

On the Reproducibility of Experiments of Indexing Repetitive Document Collections

Antonio Fariña^{*,a}, Miguel A. Martínez-Prieto^b, Francisco Claude^c, Gonzalo Navarro^d, Juan J. Lastra-Díaz^{**,e}, Nicola Prezza^{**,f}, Diego Seco^{**,g}

^aUniversity of A Coruña, Facultade de Informática, CITIC, Spain.

^bDataWeb Research, Department of Computer Science, University of Valladolid, Spain.

^cEscuela de Informática y Telecomunicaciones, Universidad Diego Portales, Chile.

^dMillennium Institute for Foundational Research on Data (IMFD), Department of Computer Science, University of Chile, Chile.

^eUniversidad Nacional de Educación a Distancia, Spain.

^fUniversity of Pisa, Italy.

^gDepartment of Computer Science, Faculty of Engineering, University of Concepción, Chile

Abstract

This work introduces a companion reproducible paper with the aim of allowing the exact replication of the methods, experiments, and results discussed in a previous work [5]. In that parent paper, we proposed many and varied techniques for compressing indexes which exploit that highly repetitive collections are formed mostly of documents that are near-copies of others. More concretely, we describe a replication framework, called `uiHRDC` (*universal indexes for Highly Repetitive Document Collections*), that allows our original experimental setup to be easily replicated using various document collections. The corresponding experimentation is carefully explained, providing precise details about the parameters that can be tuned for each indexing solution. Finally, note that we also provide `uiHRDC` as reproducibility package.

Keywords: *repetitive document collections, inverted index, self-index, reproducibility.*

1. Introduction

Scientific publications remain the most common way for disseminating scientific advances. Focusing on Computer Science, a scientific publication: *paper*, i) exposes new algorithms or techniques for addressing an open challenge, and ii) reports experimental numbers to evaluate these new approaches regarding a particular baseline. Thus, empirical evaluation is the stronger evidence of the achievements reported by a paper, and its corresponding setup is the only way that the scientific community has for assessing and (possibly) reusing such research proposals.

The ideal situation arises when a paper exposes enough information to make its research easily reproducible. According to the definition proposed by the Association for Computation Machinery (ACM) [1], an experiment is *reproducible* when an independent group of researchers can obtain the same results using artifacts which they develop independently. A weaker form of reproducibility is replicability. In this case,

*Corresponding author

**Reproducibility reviewer

Prepared for submission to *Information Systems* March 12, 2019
Email addresses: `antonio.farina@udc.es` (Antonio Fariña), `migumar2@infor.uva.es` (Miguel A. Martínez-Prieto), `fclaude@recoded.cl` (Francisco Claude), `gnavarro@dcc.uchile.cl` (Gonzalo Navarro), `jlastra@invi.uned.es` (Juan J. Lastra-Díaz), `nicola.prezza@gmail.com` (Nicola Prezza), `dseco@udec.cl` (Diego Seco)

33 the original setup is reused, so an experiment is *replicable* when an independent group of researchers can
34 obtain the same results using the author own artifacts. Finally, the ACM document [1] also defines the
35 concept of repeatability. An experiment is *repeatable* if the group of researchers is able to obtain the same
36 results reusing their original setup (including the same measurement procedures and the same system, under
37 the same operating conditions). Non-repeatable results do not provide solid insights about any scientific
38 question, so they are not suitable for publication. Thus, we can accept that research in Computer Science
39 must be, at least, repeatable.

40 Although reproducibility is the ultimate goal, getting it is not always easy because it requires that a
41 new experimental setup to be reconstructed from scratch. Implementing new algorithms, data structures,
42 or other computational artifacts can be extremely complex, and tuning them may require setting many
43 parameters which, in turn, can depend on particular system configurations or external tools. Besides, it may
44 be necessary to use programs in the exact versions used originally [26]. As a consequence, reproducibility is
45 time-consuming, error-prone, and sometimes just infeasible [4].

46 In practice, publishing replicable research¹ would be an important step forward in Computer Science.
47 According to Collberg et al. [8], replicating a computational experiment *only* needs that the source code
48 and test case data have to be available. Although simple, many papers do not meet this precondition. The
49 aforementioned paper [8] provides an interesting study involving 601 papers from conferences and journals.
50 It shows that only 32.3% of the papers provide enough information and resources to build their executables
51 in less than 30 minutes. Note that external dependencies must be accessible to compile these sources.
52 Regarding test case data, Vandewalle et al. [30] report that more papers make data available, but it is
53 mainly due to some of them use standard corpora in their corresponding experimental setups. However, the
54 problem goes beyond. Once compiled, proper parameter settings must be set, and sometimes it is tricky to
55 find this information in the papers. In conclusion, reproducibility research is currently challenging.

