

String Matching in Lempel-Ziv Compressed Strings*

Martin Farach[†]
Rutgers University

Mikkel Thorup[‡]
University of Copenhagen

Abstract

String matching and Compression are two widely studied areas of computer science. The theory of string matching has a long association with compression algorithms. Data structures from string matching can be used to derive fast implementations of many important compression schemes, most notably the Lempel-Ziv (LZ77) algorithm. Intuitively, once a string has been compressed – and therefore its repetitive nature has been elucidated – one might be tempted to exploit this knowledge to speed up string matching. The *Compressed Matching Problem* is that of performing string matching in a compressed text, without uncompressing it. More formally, let \mathcal{T} be a text, let \mathcal{Z} be the compressed string representing \mathcal{T} , and let \mathcal{P} be a pattern. The Compressed Matching Problem is that of deciding if \mathcal{P} occurs in \mathcal{T} , given only \mathcal{P} and \mathcal{Z} . Compressed matching algorithms have been given for several compression schemes such as LZW.

In this paper, we give the first non-trivial compressed matching algorithm for the classic adaptive compression scheme, the LZ77 algorithm. In practice, the LZ77 algorithm is known to compress more than other dictionary compression schemes, such as LZ78 and LZW, though for strings with constant per bit entropy, all these schemes compress optimally in the limit. However, for strings with $o(1)$ per bit entropy, while it was recently shown that the LZ77 gives compression to within a constant factor of optimal, schemes such as LZ78 and LZW may deviate from optimality by an exponential factor. Asymptotically, compressed matching is only relevant if $|\mathcal{Z}| = o(|\mathcal{T}|)$, i.e. if the compression ratio $|\mathcal{T}|/|\mathcal{Z}|$ is more than a constant. These results show that LZ77 is the appropriate compression method in such settings.

We present an LZ77 compressed matching algorithm which runs in time $O(N \log^2 U/N + P)$ where $N = |\mathcal{Z}|$, $U = |\mathcal{T}|$, and $P = |\mathcal{P}|$. Compare with the naïve “decompression” algorithm, which takes time $\Theta(U + P)$ to decide if \mathcal{P} occurs in \mathcal{T} . Writing $U + P$ as $N \cdot U/N + P$, we see that we have improved the complexity, replacing the compression factor U/N by a factor $\log^2 U/N$. Our algorithm is *competitive* in the sense that $O(N \log^2 U/N + P) = O(U + P)$, and *opportunistic* in the sense that $O(N \log^2 U/N + P) = o(U + P)$ if $N = o(U)$ and $P = o(U)$.

*A preliminary version of the paper appeared in STOC '95.

[†]Department of Computer Science, Rutgers University, Piscataway, NJ 08855, USA. (farach@cs.rutgers.edu; <http://www.cs.rutgers.edu/~farach>). Part of this work was done while the author was visiting the University of Copenhagen; Supported in part by DIMACS (Center for Discrete Mathematics and Theoretical Computer Science), a National Science Foundation Science and Technology Center under NSF contract STC-8809648.

[‡]Department of Computer Science, University of Copenhagen, Universitetsparken 1, 2100 Kbh. Ø, Denmark. (mthorup@diku.dk, <http://www.diku.dk/~mthorup>). This work was done while visiting DIMACS.

1 Introduction

String matching and Compression are two widely studied areas of computer science. The String Matching Problem is that of detecting if pattern \mathcal{P} occurs in text \mathcal{T} . The Compression Problem is that of developing schemes which can be used to give succinct representations for large files. Each of these fields has its own very rich theory. The theory of string matching has a long association with compression algorithms. Data structures from string matching can be used to derive fast implementations of many important compression schemes [6, 13, 14], most notably the Lempel-Ziv algorithm [11]. The reason for this connection is clear; string algorithms and data structures most commonly detect repetitions of various sorts within strings. The elimination of repetition is the key to compression.

Intuitively, once a string has been compressed – and therefore its repetitive nature has been elucidated – one might be tempted to exploit this knowledge to speed up string matching. This idea is inherent in the work of Amir and Benson [2], where they introduced the *Compressed Matching Problem*. They were motivated by the practical consideration that, increasingly, files are stored in a compressed format, and that current string matching technology requires the files be uncompressed before string matching takes place.

The Compressed Matching Problem

Instance: A compressed string \mathcal{Z} and a pattern \mathcal{P} .

Question: Does \mathcal{P} occur in the text \mathcal{T} represented by \mathcal{Z} ?

Assuming $|\mathcal{Z}| \ll |\mathcal{T}|$, we hope for substantial savings over the naïve $\Omega(|\mathcal{T}| + |\mathcal{P}|)$ algorithm which first decompresses, and then performs standard string matching. An *optimal* algorithm takes time $\Theta(|\mathcal{Z}| + |\mathcal{P}|)$. Of course, the nature of the compression algorithm determines the difficulty of performing efficient compressed matching. For the remainder of the paper, let U be the size of the uncompressed text, N the size of its compressed representation, and P the size of the pattern. Note that compressed matching is only relevant if the *compression ratio* is more than a constant, i.e. if $U/N = \omega(1)$; otherwise the naïve decompressing algorithm is itself optimal for compression schemes which are decodable in linear time. We will say that a compressed matching algorithm is *competitive* if its running time is $O(U + P)$, that is, if it is never worse than the naïve “decompression” algorithm. We will say that an algorithm is *opportunistic* if its running time is $o(U + P)$ whenever $N = o(U)$ and $P = o(U)$. Note that any optimal algorithm is both competitive and opportunistic. Generally a competitive and opportunistic algorithm has a time complexity of the form $O(N\alpha(U/N) + P)$. Thus, in the absence of optimal algorithms, it is natural to look for a compressed matching algorithm with time complexity $O(N(\text{polylog } U/N) + P)$, and we will refer to such algorithms as *pseudo-optimal*.

In this paper, we give a pseudo-optimal compressed matching algorithm for the classic adaptive compression scheme, the so-called LZ77¹ algorithm [11]. To the best of our knowledge this is the first non-trivial LZ77 compressed matching algorithm, and the search for such an algorithm was advocated in [3], in which strings compressed by LZW are considered. In practice, the LZ77 algorithm is known to compress more than other dictionary compression schemes, such as LZ78 and LZW, though for strings with constant per bit entropy, all these schemes compress optimally

¹The LZ77 algorithm is sometimes called LZ1 and the LZ78 algorithm is sometimes called LZ2

in the limit. However, for strings with $o(1)$ per bit entropy, while Kosaraju [9] recently showed that the LZ77 gives compression to within a constant factor of optimal, schemes such as LZ78 and LZW can deviate from optimality by an exponential factor.

