# Lempel-Ziv Parsing in External Memory[⋆]

Juha Kärkkäinen, Dominik Kempa, and Simon J. Puglisi

Department of Computer Science,
University of Helsinki
Helsinki, Finland
{firstname.lastname}@cs.helsinki.fi

**Abstract.** For decades, computing the LZ factorization (or LZ77 parsing) of a string has been a requisite and computationally intensive step in many diverse applications, including text indexing and data compression. Many algorithms for LZ77 parsing have been discovered over the years; however, despite the increasing need to apply LZ77 to massive data sets, no algorithm to date scales to inputs that exceed the size of internal memory. In this paper we describe the first algorithm for computing the LZ77 parsing in external memory. Our algorithm is fast in practice and will allow the next generation of text indexes to be realised for massive strings and string collections.

## 1  Introduction

For over three decades the Lempel-Ziv (LZ77) factorization [21] has been a fundamental tool for compressing data [7,19,8,17] and for string processing – in particular for the efficient detection of periodicities [2,16,15,14]. Recently the factorization has become the basis for several compressed full-text self-indexes [18,10,9]. These indexes are designed to support efficient storage and fast searching of massive, highly repetitive data sets such as web documents, whole genome collections, and versioned collections of source code and multi-author documents, such as Wikipedia.

In order for these LZ77-based self-indexes to be constructed, whole collection LZ77 factorizations need to be computed. However, to our knowledge, all current LZ77 algorithms require large amounts of memory and are essentially "in-memory" algorithms: they have poor locality of memory reference, do not scale to external memory (disk), and so are incapable of factorizing massive strings. In this paper we address this shortcoming and design an LZ77 factorization algorithm capable of scaling to data that substantially exceeds the size of RAM.

*Our contribution.* We have designed and implemented the first external memory algorithm for LZ77 factorization. It is based on the recent LZscan algorithm [12],

---

which itself is a space-efficient algorithm, able compute the LZ factorization of a string of length $n$ using less than $2n$ bytes of memory. The new external memory version, EM-LZscan scales beyond the main memory size. Furthermore, our experiments show that EM-LZscan is significantly faster than LZscan already for much smaller files.

Theoretically, the time complexity of EM-LZscan is $O(n^2 t_{rank}/M)$, where $M$ is the size of the main memory and $t_{rank}$ is the time complexity of the rank operation on strings (see, e.g., [3]). The disk accesses are sequential with a total I/O volume of $O(n^2/M)$. The quadratic time complexity means that the practical scalability of the algorithm is limited but still beyond any previous algorithm. A couple of days on an ordinary desktop PC with 4GiB of RAM is sufficient for computing the LZ factorization of any 10GiB file or a 40GiB highly repetitive file, which is the particularly relevant case for LZ77-based indexes.

*Related work.* A recent survey [1] and some even more recent papers [13,12] outline the many algorithms for LZ77 factorization, all of which operate in RAM and most of which make use of the suffix array (SA) and longest-common-prefix (LCP) array, as intermediate data structures. Many of the algorithms compute the *LPF tables*, which contain all the longest previous factors, not just the ones needed for the LZ parsing. Several LPF algorithms ([6] is a notably simple algorithm in this category) compute the LPF tables in a single left-to-right pass over the SA and LCP arrays, with sublinear (in fact only $O(\sqrt{n})$ words) extra memory. These algorithms would thus seem to be candidates for external memory LZ77 factorization, however it is the computation of SA and LCP arrays that forms the bottleneck. There exists an external memory implementation for computing the SA and LCP arrays [4], but it is limited by the large disk space requirement of $54n$ bytes (compared to less than $2n$ bytes for EM-LZscan).

Another way to reduce working space is to use compressed suffix arrays, augmented with some auxilliary data structures [20,17,12]. However, the experiments in [12] show that such algorithms are inferior to LZscan .

