

$$r = \sqrt{\frac{\mathcal{S}}{\lg n} \left(1 + \frac{\mathcal{S}}{\lg n}\right)} - \frac{\mathcal{S}}{\lg n},$$

giving us the bound

$$\begin{aligned} |I| + \Theta(\lg |I|) &\geq \left(1 + 2 \frac{\mathcal{S}}{\lg n} - 2 \sqrt{\frac{\mathcal{S}}{\lg n} \left(1 + \frac{\mathcal{S}}{\lg n}\right)}\right) n \lg n \\ &\quad - \left(1 - \frac{\sqrt{\mathcal{S}/\lg n}}{\sqrt{1 + \mathcal{S}/\lg n}}\right) n(\lg \mathcal{S} + \lg \lg n) + O(n) \\ &= \left(1 + 2 \frac{\mathcal{S}}{\lg n} - 2 \sqrt{\frac{\mathcal{S}}{\lg n} \left(1 + \frac{\mathcal{S}}{\lg n}\right)}\right) n \lg n + O(n \lg \mathcal{S} + n \lg \lg n) \quad (9) \end{aligned}$$

Theorem 2.2 *If there is an algorithm supporting substring searches of length $|P| = \lg n + o(\lg n)$ using at most $\mathcal{S} = \mathcal{S}(|P|, |T|)$ bit probes to a text of size $|T| = n \lg n + o(n \lg n)$, then the following lower bound on the index size $|I|$ must hold:*

$$|I| \geq \left(1 + 2 \frac{\mathcal{S}}{\lg n} - 2 \sqrt{\frac{\mathcal{S}}{\lg n} \left(1 + \frac{\mathcal{S}}{\lg n}\right)}\right) n \lg n + O(n \lg \lg n + n \lg \mathcal{S})$$

Corollary 2.1 *If $\mathcal{S} = O(\lg n)$, then $|I| = \Omega(|T|)$. For example, if $\mathcal{S} = \lg n$, then*

$$\begin{aligned} |I| &\geq (3 - 2\sqrt{2}) |T| + o(|T|) \\ &\approx 0.17157 |T| + o(|T|). \end{aligned}$$

Proof: Suppose $\mathcal{S} \leq c \lg n$. Then the coefficient of $n \lg n$ in Theorem 2.2 is $1 + 2c - 2\sqrt{c(1+c)}$. We claim that $1 + 2c - 2\sqrt{c(1+c)} > 0$ for all $c > 0$, or equivalently that $\sqrt{c(1+c)} < c + 1/2$, i.e., $c(1+c) < (c + 1/2)^2$, i.e., $c^2 + c < c^2 + c + 1/4$, which is true for all c . \square

3 Extensions

This section describes three interesting straightforward generalizations of our main result described in the previous section.

3.1 Space-Time Trade-off

Our bound extends to query time bounds \mathcal{S} that are slightly beyond $O(\lg n)$ at the cost of slightly weakening the lower bound on the index size $|I|$:

Corollary 3.1 *If there is an algorithm supporting substring searches of length $|P| = \lg n + o(\lg n)$ using at most $\mathcal{S} = \mathcal{S}(|P|, |T|) = o(\lg^2 n / \lg \lg n)$ bit probes to a text of size $|T| = n \lg n + o(n \lg n)$, then the following lower bound on the index size $|I|$ must hold:*

$$|I| = \Omega\left(\frac{\lg n}{\mathcal{S}} |T|\right)$$

Proof: Apply Equation (9). The range $\mathcal{S} = o(\lg^2 n / \lg \lg n)$ is precisely the range for which the bound

$$\Omega\left(\frac{\lg n}{\mathcal{S}} |T|\right) = \omega\left(\frac{\lg \lg n}{\lg n} |T|\right) = \omega(n \lg \lg n)$$

dominates the error factor $O(n \lg \lg n + n \lg \mathcal{S})$. □

3.2 Average-Case Query Bounds

Our results generalize to an index supporting substring searches that take linear time on average over all word queries:

Theorem 3.1 *If there is an algorithm supporting many substring searches of length $|P| = \lg n + o(\lg n)$ using at most $\mathcal{S} \leq c \lg n$ bit probes on the average to a text of size $|T| = n \lg n + o(n \lg n)$, then the following lower bound on the index size $|I|$ must hold:*

$$|I| \geq \Omega(n \lg n).$$

Proof: We modify the encoding strategy as follows. For a constant $t > 1$, patterns causing more than $t\mathcal{S}$ bit probes to the text are left absent from A , to later be encoded in table B if not already known. All other patterns are treated the same as Section 2, except that they can make as many as $t\mathcal{S}$ probes instead of just \mathcal{S} probes.