It should be noted that for most of today's applications of general compression, the compression ratio is not large enough to justify an asymptotic analysis. However, our study of compressed matching should be seen as an approach to developing general techniques for dealing with low-entropy strings. Consider, for example, a long text \mathcal{T} which is essentially just a repetition of a some rather long patterns, but where various error-patterns may occur sporadically. The text \mathcal{T} will then have low entropy, and hence, space-wise, we will benefit from compressing it into a compressed string \mathcal{Z} , where $|\mathcal{Z}| \ll |\mathcal{T}|$ if the repeated patterns are long. If we later want to look for some error pattern \mathcal{P} , compressed matching finds it directly from \mathcal{Z} , capitalizing on the small size of \mathcal{Z} both in time and space. Alternatively, if one knew the patterns that \mathcal{T} was made up from, one could easily make a specialized compression facilitating fast matching. Indeed, without general compressed matching routines, if we wanted to facilitate matching, we would have been forced to make a specialized compression, and this may exactly be the type of reason why general compression is not more frequently used on low-entropy texts. The point of general compressed matching is to avoid specialized programming; the more operations we can perform directly on generally compressed text, the less specialized programming is needed in dealing with low-entropy objects.

1.1 Previous Results

The theory of compressed matching was initiated with the study of two dimensions run-length compression, the compression algorithm used for fax transitions. In [2, 4], an optimal $O(N + P^2)$ algorithm was derived for finding a $P \times P$ two dimensional pattern within a file of length N .

In one dimension, the formerly most commonly used compression algorithm was the LZW algorithm [16] as implemented by the UNIX compress program [15]. It seems to have been replaced in popularity by a variant on LZ77, as implemented in GNU's gzip. In [3], it was shown that a length P pattern can be found in a length N text in time $O(N + P^2)$, or in $O(N \log P + P)$ in LZW compressed texts. In [10], Kosaraju improved this complexity to $O(N + P^{1+\epsilon})$. The techniques in [3] show that it is trivial to perform linear time compressed matching on LZW files if the pattern is of constant length. In contrast, it seems quite difficult to perform any non-trivial compressed matching on LZ77 files, *even for patterns of length 2!*

1.2 Our Results

We will present an algorithm for the LZ77 compressed matching problem showing:

Theorem 1.1 *Given an LZ77 compressed string \mathcal{Z} of size N representing a text \mathcal{T} of size U , and given a pattern \mathcal{P} of size P there is a randomized (Las Vegas) algorithm to decide if \mathcal{P} occurs in \mathcal{T} which runs in time $O(N \log^2 U/N + P)$ w.h.p.*

Technical Issues and Contribution: The centerpiece of our algorithm is an entropy based potential function argument. In high entropy strings, the "information" can be thought of as being more or less evenly distributed through the string. For lower entropy strings, the information

tends to clump in certain regions. We take advantage of this clumpiness to efficiently uncompress only the relevant parts of the underlying text string, as described in Sections 3 and 4. The ideas presented therein can be used directly to arrive at a deterministic algorithm with running time $O(N \log(U/N)(\log(U/N) + \log P))$. Alternatively, combining with ideas from [3, 10], we get an $O(N(\log^2 U/N) + P^{1+\epsilon})$ deterministic algorithm.

Neither of these deterministic bounds is fully competitive. They both perform worse than the naïve decompressing $\Theta(U + P)$ algorithm if the pattern is large and the compression ratio is low. However, we give a novel application of Karp and Rabin’s fingerprint method [8], which gives us a randomized $O(N \log^2(U/N) + P)$ time algorithm for LZ77 compressed matching. Thus we get an algorithm which is both opportunistic and competitive. In fact this last technique also applies to the best known compressed matching algorithm for LZW [3] to yield a $O(N \log U/N + P)$ time algorithm. Thus a side-effect of this paper is the first competitive and opportunistic LZW compressed matching algorithm.

1.3 The LZ77 Algorithm

The LZ77 algorithm gives a very natural way of representing a string and is defined as follows. Given some alphabet Σ , an *LZ77 compressed string* is a string of the form: $\mathcal{Z} = (P_0, L_0, C_0) \cdots (P_i, L_i, C_i) \cdots (P_{N-1}, L_{N-1}, C_{N-1}) \in (\mathbb{N}_0, \mathbb{N}_0, \Sigma)^N$. Let $U_0 = 0$ and $U_i = U_{i-1} + L_{i-1} + 1$ for $i > 0$. For \mathcal{Z} to be well-formed, it is required that, for all $L_i > 0, P_i < U_i$. Now, \mathcal{Z} represents \mathcal{T} if for $i = 0, \dots, N-1$, $\mathcal{T}[U_i, U_i + L_i - 1] = \mathcal{T}[P_i, P_i + L_i - 1]$, $\mathcal{T}[U_i + L_i] = C_i$. Here $\mathcal{T}[I]$ denotes the $(I + 1)$ th character of \mathcal{T} , and $\mathcal{T}[I, J]$ denotes the segment $\mathcal{T}[I] \cdots \mathcal{T}[J]$. Clearly, \mathcal{T} is uniquely defined in terms of \mathcal{Z} . For example, $\mathcal{Z} = (0, 0, a)(0, n-2, a)$ represents $\mathcal{T} = a^n$. For comparison, LZW can never compress a text to less than the square root of its original length.

Simple transformations to LZ77: First, note that each (uncompressed) text fragment $\mathcal{T}[U_i, U_i + L_i]$ consists of a region in which the string is unknown (without uncompressing) followed by a single known character. In order to facilitate a more uniform treatment, we simplify the scheme by the following trivial transformation. First, map Σ into $[1, |\Sigma|]$. Then take each triple (P_i, L_i, C_i) and replace it by the pairs $(P_i, L_i)(-C_i, 1)$. In the special case where $L_i = 0$, replace (P_i, L_i, C_i) with $(-C_i, 1)$ so all specified intervals have positive length. Now, given a compressed text $\mathcal{Z} = (P_0, L_0) \cdots (P_{N-1}, L_{N-1})$, the uncompressed text \mathcal{T} is defined by

$$\begin{aligned} \forall i = 0, \dots, N-1 : \\ \mathcal{T}[U_i, U_i + L_i - 1] = \mathcal{T}[P_i, P_i + L_i - 1] \end{aligned} \tag{1}$$

where $U_i = \sum_{j < i} L_j$ and where, for any integer C representing a character, $\mathcal{T}[-C] = C$. Thus, we imagine that \mathcal{T} has been prefixed by Σ in the negative positions. For example, a^n is compressed to $\text{“}a\text{”}(-1, 1)(0, n-1)$. For the moment, we will make the simplifying assumption that \mathcal{Z} is not *self-referential* meaning that we assume $L_i \leq U_i - P_i$ for all i . In practice, this means that if we decompress \mathcal{Z} using Eq. 1, then $\mathcal{T}[U_i, U_i + L_i - 1]$ can be copied directly from the previously decompressed text $\mathcal{T}[0, U_i - 1]$. The same simplification is made by `gnuzip` and most other implementations of LZ77. In the worst case, it gives a blow-up in the size of the compressed text which is logarithmic in the compression factor. For example, our representation of a^n is changed to

“a” $(-1, 1)(0, 1)(0, 2)(0, 4) \cdots (0, 2^{\lfloor \log n \rfloor - 1})(0, n - 2^{\lfloor \log n \rfloor} + 1)$. However, in Section 7, we will show how self-reference can be dealt with directly, avoiding any increase in asymptotic complexity.

In Section 2, we give an outline of the algorithm and conclude with an outline of the rest of the paper.

2 General Algorithmic Outline

We can make the following simple observation which limits where we need to look for the pattern occurrence.