56 The situation is not quite different in **Information Retrieval** (IR), the field more related to the research
57 proposed in our work. IR is a broad area of Computer Science focused primarily on providing the users
58 with easy access to information of their interest [3]. As an empirical discipline, advances in IR research
59 are built on experimental validation of algorithms and techniques [20], but reproducing IR experiments is
60 extremely challenging, even when they are very well documented [13]. A recent Dagstuhl Seminar about
61 *reproducibility of data-oriented experiments in e-Science* [12] concludes that IR systems are complex and
62 depend on many external resources that cannot be encapsulated with the system. Besides, it reports that
63 many data collections are private, and their size could be an obstacle even for obtaining them. Finally, it
64 notes that some experimental assumptions are often hidden, making reproducibility a less achievable goal.

¹In this paper, we use the term reproducibility to refer experimental setups that can be replicated.

65 Following the initiative of some previous papers [33, 18], the aim of this work is to propose a detailed
66 experimental setup that replicates the methods, experiments, and results on *indexing repetitive document*
67 *collections* discussed in a previous work [5] (we will refer it as the *parent paper*). This experimental setup
68 focuses on two dimensions: i) the *space* used to preserve and manage the document collection, and the ii)
69 *time* needed to query this data.

70 The rest of the paper is organized as follows. In Section 2, we briefly describe all the inverted-index based
71 variants and the self-indexes evaluated in the parent paper. We also explain the most relevant space/time
72 tradeoffs reported from our experiments. The following sections are devoted to the reproducibility of these
73 experiments. Section 3 details the fundamentals of `uiHRDC`, our replication framework, which comprises
74 i) the *source code* of all our indexing approaches, ii) some *test data* collections, and iii) the *set of scripts*
75 that runs the main tasks of our experimental setup and generates a final report with the main figures of
76 the parent paper. We focus both on discussing the workflow that leads the process of replicating all our
77 experiments, and also in describing the structure of the `uiHRDC` framework. Section 4 describes the Docker²
78 container that allows this workflow to be easily reproduced in a tuned environment. Some brief conclusions
79 are exposed in Section 5 to motivate the need of improving research reproducibility in Computer Science.
80 Finally, Appendix A includes the actual compression ratios obtained for each technique (which complements
81 the results shown in the figures throughout the paper), and Appendix B is devoted to explain how some of
82 our self-indexes can be reused for dealing with universal (not document oriented) repetitive collections.

83 2. Universal Indexes for Highly Repetitive Document Collections

84 2.1. Background

85 Indexing highly repetitive collections has gained relevance because huge corpora of versioned documents
86 are increasingly consolidated. *Wikipedia*³ is the most recognizable example, exposing millions of articles
87 which evolve to provide the best picture of the reality around us. Each Wikipedia article comprises one version
88 per document edition, and most of them are near-duplicates of others. Apart from versioned document
89 collections, other applications perform on versioned data. For instance, the version control system *GitHub*⁴.
90 In this case, a tree of versions is maintained to control changes from millions of software projects which are
91 continuously updated by their developers. Biological databases (where many DNA or protein sequences from
92 organisms of the same or related species are maintained), or periodic technical publications (where the same
93 data, with small updates, are published over and over) are other examples of computer systems managing
94 highly repetitive data.

²<https://www.docker.com/>

³<https://www.wikipedia.org/>

⁴<https://www.github.com/>

95 The parent paper [5] focused on natural language text collections, which can be decomposed into words,
96 and queried for words or phrases. Managing these document collections is a challenge by itself due to their
97 large volume, but at the same time, they are highly compressible due to their repetitiveness. In addition,
98 version control systems require direct access to individual versions. This is also challenging because it
99 demands efficient and potentially large indexes to be deployed on top of the data collection. Thus, we need
100 to compress not only the data, but also indexes required to speed up searches.

101 The **Inverted Index** [3] has been traditionally used to index natural language text collections. The
102 inverted index maintains a list of the occurrences (also referred as *posting list* or *inverted list*) of each distinct
103 word in the text. It enables two different types of inverted indexes to be distinguished: i) *non-positional*
104 indexes store the list of document identifiers that contain each different word; and ii) *positional* indexes
105 store, in addition to the document identifiers, the word offsets of the occurrences within each document.
106 Posting lists are usually sorted by increasing document identifier, and by increasing word offset within each
107 document for positional indexes. This decision is useful for list *intersections*, a fundamental task under the
108 Google-like policy of treating bag-of-word queries as ranked AND-queries. Intersections are also relevant to
109 solve queries involving multiple words (phrase queries).