## 2   Basic Notation and Algorithmic Machinery

*Strings.* Throughout we consider a string $X = X[1..n] = X[1]X[2]\ldots X[n]$ of $|X| = n$ symbols drawn from the alphabet $[0..\sigma - 1]$. We assume $X[n]$ is a special "end of string" symbol, \$, smaller than all other symbols in the alphabet. The reverse of $X$ is denoted $\hat{X}$. For $i = 1, \ldots, n$ we write $X[i..n]$ to denote the *suffix* of $X$ of length $n - i + 1$, that is $X[i..n] = X[i]X[i + 1]\ldots X[n]$. We will often refer to suffix $X[i..n]$ simply as "suffix $i$". Similarly, we write $X[1..i]$ to denote the *prefix* of $X$ of length $i$. $X[i..j]$ is the *substring* $X[i]X[i + 1]\ldots X[j]$ of $X$ that starts at position $i$ and ends at position $j$. By $X[i..j)$ we denote $X[i..j - 1]$. If $j < i$ we define $X[i..j]$ to be the empty string, also denoted by $\varepsilon$.

*LZ77.* Before defining the LZ77 factorization, we introduce the concept of a *longest previous factor* (LPF). The LPF at position $i$ in string $X$ is a pair

$\mathsf{LPF_X}[i] = (p_i, \ell_i)$ such that, $p_i < i$, $\mathsf{X}[p_i..p_i + \ell_i) = \mathsf{X}[i..i + \ell_i)$, and $\ell_i$ is maximized. In other words, $\mathsf{X}[i..i + \ell_i)$ is the longest prefix of $\mathsf{X}[i..n]$ which also occurs at some position $p_i < i$ in $\mathsf{X}$. Note also that there may be more than one potential source (that is, $p_i$ value), and we do not care which one is used.

The LZ77 factorization (or LZ77 parsing) of a string $\mathsf{X}$ is then just a greedy, left-to-right parsing of $\mathsf{X}$ into longest previous factors. More precisely, if the $j$th LZ factor (or *phrase*) in the parsing is to start at position $i$, then we output $(p_i, \ell_i)$ (to represent the $j$th phrase), and then the $(j + 1)$th phrase starts at position $i + \ell_i$. The exception is the case $\ell_i = 0$, which happens iff $\mathsf{X}[i]$ is the leftmost occurrence of a symbol in $\mathsf{X}$. In this case we output $(\mathsf{X}[i], 0)$ (to represent $\mathsf{X}[i..i]$) and the next phrase starts at position $i + 1$. When $\ell_i > 0$, the substring $\mathsf{X}[p_i..p_i + \ell_i)$ is called the *source* of phrase $\mathsf{X}[i..i + \ell_i)$. We denote the number of phrases in the LZ77 parsing of $\mathsf{X}$ by $z$.

*Matching Statistics.* Given two strings $\mathsf{Y}$ and $\mathsf{Z}$, the matching statistics of $\mathsf{Y}$ w.r.t. $\mathsf{Z}$, denoted $\mathsf{MS_{Y|Z}}$, is an array of $|\mathsf{Y}|$ pairs, $(p_1, \ell_1)$, $(p_2, \ell_2)$, ..., $(p_{|\mathsf{Y}|}, \ell_{|\mathsf{Y}|})$, such that for all $i \in [1..|\mathsf{Y}|]$, $\mathsf{Y}[i..i + \ell_i) = \mathsf{Z}[p_i..p_i + \ell_i)$ is the longest substring starting at position $i$ in $\mathsf{Y}$ that is also a substring of $\mathsf{Z}$. The observant reader will note the resemblance to the LPF array. Indeed, if we replace $\mathsf{LPF_Y}$ with $\mathsf{MS_{Y|Z}}$ in the computation of the LZ factorization of $\mathsf{Y}$, the result is the relative LZ factorization of $\mathsf{Y}$ w.r.t. $\mathsf{Z}$ [19].

## 3    LZ77 Factorization in External Memory

In this section we first describe the scanning-based, block-oriented LZ77 factorization algorithm called LZscan that was introduced in [12]. We then present the modifications to LZscan to make it run efficiently in external memory.