Observation 2.1 ([3]) *If u is the least position such that $\mathcal{P} = \mathcal{T}[u, u + |\mathcal{P}| - 1]$, then $U_i \in [u, u + |\mathcal{P}|]$ for some $i > 0$.*

Algorithmically we can formulate Observation 2.1 as follows. Suppose for $i = 0, \dots, N$ we have:

1. The longest pattern substrings \mathcal{P}_i^+ which is a prefix of $\mathcal{T}[U_i, U]$
2. The longest pattern substring \mathcal{P}_i^- which is a suffix of $\mathcal{T}[0, U_i + L_i - 1]$.

Then the compressed pattern matching problem reduces to deciding whether \mathcal{P} is a substring of $\mathcal{P}_i^- \mathcal{P}_{i+1}^+$ for some $i > 0$.

The brute force method first decompress the whole text \mathcal{T} in time $\Theta(U)$ using Eq. 1, and then performs standard string matching in time $\Theta(P + U)$. The conclusion we may draw from Observation 2.1 is that instead of uncompressing the whole text, we only need to uncompress the $2N$ $\mathcal{P}_i^{+/-}$ s. Note that we cannot afford to identify all the $\Omega(PN)$ individual characters in the $\mathcal{P}_i^{+/-}$ s. However, each $\mathcal{P}_i^{+/-}$ can be represented by a beginning and ending pointer into the pattern, i.e. in constant space. The pattern will be preprocessed so that the $\mathcal{P}_i^{+/-}$ s can be manipulated in terms of fast concatenations and truncations of pattern substrings in their succinct representations.

In order to focus on the decompression around the $\mathcal{P}_i^{+/-}$ s, we will first use Eq. 1 in reverse finding out which positions the $\mathcal{P}_i^{+/-}$ s correspond to earlier in the text and alphabet. Afterwards we know exactly which parts of the text it is relevant to decompress. Formally, we introduce what we will call *informers*. An informer x has a variable position $u(x) \in \{0, \dots, U - 1\}$ in the uncompressed text \mathcal{T} . When the compressed pattern matching starts, for $i = 0, \dots, N - 1$, we have a *right informer* r_i placed at position U_i , and our goal is to load it with \mathcal{P}_i^+ . Also, we have a *left informer* l_i placed at position $U_i + L_i - 1$ which we want to load with \mathcal{P}_i^- . The compressed pattern matching now proceeds in two phases.

The Winding Phase For $s := N - 1$ down to 0, we take all informers between U_s and $U_s + L_s - 1$ and shift them $U_s - P_s$ to the left. Since \mathcal{Z} is not self-referential, $U_s - P_s \geq L_s$. Hence, inductively, we may conclude that step s does not leave any informers to the right of U_s . As a consequence, we may reformulate step s as, for all x such that $u(x) \geq U_s$, set $u(x) := u(x) - (U_s - P_s)$. When the winding ends, all informers will have been shifted into the alphabet at negative positions. Immediately this allows us to identify the first character of \mathcal{P}_i^+ and the last character of \mathcal{P}_i^- , so in particular, the winding solves the compressed pattern matching problem for patterns of length 2. In each of the above N iterations, we risk shifting $\Omega(N)$ informers to the left. Nevertheless, as the cornerstone of this paper, in Section 3 we will implement the winding in $O(N \log^2(U/N))$ total time.

Figure 2.1: In (a), the interval to the right of the cut has informers marked with a circle. All other informers are marked with a plus. At (b), we merge the informers. After (c), we have bottomed out the computation and have partial information loaded in the informers. We un-merge the circle informers in (d) and extend the pattern substrings of all those that overlap the cut boundary. Above, it looks as if all the informers are both left and right informers, but this is to cover the case of overlap.

The Unwinding Phase In this phase we move the informers back, reversing shifts of the winding precisely. However, with each right informer, we will associate a variable pattern substring $\mathfrak{t}(x)$ which, relative to the current position $u(x)$, satisfies the condition that it is the longest pattern substring which is a prefix of the “known” part of $\mathcal{T}[u(x), u-1]$; after unwinding step s , we consider $\mathcal{T}[0, U_s + L_s - 1]$ known. Similarly, if x is a left informer, $\mathfrak{t}(x)$ is the longest pattern substring which is a suffix of $\mathcal{T}[0, u(x)]$. Satisfying these invariants implies that when the informers come back to their original positions $\mathfrak{t}(r_i) = \mathcal{P}_i^+$ and $\mathfrak{t}(l_i) = \mathcal{P}_i^-$, for $i = 0, \dots, N - 1$, as desired. If $u(x) < 0$, the informer is still in the alphabet, and it is then just loaded with the single character it is positioned at.

More constructively, as soon as the winding has ended, for each informer x , set $\mathfrak{t}(x)$ to the string consisting of the character numbered $-u(x)$. Recall that in step s of the winding, we took all informers x at positions $u(x) > U_s$ and shifted them $U_s - P_s$ to the left. Assume, for each s , that we have recorded the set X_s of informers that were shifted by step s . Then, for $s := 0$ to $N - 1$, we unwind the informers as follows. First, for all $x \in X_s$, set $u(x) := u(x) + U_s - P_s$, returning them to their position before winding step s . Concerning the loaded pattern substrings, for all right informers $x \in X_s$, truncate $\mathfrak{t}(x)$ so that it does not extend beyond $U_s + L_s - 1$, i.e. if $u(x) + |\mathfrak{t}(x)| > U_s + L_s$, delete the last $u(x) + |\mathfrak{t}(x)| - U_s + L_s$ characters of $\mathfrak{t}(x)$. Then $\mathfrak{t}(x)$ satisfies that it is the longest pattern substring which is a prefix of $\mathcal{T}[u(x), U_s + L_s - 1]$. In particular, this holds for the right informer r_s , which has now reached its original position U_s . For all the other right informers $x \notin X_s$, if $\mathfrak{t}(x)$ extends to $U_s - 1$, i.e. if $u(x) + |\mathfrak{t}(x)| - 1 = U_s - 1$, make $\mathfrak{t}(x)$ the longest pattern substring which is a prefix of $\mathfrak{t}(x)\mathfrak{t}(r_s)$. Now all right informers x satisfy that $\mathfrak{t}(x)$ is the longest pattern substring which is a prefix of $\mathcal{T}[u(x), U_s + L_s - 1]$.

Concerning the loading of the left informers, first notice that all left informers that are not shifted satisfy the invariant that $\mathfrak{t}(x)$ is the longest pattern substring which is a suffix of $\mathcal{T}[0, u(x)]$. In particular, this holds for l_{i-1} that has received its final loading, i.e. $\mathfrak{t}(l_{i-1}) = \mathcal{P}_{i-1}^-$. For the shifted left informers $x \in X_s$, we need to do as follows. If $\mathfrak{t}(x)$ extends back to U_s , i.e. if $u(x) - |\mathfrak{t}(x)| + 1 \leq U_s$, then first we delete the first $U_s - (u(x) - |\mathfrak{t}(x)| + 1)$ characters of $\mathfrak{t}(x)$, second we make $\mathfrak{t}(x)$ the longest pattern substring which is a suffix of $\mathfrak{t}(l_{s-1})\mathfrak{t}(x)$. Now all left informers x satisfy that $\mathfrak{t}(x)$ is the longest pattern substring which is a suffix of $\mathcal{T}[0, u(x)]$. Thus, we can proceed to the unwinding of winding step $s + 1$. Figure 2.1, page 5, illustrates the work of the algorithm.