Consequently, the number of bits encoded by B may grow to incorporate these new entries left absent from B . Because the average number of probes is \mathcal{S} , there are at most n/t patterns P causing $t\mathcal{S}$ probes. Hence, Equation (7) generalizes to

$$|B| \leq (b(1-r) + n/t) \lg n.$$

The number of bits encoded by A remains bounded by Equation (6). Equation (6) remains nearly the same:

$$\text{number of bits encoded by } A \leq a(\lg n + t\mathcal{S}).$$

Combining these two equations, we obtain

$$n \lg n = \text{total number of bits encoded} \leq a(\lg n + t\mathcal{S}) + (b(1-r) + n/t) \lg n.$$

Substituting $b = n - a$ and solving for a , we obtain

$$a \geq \frac{r - 1/t}{r + t\mathcal{S}/\lg n} n.$$

Applying this bound on a to Equation (5) as before, we have

$$|I| + \Theta(\lg |I|) \geq (1 - r) \frac{r - 1/t}{r + t\mathcal{S}/\lg n} n \lg n - \frac{r - 1/t}{r + t\mathcal{S}/\lg n} (\lg t\mathcal{S} + \lg \lg n) + O(n).$$

The lead $n \lg n$ term can again be maximized with respect to the free parameters t and r , but the result is rather technical and unenlightening. To establish the theorem, it suffices to consider $r = 2/3$ and $t = 2$ which yield the bound

$$|I| \geq \frac{1}{12(1 + 3\mathcal{S}/\lg n)} n \lg n + O(n \lg \lg n + n \lg \mathcal{S}).$$

which is $\Omega(n \lg n)$ as desired. □

3.3 Nonbinary Alphabet

Our results also generalize to texts T with a nonbinary alphabet Σ , $|\Sigma| \geq 2$ and $|\Sigma|$ not necessarily a constant. In this case, the text T is constructed by writing the permutation in base $|\Sigma|$, and all bounds simply change to use $\log_{|\Sigma|}$ instead of $\lg = \log_2$. Thus we have

Corollary 3.2 *If there is an algorithm supporting substring searches of length $|P| = \log_{|\Sigma|} n + o(\log_{|\Sigma|} n)$ using at most $\mathcal{S} = \mathcal{S}(|P|, |T|)$ bit probes to a text of $|T| = n \log_{|\Sigma|} n + o(n \log_{|\Sigma|} n)$ letters from the alphabet Σ , then the following lower bound on the index size $|I|$ must hold:*

$$|I| \geq \left(1 + 2 \frac{\mathcal{S}}{\log_{|\Sigma|} n} - 2 \sqrt{\frac{\mathcal{S}}{\log_{|\Sigma|} n} \left(1 + \frac{\mathcal{S}}{\log_{|\Sigma|} n} \right)} \right) n \log_{|\Sigma|} n + O(n \log_{|\Sigma|} \log_{|\Sigma|} n + n \log_{|\Sigma|} \mathcal{S})$$

In particular, if $\mathcal{S} = O(\log_{|\Sigma|} n)$, then $|I| = \Omega(|T|)$.

4 Conclusion

This work initiates an important area in text retrieval of determining the size of the smallest index to support fast substring searches. We have proved the first nontrivial lower bound on this string-matching problem. Namely, we have shown that if the search algorithm makes only $O(|P|)$ bit probes into the text for a query substring P , then the size of the index must be proportional to the text.

Our bound of $\Omega(|T|)$ bits on the size of the index is tight up to a constant factor in our model: a trivial strategy is to store an exact copy of the text in the index (using exactly $|T|$ bits), and answer queries by scanning the index (which has zero cost in our model). However, it is natural to ask for the exact factor c for which the smallest index supporting fast queries uses $c|T| + o(|T|)$ bits. It seems difficult to improve past the trivial strategy described above, so we conjecture that the factor c is in fact 1:

Conjecture 4.1 *For any search algorithm that makes $O(|P|)$ bit probes to the text T , the index must have size at least $|T| + o(|T|)$.*

Another open problem is how small the index can be if the algorithm is allowed $O(|P|)$ probes into the text and each probe retrieves a *word* instead of a bit. Here we follow the standard model of a machine word having around $\lg |T|$ bits. It seems that with $\lg |T|$ clusters of $\lg |T|$ consecutive bit probes, no improvement can be made, and again the size of the index must be linear in the size of the text.