### 3.1    Basic Algorithm

Conceptually LZscan divides $\mathsf{X}$ up into $d = \lceil n/b \rceil$ fixed size blocks of length $b$: $\mathsf{X}[1..b]$, $\mathsf{X}[b + 1..2b]$, ... . In the description that follows we will refer to the block currently under consideration as $\mathsf{B}$, and to the prefix of $\mathsf{X}$ that ends just before $\mathsf{B}$ as $\mathsf{A}$. Thus, if $\mathsf{B} = \mathsf{X}[kb + 1..(k + 1)b]$, then $\mathsf{A} = \mathsf{X}[1..kb]$.

To begin, we will assume no LZ factor or its source crosses a boundary of the block $\mathsf{B}$. Later we will show how to remove these assumptions.

The outline of the algorithm for processing a block $\mathsf{B}$ is shown below.

1. Compute $\mathsf{MS_{A|B}}$
2. Compute $\mathsf{MS_{B|A}}$ from $\mathsf{MS_{A|B}}$, $\mathsf{SA_B}$ and $\mathsf{LCP_B}$
3. Compute $\mathsf{LPF_{AB}}[kb + 1..(k + 1)b]$ from $\mathsf{MS_{B|A}}$ and $\mathsf{LPF_B}$
4. Factorize $\mathsf{B}$ using $\mathsf{LPF_{AB}}[kb + 1..(k + 1)b]$

*Step 1: Computing Matching Statistics.* Similarly to most algorithms for computing the matching statistics, we first construct some data structures on $\mathsf{B}$ and then scan $\mathsf{A}$. For the details of the data structures we refer to [12]. The key properties are the space requirement of $27b$ bytes and linear time construction.

The scanning of $\mathsf{A}$ is the computational bottleneck of the algorithm in theory and practice. Theoretically, the time complexity of Step 1 is $O((|\mathsf{A}| + |\mathsf{B}|)t_{\mathsf{rank}})$, where $t_{\mathsf{rank}}$ is the time complexity of the rank operation on strings over the alphabet $\Sigma$ (see, e.g., [3]). Thus the total time complexity of $\mathsf{LZscan}$ is $O(dnt_{\mathsf{rank}})$. In practice, each step of the scan may involve a substantial amount of work in navigating the data structures. However, each character of $\mathsf{A}$ is accessed only once, and this is mostly done sequentially from right to left.

An important optimization, called skipping trick, speeds up the computation for highly repetitive inputs. It takes advantage of repetition present in $\mathsf{A}$ that was found in the previous stages of the algorithm. Consider an LZ factor $\mathsf{A}[i..i + \ell)$. Because, by definition, $\mathsf{A}[i..i+\ell)$ occurs earlier in $\mathsf{A}$ too, any source of an LZ factor of $\mathsf{B}$ that is completely inside $\mathsf{A}[i..i + \ell)$ could be replaced with an equivalent source in that earlier occurrence. Thus such factors can be skipped during the computation of $\mathsf{MS}_{\mathsf{A|B}}$ without an effect on the factorization.

More precisely, if during the scan we compute $\mathsf{MS}_{\mathsf{A|B}}[j] = (p, k)$ and find that $i \leq j < j+k \leq i+\ell$ for an LZ factor $\mathsf{A}[i..i+\ell)$, we will compute $\mathsf{MS}_{\mathsf{A|B}}[i-1]$ and continue the scanning from $i - 1$. However, we will do this only for long phrases with $\ell \geq 40$. To compute $\mathsf{MS}_{\mathsf{A|B}}[i - 1]$ from scratch, we use right extension operations implemented by a binary search on $\mathsf{SA}$. This is the only situation, where a part of $\mathsf{A}$ is scanned from left to right, but still sequentially.