In the worst case, the above naïve unwinding could give rise to $\Theta(N^2)$ operations on the loaded pattern substrings. However, many of these loaded pattern substrings can be deduced from the

others, and in Section 4 it will be shown that given the winding from Section 3, we only need to make $O(N \log(U/N))$ operations on pattern substrings for the unwinding. Having applied an $O(P)$ pattern preprocessing technique from [7], each of these operations are done in $O(\log P)$ time, giving a total complexity for the compressed pattern matching problem of $O(P + N \log(U/N)(\log(U/N) + \log P))$. For the case where $\log P = \omega(\log U/N)$, in Section 5, we present a Monte Carlo randomized algorithm which solves the compressed pattern matching problem within the announced time bound of $O(P + N \log^2(U/N))$. In Section 6, we give a checker with which we convert our Monte Carlo algorithm into a Las Vegas one. Finally, in Section 7, we will show how we can deal with self-references.

3 Fast Winding

3.1 A fast algorithm

In this section, we will speed-up the naïve $\Omega(N^2)$ winding from the last section to run in time $O(N \log^2(U/N))$. Left and right informers will be wound separately. In particular this means that if two informers end up at the same position, then we can remove one of them as a duplicate, restoring it at the appropriate time of the unwinding. First consider the pseudo code for winding described in the last section.

Algorithm A Winds a list \mathcal{L} , $|\mathcal{L}| = O(N)$, of informers sorted according to u -values.

A.1. For $s := N - 1$ downto 0 do

A.1.1. Cut \mathcal{L} at U_s into lists \mathcal{L}_0 and \mathcal{L}_1 , i.e. $\mathcal{L} = \mathcal{L}_0 \mathcal{L}_1$ and $x \in \mathcal{L}_0$ iff $u(x) < U_s$.

A.1.2. Subtract $U_s - P_s$ from u for all informers in \mathcal{L}_1 , corresponding to setting $u(x) := u(x) - U_s + P_s$ for all $x \in \mathcal{L}_1$.

A.1.3. Merge \mathcal{L}_0 with \mathcal{L}_1 into \mathcal{L} , removing informers that are duplicate with respect to u -values.

Using any standard balanced tree data structure [1] to implement the lists, all steps but the merge step (A.1.3) can be implemented in time $O(\log N)$, hence in time $O(N \log N)$ over all iterations. However, for the merge step it is possible to construct inputs for which we have both \mathcal{L}_0 and \mathcal{L}_1 of length $\Omega(N)$ in $\Omega(N)$ iterations. Thus, it seems that the merging alone has running time $\Omega(N^2)$. It turns out that what we need is an old folklore version of merging, below referred to as *segment merge* (Algorithm B). This version is normally abandoned because it has a bad worst-case performance for the isolated task of merging two lists (if the lists have lengths n_1 and n_2 it takes time $O((n_1 + n_2) \log(n_1 + n_2))$). Here we will resurrect segment merge, showing that in our context it has a very good amortized complexity. More precisely, with a potential function argument, we will show that segment merge takes $O(N \log N \log U)$ time in total, and later we will reduce this complexity to $O(N \log^2(U/N))$. In the following algorithm, let $\text{head}(L)$ be the first item on list L , and $\text{tail}(L)$ be the the list with the first item removed.

Algorithm B *Segment-Merge*($\mathcal{L}_0, \mathcal{L}_1$) implements step A.1.3.

B.1. If one of the lists is empty, return the other.

B.2. For $i := 0, 1$, set $x_i := \text{head}(\mathcal{L}_i)$.

- B.3. If $u(x_0) = u(x_1)$, return
 $x_0 \text{Segment-Merge}(\text{tail}(\mathcal{L}_0), \text{tail}(\mathcal{L}_1))$.
- B.4. For $i := 0, 1, \bar{i} = 1 - i$, if $u(x_i) < u(x_{\bar{i}})$:
- B.4.1. Cut \mathcal{L}_i at $u(x_{\bar{i}})$ into \mathcal{L}'_i and \mathcal{L}''_i .
- B.4.2. Return $\mathcal{L}'_i \text{Segment-Merge}(\mathcal{L}''_i, \mathcal{L}_{\bar{i}})$.

In the following, by a *proper* call to `Segment-Merge` we refer to one where both lists are nonempty. Note that the number of improper calls is bounded by the number of calls to `Segment-Merge` from step A.1.3, hence by N . Hence we will only be concentrating on the proper calls to `Segment-Merge`. A single call `Segment-Merge`($\mathcal{L}_0, \mathcal{L}_1$) can give rise to $|\mathcal{L}_0| + |\mathcal{L}_1| - 1 = \Omega(N)$ recursive calls. Nevertheless we will show that from all the N direct calls to `Segment-Merge` from Algorithm A, we will only get $O(N \log U)$ calls in total.

3.2 The potential function argument

Our argument is based on a simple potential function. In the above algorithm each informer is either in a list or deleted. For a “live” informer x , i.e. x has not been removed as a duplicate, in a list, we define $d^-(x) = u(x) - u(x^-)$ if x^- is the predecessor of x , and $d^-(x) = 2U$ if x is the first element in the list. Similarly, $d^+(x) = u(x^+) - u(x)$ if x^+ is the successor of x , and $d^+(x) = 2U$ if x is the last element in the list. Note that for any informer x in a list, $1 \leq d^{-/+}(x) \leq 2U$. Now the *local potential* $\pi(x)$ is defined as $1 + \log d^-(x) + \log d^+(x)$. If x has been deleted $\pi(x) = 0$. The *total potential* is $\Pi = \sum_x \pi(x) > 0$.

The idea behind Π is that if, during a merge, we insert a segment of a list between two elements x_0 and x_1 in another list, then we halve either $d^+(x_0)$ or $d^-(x_1)$. Thereby we decrease Π by at least 1. Thus, merges will be accounted for as decreases in Π . Before going into details about the decreases, first we will discuss the possible increases to Π .

Initially, the \mathcal{L} contains all N informers, and $\Pi \leq N(1 + 2 \log 2U) = O(N \log U)$. During the algorithm, Π will increase during the cut step A.1.1. More precisely, the cut will affect the last informer x_0 in \mathcal{L}_0 , setting $d^+(x_0) = 2U$, and the first informer x_1 in \mathcal{L}_1 , setting $d^-(x_1) = 2U$. Hence Π will increase by at most $2 \log 2U$, so in total Π is increased by $O(N \log U)$ in step A.1.1. The subtraction step A.1.2 does not change Π . Hence, we prove the bound if we can show that the merge step A.1.3 decreases Π with a number proportional to the number of recursive calls to `Segment-Merge`, as defined in Algorithm B.

In order to study the total decrease in Π during a merge, we will see the process as divided into two phases. In the first phase we merge the elements, allowing multiplicity. Thus, the only effect of the first phase is to reduce d^- and d^+ for some of the informers. In the second phase we delete the duplicate elements. Notice that neither phase increases the local potential of any informer.