Finally, there is work to be done on the upper-bound side. For a more reasonable model in which probes to the index are not free, is there a strategy that matches the trivial space bound of $|T|$ above, or even achieving $O(|T|)$ bits? For the text and queries used in the lower-bound proof in this paper, the optimal bounds can be achieved: an inverted word list (for words such as 0101#) can be stored in $|T|$ bits, and suffices to answer queries in $O(|P|)$ time, indeed, $O(|P|)$ bit probes to the index (with no probes to the text). Can an $O(|P|)$ time bound be obtained in the more general setting of searching for arbitrary substrings instead of just words? Currently the best-known structures use either $\Omega(|T|)$ words of space [Wei73, McC76, Ukk95, GK97], or $\Omega(|T| \lg^\varepsilon |T|)$ bits of space [FM02], or $O(|T|)$ bits of space but require an additional $\Omega(\lg^\varepsilon N)$ time for queries [GV00, FM00, Sad00]. One of these structures [GV00] would even obtain the desired $O(|P|)$ bit-probe bound if it did not have the extra $\Omega(\lg^\varepsilon N)$ factor.

Acknowledgments

This work was initiated at a Schloss Dagstuhl seminar on Data Structures, organized by Susanne Albers, Ian Munro, and Peter Widmeyer. During the open problem session, Roberto Grossi posed the problem at hand. We thank Prosenjit Bose, Andrej Brodnik, Roberto Grossi, and Ian Munro for helpful discussions during the meeting. We also thank the anonymous referees for their helpful comments.

References

- [BMRS00] Harry Buhrman, Peter Bro Miltersen, Jaikumar Radhakrishnan, and Venkatesh Srinivasan. Are bitvectors optimal? In *Proceedings of the 32nd Annual ACM Symposium on Theory of Computing*, pages 449–458, Portland, Oregon, May 2000.
- [FKS84] Michael L. Fredman, János Komlós, and Endre Szemerédi. Storing a sparse table with $O(1)$ worst case access time. *Journal of the ACM*, 31(3):538–544, July 1984.
- [FM00] Paolo Ferragina and Giovanni Manzini. Opportunistic data structures with applications. In *Proceedings of the 41st IEEE Symposium on Foundations of Computer Science*, pages 390–398, Redondo Beach, California, 2000.
- [FM02] Paolo Ferragina and Giovanni Manzini. On compressing and indexing data. Technical Report TR-02-01, Dipartimento di Informatica, Università di Pisa, 2002.
- [GK97] R. Giegerich and S. Kurtz. From Ukkonen to McCreight and Weiner: A unifying view of linear-time suffix tree construction. *Algorithmica*, 19(3):331–353, 1997.

- [GV00] Roberto Grossi and Jeffrey Scott Vitter. Compressed suffix arrays and suffix trees with applications to text indexing and string matching. In *Proceedings of the 32nd Annual ACM Symposium on Theory of Computing*, pages 397–406, Portland, Oregon, May 2000.
- [KMP77] Donald E. Knuth, James H. Morris, Jr., and Vaughan R. Pratt. Fast pattern matching in strings. *SIAM Journal on Computing*, 6(2):323–350, 1977.
- [LV97] Ming Li and Paul Vitányi. *An Introduction to Kolmogorov Complexity and its Applications*. Springer-Verlag, New York, second edition, 1997.
- [McC76] Edward M. McCreight. A space-economical suffix tree construction algorithm. *Journal of the Association for Computing Machinery*, 23(2):262–272, 1976.
- [Mil93] Peter Bro Miltersen. The bit probe complexity measure revisited. In P. Enjalbert, A. Finkel, and K. W. Wagner, editors, *Proceedings of the 10th Annual Symposium on Theoretical Aspects of Computer Science*, pages 662–671, Würzburg, February 1993.
- [Pag01] Rasmus Pagh. Low redundancy in static dictionaries with constant query time. *SIAM Journal on Computing*, 31(2):353–363, 2001.
- [Sad00] Kunihiro Sadakane. Compressed text databases with efficient query algorithms based on the compressed suffix array. In *Proceedings of the 11th International Symposium on Algorithms and Computation*, volume 1969 of *Lecture Notes in Computer Science*, pages 410–421, 2000.
- [Sad02] Kunihiro Sadakane. Succinct representations of *lcp* information and improvements in the compressed suffix array. In *Proceedings of the 13th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 225–232, San Francisco, California, January 2002.
- [Ukk95] E. Ukkonen. On-line construction of suffix trees. *Algorithmica*, 14(3):249–260, 1995.
- [Wei73] Peter Weiner. Linear pattern matching algorithm. In *Proceedings of the 14th Annual IEEE Symposium on Switching and Automata Theory*, pages 1–11, Iowa City, Iowa, 1973.
- [Yao81] Andrew Chi Chih Yao. Should tables be sorted? *Journal of the Association for Computing Machinery*, 28(3):615–628, 1981.