*Step 2: Inverting Matching Statistics.* With the help of $\mathsf{SA}_{\mathsf{B}}$ and $\mathsf{LCP}_{\mathsf{B}}$, we can *invert* $\mathsf{MS}_{\mathsf{A|B}}$ to obtain $\mathsf{MS}_{\mathsf{B|A}}$, which is what we need for LZ77 factorization. Again, we refer to [12] for details of the inversion algorithm and give only the key properties. The algorithm accesses each entry of $\mathsf{MS}_{\mathsf{A|B}}$ (except those skipped by the skipping trick) once, in an arbitrary order, and processes the entry in constant time. Thus we do not need to store $\mathsf{MS}_{\mathsf{A|B}}$ but can process each entry as soon as it is produced in Step 1. The rest of the computation takes $O(b)$ time.

*Step3: Computing LPF.* Consider the pair $(p, \ell) = \mathsf{LPF}_{\mathsf{AB}}[i]$ for $i \in [kb + 1..(k + 1)b]$ that we want to compute and assume $\ell > 0$ (otherwise $i$ is the position of the leftmost occurrence of $\mathsf{X}[i]$ in $\mathsf{X}$, which we can easily detect). Clearly, either $p \leq kb$ and $\mathsf{LPF}_{\mathsf{AB}}[i] = \mathsf{MS}_{\mathsf{B|A}}[i]$, or $kb < p < i$ and $\mathsf{LPF}_{\mathsf{AB}}[i] = (kb + p_{\mathsf{B}}, \ell_{\mathsf{B}})$, where $(p_{\mathsf{B}}, \ell_{\mathsf{B}}) = \mathsf{LPF}_{\mathsf{B}}[i - kb]$. Thus computing $\mathsf{LPF}_{\mathsf{AB}}$ from $\mathsf{MS}_{\mathsf{B|A}}[i]$ and $\mathsf{LPF}_{\mathsf{B}}$ is easy.

The above is true if the sources do not cross the block boundary, but the case where $p \leq kb$ but $p + \ell > kb + 1$ is not handled correctly. An easy correction is to replace $\mathsf{MS}_{\mathsf{A|B}}$ with $\mathsf{MS}_{\mathsf{AB|B}}[1..kb]$ in all of the steps. This does not affect the essential features of the algorithm.

*Step 4: Parsing.* We use the standard LZ77 parsing to factorize $\mathsf{B}$ except $\mathsf{LPF}_{\mathsf{B}}$ is replaced with $\mathsf{LPF}_{\mathsf{AB}}[kb + 1..(k + 1)b]$.

So far we have assumed that every block starts with a new phrase, or, put another way, that a phrase ends at the end of every block. Let $Z = X[i..(k+1)b]$ be the last factor in B after we have factorized B as described above. This may not be a true LZ factor when considering the whole X because the true LZ factor may continue beyond the end of B. If $|Z| \leq b/2$, we start the next block at $i$ instead of $(k+1)b+1$, and compute the true phrase starting at $i$ while processing that block. This at most doubles the computation. If $|Z| > b/2$, we need to do something more sophisticated. In [12], a modified constant extra space pattern matching algorithm by Crochemore [5] is used for finding the true phrase.

## 3.2   Implementation in External Memory

Next we describe an external memory adaptation of LZscan called EM-LZscan.

*Block data structures.* The structures constructed for the current block B are essentially the same in EM-LZscan as in LZscan and are kept in memory during the processing of B. There are two notable differences: B itself is read from disk and held in memory during the stage, and we replace 32-bit integers with 40-bit integers to represent positions in the whole text (but still 32-bit integers for positions in B). These changes raise the peak memory usage of the data structures from $27b$ bytes to $29b$ bytes. We have also implemented a 32-bit version of EM-LZscan that needs $28n$ bytes of space and runs slightly faster because processing 40-bit integers incurs a small overhead.

*Scanning* A. The scanning of A is performed by reading A from disk into a buffer of size 256KiB. We store X in reverse order on disk so that the backward scan of A involves reading in forward direction. This seems to make the algorithm faster even when the time for reversing X is included.