There are now two cases to consider. If we go into the case in step B.3 with $u(x_0) = u(x_1)$, then this corresponds to the deletion of x_1 in the second phase. Since any live informer has $\pi \geq 1$, this step corresponds to a decrease in Π of at least 1.

Now consider the case in step B.4 where $u(x_i) < u(x_{\bar{i}})$. We will relate this to a decrease in Π during the first phase. Let $\mathcal{L}_{\bar{i}}^o$ be the list of $x_{\bar{i}}$ from the original (non-recursive) call to `Segment-Merge` from step A.1.3. Suppose that $x_{\bar{i}}$ is not the first element in $\mathcal{L}_{\bar{i}}^o$ and let $x_{\bar{i}}^-$ be its predecessor. Then the first phase of the merge is to insert the segment \mathcal{L}'_i between $x_{\bar{i}}^-$ and $x_{\bar{i}}$,

thereby either halving $d^+(x_{\bar{i}}^-)$ or $d^+(x_{\bar{i}})$, in either case decreasing Π by at least 1. Note that if $x_{\bar{i}}^-$ or $x_{\bar{i}}$ are deleted it is not until the second phase. In the other case, where $x_{\bar{i}}$ is the first element of $\mathcal{L}_{\bar{i}}^o$, we are inserting $\mathcal{L}_{\bar{i}}^i$ before $x_{\bar{i}}$, thereby decreasing $d^-(x)$ from $2U$ to at most U , again decreasing Π by at least 1. This completes the basic potential function argument.

3.3 The final touch

We will now apply a bucketing argument, with which we reduce $O(N \log U \log N)$ to $O(N \log^2(U/N))$. Fix $F = \lceil U/N \rceil$. We have so far used a single list to manipulate all informers. However, we can partition these lists into N lists, each over a range of size U/N , by which we will reduce both the data structural time as well as tightening the analysis. We note for completeness that our complexity will become $O(N \log(U/N) \min\{\log N, \log(U/N)\})$. Specifically, before running Algorithm A we will preprocess the compressed string \mathcal{Z} in two steps:

1. Insert cuts at all positions $kF < U - 1$ where $k \in \mathbb{N}_0$. That is, if a pair (P, L) in \mathcal{Z} covers the interval $[U, U + L - 1]$, i.e. U is the sum of the lengths of the preceding pairs in \mathcal{Z} , and if $kF \in [U + 1, U + L - 1]$, we replace (P, L) by the two pairs $(P, L')(P + L', L - L')$ where $L' = kF - U$. Clearly, this step does not change the represented text, and the size of the compressed text is at most doubled by the at most $\lfloor U/F \rfloor \leq N$ new cuts.
2. For all pairs (P, L) where $kF \in [P + 1, P + L - 1]$ for some $k \in \mathbb{N}_0$, cut (P, L) into the pairs $(P, L')(P + L', L - L')$ where $L' = kF - P$.

After step 1 no length is greater than F , so step 2 can only cut each pair once. Thus our new compressed text is still of size $O(N)$, and clearly it represents the same text \mathcal{T} .

Now suppose that both step 1 and step 2 have been applied to \mathcal{Z} before the winding. A first advantage is that we can partition all our lists into buckets $B_k = [kF, (k + 1)F - 1]$. Step 1 above, ensures that for the cut step A.1.1, \mathcal{L}_1 is entirely contained in the bucket $B_{\lfloor U_i/F \rfloor}$, and step 2 ensures that after the subtraction step A.1.2, \mathcal{L}_1 is entirely contained in $B_{\lfloor P_i/F \rfloor}$. As there are only F possible informer positions in a bucket, all list operations are now on lists of length bounded by F . Thus the complexity of each list operation is changed from $O(\log N)$ to $O(\log F) = O(\log U/N)$.

In fact the above preprocessing also allows us to tighten the potential function argument, reducing the number of merges from $O(N \log U)$ to $O(N \log U/N)$. Our lists are now partitioned into buckets, and clearly, the maximum distance between two informers in the same bucket is bounded by F . Thus, if we define $d^-(x) = 2F$ if x is the left-most informer in a bucket and $d^+(x) = 2F$ if x is the right-most informer in a bucket, then we get exactly the same potential function argument as before but with U replaced by F . We may therefore conclude:

Theorem 3.1 *The winding takes $O(N \log^2(U/N))$ time. Each of the lists involved has length $O(U/N)$, and the total number of list operations is $O(N \log U/N)$.*

4 Unwinding

In this section, we will show how to unwind all left informers in $O(N \log(U/N)(\log U/N + \log P))$ time. The right informers may be treated symmetrically, and they will not be discussed further. We are

now in the situation that we have shifted all the left informers into the alphabet, and we will now shift them back into their original positions, but loaded with information about the longest pattern substring extending to the left of their final position.

During the unwinding, if at some stage an informer is at position u , then it should be loaded with the longest suffix of $\mathcal{T}[0, u]$ which is a substring of \mathcal{P} . We will use a function $\text{LEFTAPPEND}(\alpha, \beta)$ which, given substrings α, β of the pattern, returns the longest suffix of $\alpha\beta$ which is a substring of the pattern. In [7], an algorithm for LEFTAPPEND was given which preprocesses a string of length P in $O(P)$ time and answers such queries in $O(\log P)$ time.

The basic idea is to log the winding of Algorithm A, so that we can reverse it directly in the unwinding. Thus, consider step s of the winding, where we first cut \mathcal{L} at U_s into lists \mathcal{L}_0 and \mathcal{L}_1 , second subtract $U_s - P_s$ from u for all informers in \mathcal{L}_1 , and third merge \mathcal{L}_0 with \mathcal{L}_1 into \mathcal{L} , removing duplicates. In the unwinding of this step, our starting point is the list \mathcal{L} with all informers at positions to the left of U_s , and each informer x being loaded with the longest suffix of $\mathcal{T}[0, u(x) - 1]$ which is a substring of \mathcal{P} . Moreover, we have the information that allows us to split \mathcal{L} into \mathcal{L}_0 and \mathcal{L}_1 . Being a bit more precise, we will see \mathcal{L} as divided into elements L_0, \dots, L_{m_s-1} where each L_i is either a segment of \mathcal{L}_0 , a segment of \mathcal{L}_1 , or an informer that should be copied to both \mathcal{L}_0 and \mathcal{L}_1 . From our potential function argument we know that the sum of the m_s over all steps is $O(N \log U/N)$. In order to use this bound, we need to be able to process each segment L_i in a time which is independent of the number of informers it contains. This independence is achieved by sparsification of the loading of the informers as follows.

Consider informers x_1, x_2, x_3 at positions u_1, u_2, u_3 where $u_1 < u_2 < u_3$. For $i = 1, 2, 3$, let α_i be the longest suffix of $\mathcal{T}[0, u_i]$ which is a substring of \mathcal{P} . Suppose that $u_1 + |\alpha_3| \geq u_3$. Then $\alpha_2 = \text{LEFTAPPEND}(\alpha_1, \alpha_3[u_1 + |\alpha_3| - u_3, |\alpha_3| + u_2 - u_3 - 1])$, and then we say that x_1 and x_3 cover x_2 . Now in such a case, we need not explicitly keep the information at x_2 up to date at all times since we can quickly recover such information as needed by one call to LEFTAPPEND . For a sequence of informers, we will then keep some informers loaded, that is, up to date, and other informers unloaded. We will say any such sequence is *proper* if each informer is either loaded or covered by loaded informers. If a sequence of informers is properly loaded and minimal, in the sense that no loaded informer is covered by two other loaded informers, then we will say the sequence is *economical*. In an economically loaded informer sequence, we will keep the positions of loaded informers in a balanced binary tree. This will allow us to update the list (with insertions, splits, concatenations, etc.) in $O(\log U/N)$ time, since by the arguments above, no list will ever be of length greater than U/N . Also, each unloaded informer will be able to find its covering informers in time $O(\log U/N)$. Thus, we will be able to load any informer with one tree lookup and one call to LEFTAPPEND .