110 There is a burst of recent activity in exploiting repetitiveness at the indexing structures, in order to
111 provide fast searches in the collection within little space. These approaches are summarized in the following.
112 On the one hand, Section 2.1.1 discusses the fundamentals of inverted index compression, and summarizes
113 all our approaches. On the other hand, Section 2.1.2 introduces the concept of self-index [23], an innovative
114 succinct data structure which integrates data and index into a single compressed representation. Besides,
115 we explain how self-indexes can be adapted to perform on document versioned collections attending to our
116 original work. A detailed review of related literature, and a full explanation of our compressed inverted
117 indexes and self-indexes can be found in the parent paper.

118 2.1.1. Compressed Inverted Indexes

119 Since posting lists are sequences of increasing numbers, traditional compression techniques typically
120 compute the difference between consecutive values and then encode those *gaps* with a variable-length encoding
121 that favors small numbers. This is the basis of techniques traditionally used to represent posting lists: those
122 using **Rice** codes or **Vbyte** codes, that assign one codeword to each *gap*, or others such as **Simple9**, that
123 packs several *gaps* within a unique integer, or **PforDelta** which compresses blocks of k *gaps*. Therefore,
124 all aforementioned gap encoding techniques take advantage of the fact that the values within posting lists
125 (document identifiers or position values) are probably very close, and consequently they require few bits for
126 their encoding. We have revisited all the previous techniques and we used them in this repetitive scenario.
127 In addition, we have considered other state-of-the-art posting list representations including: **QMX**, which is a
128 good representative of the last generation of SIMD-based techniques [29], or the recent Elias-Fano technique

129 from [24] (**EF-opt**) and the implementations from [24] of some well-known representations such as **OPT-PFD**
130 [34], **Interpolative** coding [22], and **varintG8IU** [27].

131 Unfortunately, traditional techniques are not able to detect the repetitiveness that arises in versioned
132 collections. As one of the major contributions of our parent paper, we proposed some different techniques
133 which exploit the types of repetitiveness that arise in versioned collections. We applied them both for non-
134 positional and positional inverted indexes. They are based on run-length encoding (**Rice-Runs**), grammar-
135 based compression (**RePair**), and Lempel-Ziv compression (**Vbyte-LZMA** and **Vbyte-Lzend**).

136 Even though the use of compressed posting list representations permits to drastically reduce inverted
137 index size, it also has implications in query time. As we indicated in the parent paper, intersections can be
138 performed by traversing the lists sequentially in a *merge*-type fashion. However, if one of the lists to intersect
139 is much shorter than the other/s it is preferred to provide direct access (using sampling) so that the elements
140 of the smallest list can be searched for within the longer list. This type of intersection is commonly referred
141 to as *Set-vs-Set (svs)* in the literature. In practice, the best choice is to sort the lists by length. Then, we
142 take the shortest one as the “candidate” list, and it is iteratively intersected with longer and longer lists,
143 hence shortening the candidate list at each step.

144 There are two main sampling structures to provide direct access in the literature. Culpepper and Moffat
145 [9] propose the first one, referred to as **CM**. It regularly samples the compressed list and stores separately an
146 array of samples, which is searched with exponential search. Given a sampling parameter k , a list of length
147 ℓ is sampled every $k \log_2 \ell$ positions. The second method, by Transier and Sanders [28] (**ST**), applies domain
148 sampling. It regularly samples the universe of positions, so that the exponential search can be avoided. Given
149 a parameter B , it samples the universe of size u at intervals $2^{\lceil \log_2(uB/\ell) \rceil}$. To perform intersections, *lookup*
150 algorithm was defined [28]. It somehow resembles a *merge*-type algorithm but takes advantage of the domain
151 sampling. As we will see below, we combined both **CM** and **ST** sampling with posting list representations
152 based on **Vbyte** and with our **Re-Pair**-based alternatives.

153 *Posting List Representations.* We include here a brief description of the different posting list representa-
154 tions evaluated in our original paper and the tuning options they support (if any). All this information is
155 summarized in Table 1.

- 156 • **Vbyte.** We included a simple posting list representation based on **Vbyte** [31] which uses no sampling
157 and consequently performs intersections in a *merge*-wise fashion. We also included two alternatives
158 using **Vbyte** coupled with sampling [9] (called **Vbyte-CM**), with $k = \{4, 32\}$, or domain sampling [28]
159 (called **Vbyte-ST**), with $B = \{16, 128\}$. In the former case, as indicated above, intersections follow
160 a *Set-vs-Set approach* where the smallest list is decompressed and its values are searched for within
161 the other lists using exponential search. For **Vbyte-ST** we used *lookup* intersection algorithm [28]. In
162 addition, we included a hybrid variant of **Vbyte-CM** that uses bitmaps to represent the longest posting