To implement the skipping trick, we need to identify phrases of length 40 or more during the scan. They are stored in a separate file in reverse sequential order, which is then scanned in synchrony with A. Since the file grows backwards during the computation, we create a file of the maximum size, which is $n/5$ bytes, in the beginning, and fill it starting from the end.

*Long incomplete phrases.* A potentially incomplete phrase at the end of a block is handled the same way as in LZscan. If the incomplete phrase is short, we start the next block at the beginning of the phrase, and if it is long, we use the modified pattern matching algorithm by Crochemore to compute the full phrase. Being a constant extra space algorithm, Crochmore's algorithm works in the external memory setting as long as the full phrase fits in memory. Even longer phrases can be handled as follows. Find the occurrences of the longest prefix of the phrase that fits in memory and write the end positions of the occurrences to disk. Then read the next part of the phrase and find its occurrences that start at those end positions and so on. Crochmore's algorithm can be modified to handle this too [11]. Note that this does not change the complexity of the algorithm since each additional round of computation advances the factorization by $\Omega(b)$ steps.

| Name | $n$ | $\sigma$ | $n/z$ | Description |
|------|-----|----------|-------|-------------|
| hg | 5.85 GiB | 31 | 18.98 | $2 \times$ Human genome |
| enwik | 8 GiB | 209 | 20.45 | English Wikipedia XML |
| countries | 40.5 GiB | 203 | 3185 | Wikipedia version database |
| cere | 31 GiB | 5 | 4849 | Yeast DNA |

**Table 1.** Files used in the experiments. The value of $n/z$ (the average length of a phrase in the LZ factorization) is included as a measure of repetitiveness.

*Complexity.* The CPU time complexity of EM-LZscan is the same as LZscan $O(dnt_{\mathsf{rank}})$, where $d = O(n/M)$ and $M$ is the size of the main memory. Thus the time complexity is $O(n^2 t_{\mathsf{rank}}/M)$. The I/O complexity is dominated by the scans of A. Thus the total I/O volume is $O(n^2/M)$.

## 4   Experimental Results

We implemented two versions of the algorithm described in this paper: the first (32-bit) can parse files up to 4 GiB, the second (40-bit) is capable of handling texts up to 1 TiB. We simulate 40-bit integers as pairs of 32- and 8-bit integers. This slightly deteriorates the speed but compared to 32-bit version increases the space usage only by $b$ bytes. The implementations will be made available at `http://www.cs.helsinki.fi/group/pads/`

*Data set.* In our experiments we used the following files:
- hg: a concatenation of two different Human genomes [1,2],
- enwik: a prefix of the latest (20130403) English Wikipedia dump[3],
- countries: a concatenation of all versions (as of April 16, 2013) of Wikipedia articles about 40 large countries[4],
- cere: a concatenation of multiple copies of *cere* testfile from Pizza&Chili repetitive corpus[5], each randomly mutated with respect to original with rate 0.01%.
   Statistics about the files are summarized in Table 1.

*Setup.* We performed experiments on a set of three identical machines, each equipped with a 3.16GHz Intel Core 2 Duo CPU with 6144KiB L2 cache and 4GiB of main memory. Each machine had two 320GiB hard drives, one held the input text used during experiments and the other stored the operating system as well as all auxiliary files created by algorithms. The machines had no other significant CPU tasks running and only a single thread of execution was used. The OS was Linux (Ubuntu 12.04, 64bit) running kernel 3.2.0. All programs were