The following lemmas allow us to manipulated economically loaded informer sequences.

Lemma 4.1 *Let \mathcal{L}_0 and \mathcal{L}_1 be lists of informers, each of which has economical coverings. Then we can produce an economical covering of their concatenation $\mathcal{L}_0\mathcal{L}_1$ by changing $O(1)$ informers.*

Proof: We may need to unload the last informer of \mathcal{L}_0 and the first in \mathcal{L}_1 . But we can do so by checking to see if they are covered in the new concatenated list. So we change the status of at most 2 informers. ■

Lemma 4.2 *Let \mathcal{L} be a list of informers with an economical covering. Let \mathcal{L}_0 and \mathcal{L}_1 be the lists produced by splitting \mathcal{L} at u . Then we can produce economical coverings for \mathcal{L}_0 and \mathcal{L}_1 changing $O(1)$ informers.*

Proof: Let x be the first informer in \mathcal{L}_1 . Suppose that x is not loaded. Load x using its covering in \mathcal{L} . Let x' be the loaded neighbor to the right of x , and let x'' be the loaded neighbor to the right of x' . If x' is covered by x and x'' , unload x' . Finally, update the list of loaded informers. The last informers in \mathcal{L}_0 is dealt with symmetrically, and in conclusion, we only make a constant number of informer tree operations and calls to `LEFTAPPEND` in the processing of \mathcal{L}_0 and \mathcal{L}_1 . ■

Let us return to our sequence \mathcal{L} which was divided into elements L_0, \dots, L_{m_s-1} where each L_i is either a segment of \mathcal{L}_0 , a segment of \mathcal{L}_1 , or a single informer that should be copied to both \mathcal{L}_0 and \mathcal{L}_1 . Using Lemma 4.2 we get an economical loading for each segment L_i , and using Lemma 4.1 we concatenate the L_i s into \mathcal{L}_0 and \mathcal{L}_1 . Thereby, we get a total of $O(m_s)$ informer load changes.

The next part of redoing step s is to add $U_s - P_s$ to u for the informers in \mathcal{L}_1 . This should change the specified contents of some of the informers, but we claim that among the loaded informers, it is only the two leftmost that get affected. The loading of an informer x in \mathcal{L}_1 should only change if it extends to U_s or beyond after the addition to the u -values, that is if $u(x) - |t(x)| + 1 \leq U_s$. In this case, recall that we already have a correctly loaded left informer y at position $U_s - 1$. Thus, in order to get a correct loading of x , we set $t(x) := \text{LEFTAPPEND}(t(y), \alpha)$ where α is the suffix of $t(x)$ of length $u(x) - U_s + 1$. Now, let x, x' be the two leftmost loaded informers in \mathcal{L}_1 , and let x'' be any loaded informer further to the right. Since \mathcal{L}_1 is economic, x' is not covered by x and x'' , so $u(x'') - |t(x'')| > u(x) = U_s$. Thus the loading of x'' need not be changed.

The final part in redoing step s is to concatenate \mathcal{L}_0 and \mathcal{L}_1 , again by Lemma 4.1. In conclusion, the unwinding is done with $O(N \log U/N)$ list operations and calls to `LEFTAPPEND`.

Theorem 4.3 *We can load all informers in time $O(N \log(U/N)(\log U/N + \log P))$.*

In [7], an algorithm was given to compute if a pattern occurs within the concatenation of two of its substrings. This operation takes $O(P)$ preprocessing on a string of length P and answers such queries in $O(\log P)$ time. This result suffices for the following.

Corollary 4.4 *We can do LZ77 compressed matching on a string of length N which represents a string of length U with a pattern of length P in time $O(N \log(U/N)(\log U/N + \log P))$.*

In the following section, we show what to do if P is so long that $\log P = \omega(\log U/N)$.

5 Unwinding Long Patterns

In this section, we will address the problem of unwinding for long patterns. We will consider the case where $P \geq F^5$. Note that Corollary 4.4 settles Theorem 1.1 without randomization for cases where $P = O(F^5)$.

Recall that in §3.3, step 2, we introduced an artificial cut at each position which is a multiple of F . Define k such that $(k+1)F \leq P < (k+2)F$. Consider some pattern occurrence $\mathcal{P} = \mathcal{T}[j; j+P-1]$, and let i be minimal such that $iF \geq j$, i.e. $i = F \lceil j/F \rceil$. Then \mathcal{P} can be written as BME where

following the definition of the compression scheme, $\text{char}(x)$ is not changed once it has been set.

Our task is to reverse steps A.1.1–A.1.3 in Algorithm A, or, in fact, we only need to change the contents of the informers in connection with steps A.1.1–A.1.2. Let \mathcal{L} , \mathcal{L}_1 , \mathcal{L}_2 , U_s and P_s be as in Algorithm A. Thus, we want to add $U_s - P_s$ to the positions of all informers in \mathcal{L}_2 and then set $\mathcal{L} := \mathcal{L}_1\mathcal{L}_2$. First, concerning the base case, where \mathcal{L}_2 consists of a single informer x which is moved from the alphabet at the negative position $-C$. Hence we set $\text{char}(x) := C \bmod \wp$. If $s = 0$, we set $\text{pref}(x) := \text{char}(x)$. Otherwise, we note that we have a left informer l_{s-1} in \mathcal{L}_1 at position $U_s - 1$, and hence we set

$$\text{pref}(x) := \text{pref}(l_s) + \text{char}(x)\sigma^{U_s-1} \bmod \wp.$$

We are now ready for the real case where \mathcal{L}_2 is a sequence of informers in the text with positive u-values. We will exploit that we have a left informer l_{s-1} at position $U_s - 1$ in \mathcal{L}_1 and that the right-most informer r_s of \mathcal{L}_2 is going to be put at position U_s . The needed transformation for each informer x in \mathcal{L}_2 is now:

- $\text{pref}(x) := \text{pref}(x) - \text{pref}(r_s) + \text{char}(r_s)\sigma^{P_s-1} \bmod \wp.$
- $\text{pref}(x) := \text{pref}(x)\sigma^{U_s-P_s} \bmod \wp.$
- $\text{pref}(x) := \text{pref}(l_{s-1}) + \text{pref}(x) \bmod \wp.$

This procedure depends on a data structure on trees which supports the following operations: in time $O(\log \mathbb{F})$ it should be able to multiply or add to pref for all elements in a list of length \mathbb{F} , and it should support cuts and concatenations on such lists. See cutting and linking trees in [12]. As a result, the work connected with maintaining pref and char is $O(N \log \mathbb{F})$. Combining with Theorem 3.1 and Corollary 4.4, the whole compressed matching is done in time $O(N \log^2 \mathbb{F} + P)$, thus establishing Theorem 1.1 with a Monte Carlo algorithm for the case without self-references.

Finally, we need to show how to check our output, thus changing the Monte Carlo algorithm above to a Las Vegas algorithm. We also need to show how to deal with self-references. Once these technical steps are performed, we will have completed the proof of Theorem 1.1.

6 Checking the Answer

The randomized algorithm from Section 5 for the case of long patterns is of the Monte Carlo type, that is, it has a guaranteed running time but may sometimes report pattern occurrence where none exists. In order to guarantee the correctness of the output, at the risk of increasing the running time, we have to test any pattern occurrence found by the algorithm. A simple $O((P + N) \log^2(U/N))$ solution to the checking problem would be to place an informer on each of the P positions of the potential pattern occurrence. We could then wind these informers back into the alphabet in $O((P + N) \log^2(U/N))$ time, thereby checking each character of the potential occurrence. However, it is possible to get an $O(P + N \log^2(U/N))$ checking algorithm. Generally the idea is that we load the informers with pattern substrings corresponding to the potential pattern occurrence. Then we start winding. As the winding proceeds, some of these substrings will be cut and distributed over several informers, and others will start overlapping. In the latter case, we have to check that the

pattern substrings agree in the overlap. Such a check can be done in constant time, after an $O(\mathbb{P})$ time preprocessing, via standard suffix tree and least common ancestors techniques. At the end, the surviving informers will end up in the alphabet, where we have to check that the characters match.

To be more concrete, we want to check the long middle of a pattern occurrence $\mathcal{P}[j, j + k\mathbb{F}] = \mathcal{T}[i\mathbb{F}, (i + k)\mathbb{F} - 1]$ suggested by technique from Section 5. First, before the winding, for $a = 1, \dots, k$, we load the left informer x at position $(i + a)\mathbb{F} - 1$ with $\mathfrak{t}(x) = \mathcal{P}[j + (a - 1)\mathbb{F}, j + a\mathbb{F}]$. No other left informers are loaded. Then we start winding following Algorithm A from Section 3. At the beginning of any iteration, we have a sequence x_1, \dots, x_l of loaded informers. Our suggested pattern occurrence is correct if and only if $\mathfrak{t}(x_i)$ is a suffix of $\mathcal{T}[0, u(x_i)]$ for all i . The covering is kept *non-overlapping* in the sense that for all $i > 1$, $u(x_i) - |\mathfrak{t}(x_i)| \geq u(x_{i-1})$.

Consider a cut in step A.1.1. Let x be the first informer in \mathcal{L}_1 . Suppose $\mathfrak{t}(x)$ extends past U_s , i.e. $u(x) - |\mathfrak{t}(x)| + 1 < U_s$. Since the covering is non-overlapping, this implies that the left informer y at position $U_s - 1$ is unloaded. Let $\alpha\beta = \mathfrak{t}(x)$ where $|\alpha| = u(x) - U_s + 1$. Set $\mathfrak{t}(x) = \beta$ and $\mathfrak{t}(y) = \alpha$. In this way, each cut loads at most one informer. The above cut is implemented in time $O(\log \mathbb{F})$, giving a total cost of $O(N \log \mathbb{F})$.

The merging of \mathcal{L}_0 and \mathcal{L}_1 in step A.1.3, consists of a series of segment merges, where we take a segment (x_a, \dots, x_b) of loaded informers from one list, and insert after a loaded informer y from the other list. Now *match* x_a with y as follows. If $\mathfrak{t}(x_a)$ extends beyond $u(y)$, i.e. if $u(x_a) - |\mathfrak{t}(x_a)| < u(y)$, let $\alpha\beta = \mathfrak{t}(x_a)$ where $|\beta| = u(x_a) - u(y)$. Check in constant time that the shorter of α and $\mathfrak{t}(y)$ is a suffix of the other. If not, report that the pattern occurrence was incorrect. Otherwise, if α is shorter than $\mathfrak{t}(y)$, set $\mathfrak{t}(x_a) = \beta$. If α is longer than $\mathfrak{t}(y)$, unload y and match x_a with y 's loaded predecessor.

Ignoring the repeated unloading from the recursive matching, each segment merge is implemented in time $O(\log \mathbb{F})$, and from our potential function argument, we know that there are at most $O(N \log \mathbb{F})$ segment merges, amounting to a total cost of $O(N \log^2 \mathbb{F})$.

Attributing a cost of $O(\log \mathbb{F})$ to the segment merge, the repeated unloading can be done in time linear in the number of unloads. We start with $< \mathbb{P}$ loaded informers, and cuts result in the loading of $O(N)$ informers. Hence the number of unloadings is bounded by $O(\mathbb{P} + N)$, so total cost of merging is $O(\mathbb{P} + N \log^2 \mathbb{F})$.

At the end of the winding, we have a sequence x_1, \dots, x_l of informers in the alphabet. Recall that they are all loaded with sub-patterns of \mathcal{P} , and we now have to check that this loading match the alphabet. Since the alphabet does not repeat characters and since the covering is non-overlapping, the loadings must correspond to disjoint sub-strings of \mathcal{P} . In particular, this implies that if $\sum_i |\mathfrak{t}(x_i)| > |\mathcal{P}|$, we may report that the suggested pattern occurrence is incorrect. Otherwise, in time $O(\mathbb{P})$ it is trivial to check the loadings. That is, the last checking is done in time $O(\mathbb{P})$.

In conclusion, we have checked our pattern occurrence in time $O(\mathbb{P} + N \log^2 \mathbb{F}) = O(\mathbb{P} + N \log^2(U/N))$, thus achieving a Las Vegas type randomization.

7 Self-reference

Referring to the symbols as they were defined in Sections 1–2, for $i = 0, \dots, N - 1$, let $D_i = U_i - P_i$. So far we have required that our compressed text \mathcal{Z} is non-referential, i.e. that $D_i \geq L_i$, but now

we will show how to deal with the general case without this restriction. Thus, our only requirement is that $D_i > 0$, as in the original definition of Lempel-Ziv compressed texts.

Note the following alternative but equivalent statement of Eq. 1 from Section 1:

$$\begin{aligned} \forall i = 0, \dots, N-1 : & \\ \mathcal{T}[U_i, U_i + L_i - 1] \text{ is a prefix of } \mathcal{T}[P_i, P_i + D_i - 1]^\infty & \end{aligned} \quad (2)$$

Winding

With Eq. 2 in mind, we can easily generalize our winding algorithm to deal with the general case.

Algorithm C Winding algorithm allowing self-references.

C.1. For $s := N - 1$ downto 0 do

C.1.1. While the last informer x in \mathcal{L} has $u(x) \geq U_s$ do:

C.1.1.1. Cut \mathcal{L} at U_s into lists \mathcal{L}_0 and \mathcal{L}_1 .

C.1.1.2. Let x be the first informer in \mathcal{L}_1 , and set $k = \lfloor (u(x) - U_s) / D_s \rfloor$.

C.1.1.3. Subtract $(k + 1)D_s$ from u for all informers in \mathcal{L}_1 .