---

[1] `http://hgdownload.soe.ucsc.edu/goldenPath/hg19/chromosomes/`

[2] `ftp://public.genomics.org.cn/BGI/yanhuang/fa/`

[3] `http://dumps.wikimedia.org/enwiki/`

[4] `http://en.wikipedia.org/wiki/List_of_countries_by_GDP_(nominal)`
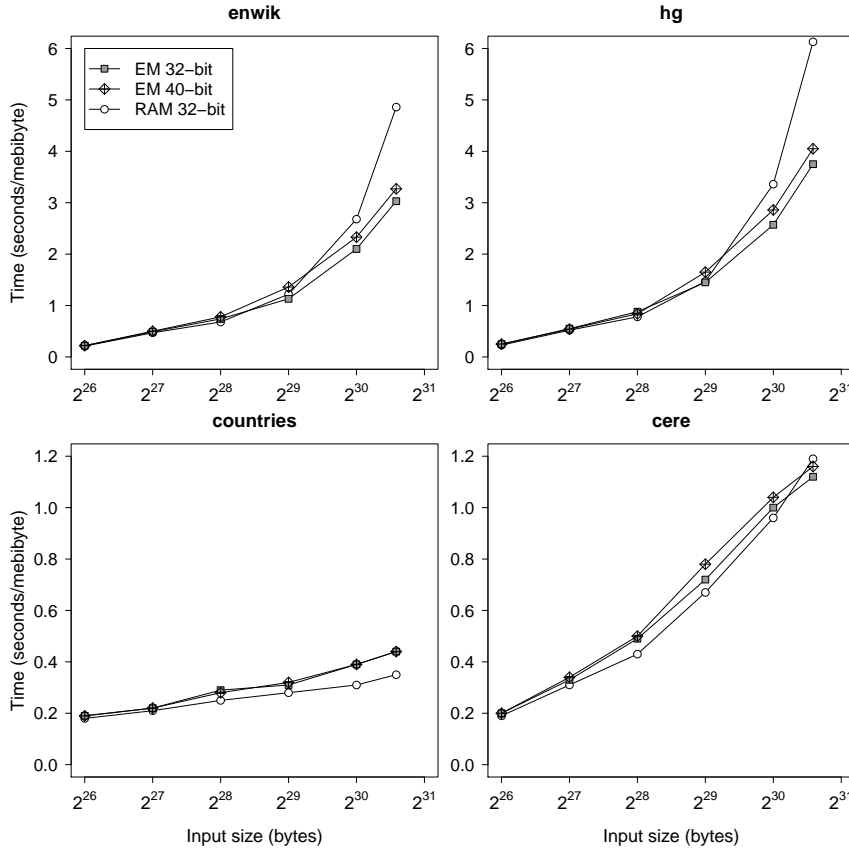
[5] `http://pizzachili.dcc.uchile.cl/repcorpus.html`

**Fig. 1.** Comparison of EM and RAM variants of LZscan.

compiled using `g++` version 4.6.4 with `-O3 -static -DNDEBUG` options. All reported runtimes are wallclock (real) times, recorded with the Unix `gettimeofday` function.

*Description.* Our first experiment compares the speed of two EM-LZscan versions to the in-memory variant of LZscan, described in [12]. The parsing was computed for increasing length prefixes of testfiles. All algorithms were limited to use 3GiB of RAM in all runs. The results are given in Fig. 1.