C.1.1.4. Merge \mathcal{L}_0 with \mathcal{L}_1 into \mathcal{L} , removing informers that are duplicate with respect to u -values.

Above, if $k = 0$, we are following Eq. 1 directly. If $k > 0$, we are just skipping loops with superfluous merges where the first informer in \mathcal{L}_1 lands strictly before the last informer in \mathcal{L}_0 . Extending the potential function argument from Section 3, we will show that with this little optimization, the general winding of Algorithm C is performed within the same time bounds as we had for the winding without self-references in Algorithm A.

For each of the $O(N)$ iterations of the outer for-loop (step C.1), we will show that Π increases by at most $O(\log \mathbb{F})$. Moreover, we will show that each iteration of the inner while-loop (step C.1.1) decreases Π by $\Omega(1)$. Since Π starts at $O(N \log \mathbb{F})$, this will allow us to conclude that we have a total of $O(N \log \mathbb{F})$ while-loop iterations.

Consider a specific for-loop iteration (step C.1), fixing some value of s . To facilitate the proof, we insert two auxiliary informers a and b where $u(a) = P_s$ and $u(b) = U_s - 1$. These will be removed again at the end of the for-loop, if they have not already been removed as duplicates. Note that b will always be the last informer in \mathcal{L}_0 . Technically, we fix $d^+(b) = 2\mathbb{F}$. Hence $\pi(b)$ is not affected when cutting \mathcal{L} into \mathcal{L}_0 and \mathcal{L}_1 . Inserting or removing an informer affects the potential of at most two other informers; namely its neighbors. Thus our insertion and removal of a and b increases Π by $O(\log \mathbb{F})$, and this will capture the whole increase to Π during the for-loop.

Consider an iteration of the while-loop (step C.1.1). Let x be the first informer in \mathcal{L}_1 . Since $d^+(b)$ is fixed, $d^-(x)$ is the only $d^{+/-}$ distance that could possibly increase during the while-loop. However, before the while-loop, $d^+(x)$ is the distance to b , but the distance to a is at most $D_s - 1$ bigger, and step C.1.1.3 moves x a positive multiple of D_s closer to a . Thus $d^+(x)$ is, in fact, decreased during the while-loop, that is, no informer increases potential during the while-loop.

If x lands on another informer, an informer gets removed, and Π decreases by $\Omega(1)$, as desired. Otherwise, x lands between a and b . Let y and z , $u(y) < u(z)$, be its nearest neighbors from \mathcal{L}_0

after the merge. Then, either $d^+(y)$ or $d^-(z)$ is halved by the merge, so we conclude that Π is decreased by $\Omega(1)$ during the while-loop.

The argument from Section 3, that each recursive call to `Segment-Merge` decreases Π by $\Omega(1)$, still holds. Thus, the total number of while-loops and recursive calls to `Segment-Merge` is bounded by the initial value plus the total increase of Π , which is $O(N \log F)$. As in Section 3, this again bounds the total number of list operations, each of which takes time $O(\log F)$. We may therefore conclude that the general winding of Algorithm C, which allows self-reference, satisfies the description in Theorem 3.1.

Unwinding

For the unwinding, we can essentially follow the pattern from Section 4. Again, we are only focussing on left informers as right informers may be treated symmetrically. The only critical point is what happens to the loading of the first left informer x in \mathcal{L}_1 when redoing the shift in step C.1.1.3. Let s , x , \mathcal{L} , \mathcal{L}_0 , \mathcal{L}_1 , and k be fixed according to a specific iteration of the shift in step C.1.1.3.

As in Section 4, we already have a correctly loaded left informer y at position $U_s - 1$. Let $u^o(x)$ and $u^n(x)$ denote the position of x before and after the shift. If $|\mathfrak{t}(x)| < u^o(x) - P_s + 1$, $\mathfrak{t}(x)$ does not need to be changed, so suppose that this is not the case. Let α be the suffix of $\mathfrak{t}(x)$ of length $u^o(x) - P_s + 1$. If $|\mathfrak{t}(y)| < D_s$, as in Section 4, we set $\mathfrak{t}(x) := \text{LEFTAPPEND}(\mathfrak{t}(y), \alpha)$. However, suppose $|\mathfrak{t}(y)| \geq D_s$. Let $\gamma\beta = \mathfrak{t}(y)$ where $|\beta| = D_s$. Then $\mathfrak{t}(x) := \text{LEFTAPPEND}(\gamma\beta^{k+1}, \alpha)$. The pattern preprocessing for `LEFTAPPEND` from [7] allows queries in time $O(\log P)$ despite the β^{k+1} in the argument.

The unwinding of long patterns from Section 5 is changed similarly, so this completes the proof of Theorem 1.1.

References

- [1] A.V. Aho, J.E. Hopcroft, and J.D. Ullman. *The Design and Analysis of Computer Algorithms*. Addison-Wesley, 1974.
- [2] A. Amir and G. Benson. Efficient two dimensional compressed matching. *Proc. of the 2nd IEEE Data Compression Conference*, pages 279–288, Mar 1992.
- [3] A. Amir, G. Benson, and M. Farach. Let sleeping files lie: Pattern matching in Z-compressed files. *Proc. of the 5th Annual ACM-SIAM Symposium on Discrete Algorithms*, 1994.
- [4] A. Amir, G. Benson, and M. Farach. Optimal two-dimensional compressed matching. *Proc. of 21st International Colloquium on Automata Languages and Programming*, 1994.
- [5] Corman, Leiserson, and Rivest. *Introduction to Algorithms*. MIT Press, 1990.
- [6] M. Farach and S. Muthukrishnan. Optimal parallel dictionary matching and compression. *7th Annual ACM Symposium on Parallel Algorithms and Architectures*, 1995.

- [7] M. Gu, M. Farach, and R. Beigel. An efficient algorithm for dynamic text indexing. *Proc. of the 5th Annual ACM-SIAM Symposium on Discrete Algorithms*, 1994.
- [8] R.M. Karp and M.O. Rabin. Efficient randomized pattern-matching algorithms. *IBM Journal of Research and Development*, 31:249–260, 1987.
- [9] S. Rao Kosaraju. On the entropy estimation of low entropy sources. Personal Communication, 1994.
- [10] S. Rao Kosaraju. Pattern matching in compressed texts. *Proc. of the 10th Annual Conference on Foundations of Software Technology and Theoretical Computer Science (FST&TCS)*, 1995.
- [11] A. Lempel and J. Ziv. On the complexity of finite sequences. *IEEE Transactions on Information Theory*, 22:75–81, 1976.
- [12] D. Sleator and R. Tarjan. A data structure for dynamic trees. *Journal of Computer and System Sciences*, 24, 1983.
- [13] L. Stauffer and D. Hirschberg. Pram algorithms for static dictionary compression. *8th International Parallel Processing Symposium*, 1994.
- [14] J. Storer. *Data compression: methods and theory*. Computer Science Press, Rockville, Maryland, 1988.
- [15] T. A. Welch. A technique for high-performance data compression. *IEEE Transactions on Computers*, 17:8–19, 1984.
- [16] J. Ziv and A. Lempel. A universal algorithm for sequential data compression. *IEEE Transactions on Information Theory*, IT-23(3):337–343, 1977.