The second experiment measures the scalability of EM-LZscan. Similarly to previous experiment we perform the computation for various length prefixes of input files. The results are presented in Fig. 2. For each run we report the runtime (scaled to prefix length) and the I/O volume, that is, the total number of bytes transferred between RAM and disk normalized to bytes per text symbol (here also byte).
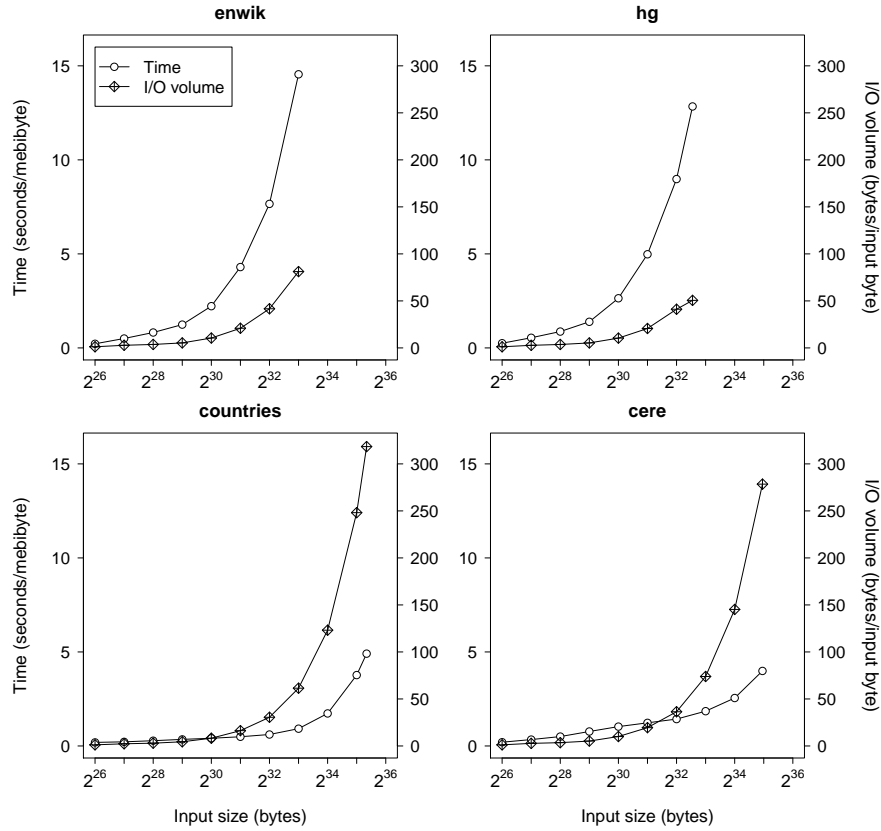
**Fig. 2.** Scalability of 40-bit EM-LZscan.

*Discussion.* The comparison of EM and RAM LZscan clearly exhibits dependence on the repetitiveness of the input. For non-repetitive files, where the runtime is dominated by the CPU computation, the EM algorithms are faster due to the use of larger blocks (no space is necessary to hold the input text, as in the RAM version). Bigger blocks imply less streaming phrases during which the matching statistics computation is performed (which is the bottleneck in the parsing of non-repetitive files – the streaming speed does not exceed 3 MiB/s).

On repetitive inputs the RAM version is slightly faster because of the streaming speed, which (mostly due to the skip trick) exceeds 200 MiB/s, whereas the disk latency is limiting the streaming speed of EM variants to ∼70MiB/s.

In all cases the 32-bit version of EM-LZscan is slightly faster than the 40-bit version. This is caused by the memory layout of two arrays holding 40-bit integers. Accessing each integer requires reaching two distant memory locations possibly attracting two cache misses, unlike the 32-bit version where at most one cache miss can occur.

As observed from Figure 2, the 40-bit EM-LZscan scales really well for highly repetitive inputs (which are the most common targets of LZ77 factorization), e.g. parsing of *countries* (40.5 GiB) file took 2.3 days on our commodity hardware. The usability of EM-LZscan on non-repetitive inputs cannot extend much beyond 2 to 3 times the size of RAM due to its quadratic complexity, but in practice non-repetitive inputs are less often treated with LZ77 because of its limited compression for such texts.

## References

1. Al-Hafeedh, A., Crochemore, M., Ilie, L., Kopylova, E., Smyth, W., Tischler, G., Yusufu, M.: A comparison of index-based Lempel-Ziv LZ77 factorization algorithms. ACM Computing Surveys 45(1), 5:1–5:17 (2012)
2. Badkobeh, G., Crochemore, M., Toopsuwan, C.: Computing the maximal-exponent repeats of an overlap-free string in linear time. In: Calderón-Benavides, L., González-Caro, C.N., Chávez, E., Ziviani, N. (eds.) SPIRE 2012. LNCS, vol. 7608, pp. 61–72. Springer (2012)
3. Barbay, J., Gagie, T., Navarro, G., Nekrich, Y.: Alphabet partitioning for compressed rank/select and applications. In: Cheong, O., Chwa, K.Y., Park, K. (eds.) ISAAC 2010. LNCS, vol. 6507, pp. 315–326. Springer (2010)
4. Bingmann, T., Fischer, J., Osipov, V.: Inducing suffix and lcp arrays in external memory. In: Sanders, P., Zeh, N. (eds.) ALENEX 2013. pp. 88–102. SIAM (2013)
5. Crochemore, M.: String-matching on ordered alphabets. Theoretical Computer Science 92, 33–47 (1992)
6. Crochemore, M., Ilie, L., Smyth, W.F.: A simple algorithm for computing the Lempel Ziv factorization. In: DCC 2008. pp. 482–488. IEEE Computer Society (2008)
7. Ferragina, P., Manzini, G.: On compressing the textual web. In: Davison, B.D., Suel, T., Craswell, N., Liu, B. (eds.) WSDM 2010. pp. 391–400. ACM (2010)
8. Gagie, T., Gawrychowski, P.: Grammar-based compression in a streaming model. In: Dediu, A.H., Fernau, H., Martín-Vide, C. (eds.) LATA 2010. LNCS, vol. 6031, pp. 273–284. Springer (2010)
9. Gagie, T., Gawrychowski, P., Kärkkäinen, J., Nekrich, Y., Puglisi, S.J.: A faster grammar-based self-index. In: Dediu, A.H., Martín-Vide, C. (eds.) LATA 2012. LNCS, vol. 7183, pp. 240–251. Springer (2012)
10. Gagie, T., Gawrychowski, P., Puglisi, S.J.: Faster approximate pattern matching in compressed repetitive texts. In: Asano, T., Nakano, S.I., Okamoto, Y., Watanabe, O. (eds.) ISAAC 2011. LNCS, vol. 7074, pp. 653–662. Springer (2011)
11. Kärkkäinen, J., Kempa, D., Puglisi, S.J.: Crochemore's string matching algorithm: simplification, extensions, applications. In: PSC 2013. To appear.
12. Kärkkäinen, J., Kempa, D., Puglisi, S.J.: Lightweight Lempel-Ziv parsing. In: Bonifaci, V., Demetrescu, C., Marchetti-Spaccamela, A. (eds.) SEA 2013. LNCS, vol. 7933, pp. 139–150. Springer (2013)
13. Kärkkäinen, J., Kempa, D., Puglisi, S.J.: Linear time Lempel-Ziv factorization: Simple, fast, small. In: Fischer, J., Sanders, P. (eds.) CPM 2013. LNCS, vol. 7922, pp. 189–200. Springer (2013)
14. Kolpakov, R., Bana, G., Kucherov, G.: mreps: efficient and flexible detection of tandem repeats in DNA. Nucleic Acids Research 31(13), 3672–3678 (2003)

15. Kolpakov, R., Kucherov, G.: Finding approximate repetitions under Hamming distance. Theoretical Computer Science 303(1), 135–156 (2003)
16. Kolpakov, R., Kucherov, G.: Finding maximal repetitions in a word in linear time. In: FOCS 1999. pp. 596–604. IEEE Computer Society (1999)
17. Kreft, S., Navarro, G.: LZ77-like compression with fast random access. In: Storer, J.A., Marcellin, M.W. (eds.) DCC. pp. 239–248. IEEE Computer Society (2010)
18. Kreft, S., Navarro, G.: Self-indexing based on LZ77. In: Giancarlo, R., Manzini, G. (eds.) CPM 2011. LNCS, vol. 6661, pp. 41–54. Springer (2011)
19. Kuruppu, S., Puglisi, S.J., Zobel, J.: Relative Lempel-Ziv compression of genomes for large-scale storage and retrieval. In: Chávez, E., Lonardi, S. (eds.) SPIRE 2010. LNCS, vol. 6393, pp. 201–206. Springer (2010)
20. Ohlebusch, E., Gog, S.: Lempel-Ziv factorization revisited. In: Giancarlo, R., Manzini, G. (eds.) CPM 2011. LNCS, vol. 6661, pp. 15–26. Springer (2011)
21. Ziv, J., Lempel, A.: A universal algorithm for sequential data compression. IEEE Transactions on Information Theory 23(3), 337–343 (1977)