| Method | Variants | Description |
|--|-------------------------------|---|
| Vbyte | Vbyte ^[31] | Simple Vbyte encoding with no sampling. Intersections are performed in a <i>merge-wise</i> fashion. |
| | VbyteB | Vbyte enhanced with bitmaps to represent the longest posting lists. |
| | Vbyte-CM ^[9] | Vbyte coupled with list sampling: $k = \{4, 32\}$. Intersections are performed in a <i>set-vs-set</i> approach. |
| | Vbyte-CMB ^[9] | Vbyte-CM enhanced with bitmaps to represent the longest posting lists. |
| | Vbyte-ST ^[28] | Vbyte coupled with domain sampling: $B = \{16, 128\}$ ($B = \{64\}$ is also tested for the positional scenario). Intersections follow a <i>lookup</i> approach. |
| | Vbyte-STB | Vbyte-ST enhanced with bitmaps to represent the longest posting lists. |
| Rice | Rice ^[32] | Simple Rice encoding with no sampling. Intersections are performed using a <i>merge</i> algorithm. |
| | RiceB | Rice enhanced with bitmaps to represent the longest posting lists. |
| Simple9 ^[2] | No variants | Word-aligned Simple9 that packs consecutive <i>gaps</i> into a 32-bit words. Intersections are performed using a <i>merge</i> algorithm. |
| PforDelta ^[15, 35] | No variants | Extends Simple9 <i>gaps</i> to pack many <i>gaps</i> together (up to 128). A variant of Simple9 is used to encode exceptions. Intersections are performed in a <i>merge-wise</i> fashion. |
| QMX ^[29, 19] | No variants | Exploits SIMD-instructions to boost the decoding of long lists. Intersections are performed using an specific algorithm that also benefits from such instructions. |
| Rice-Runs ^[32] | No variants | Rice coupled with <i>run-length</i> encoding. Intersections are performed in a <i>merge-wise</i> fashion. |
| Vbyte-LZMA ^[31] | No variants | Encodes <i>gaps</i> with Vbyte and, if the size of the resulting Vbyte-sequence is ≥ 10 bytes, then it is further compressed with LZMA. A bitmap indicates which posting lists are compressed with LZMA. Intersections are performed using <i>merge</i> algorithm. |
| Vbyte-Lzend ^[31, 16] | No variants | Encodes <i>gaps</i> with Vbyte and then all posting lists are compressed as a whole with <i>LZ-End</i> . It can be parameterized: $ds = \{4, 16, 64, 256\}$. Intersections first obtain Vbyte-encoded sequences and then perform in a <i>merge-wise</i> fashion. |
| RePair-based ^[17] | RePair | Compresses posting lists as a whole using RePair. Intersections are performed using <i>merge</i> algorithm over the compressed lists. |
| | RePair-Skip | Adds <i>skipping</i> data to RePair to improve performance. |
| | RePair-Skip-CM | RePair-Skip coupled with CM-type (list) sampling: $k = \{1, 64\}$. |
| | RePair-Skip-ST | RePair-Skip coupled with ST-type (domain) sampling: $B = \{1024\}$ for the non-positional scenario; $B = \{256\}$ for the positional scenario. |
| Indexes from Ottaviano and Venturini's Framework ^[24] | EF-opt ^[24] | Partitioned Elias-Fano index. |
| | OPT-PFD ^[34] | Optimized PforDelta variant. |
| | Interpolative ^[22] | Binary interpolative coding. |
| | varintG8IU ^[27] | SIMD-optimized Variable Byte code. |

Table 1: Summary of posting list representations evaluated in the parent paper.

163 lists (**Vbyte-CMB**) [9]. In practice, lists longer than $u/lenBitmapDiv$ (where u is the largest document
164 identifier, and $lenBitmapDiv = 8$) are replaced by a bitmap [9] that marks which documents are
165 present in the list. For completeness, we used the hybrid approach over **Vbyte-ST** to build **Vbyte-STB**,
166 and also included a variant **VbyteB** with no sampling.

167 • **Rice**. We included a representation based on Rice codes [32] and also implemented for completeness a
168 **RiceB** variant using bitmaps for the longest lists. In both cases they use no sampling and intersections
169 are done using *merge* algorithm.

170 • **Simple9**. We included a representation based on the word-aligned **Simple9** technique, by Anh and
171 Moffat [2]. It packs consecutive *gaps* into a 32-bit word. It uses the first 4 bits to signal the type of
172 packing done, depending on how many bits the next *gaps* need: we can pack 28 1-bit numbers, or 14
173 2-bit numbers, and so on. Intersections are done using *merge* algorithm.

174 • **PforDelta**. This representation [15, 35] uses the same idea of packing many *gaps* together, typically
175 up to 128 (parameter $pdfThreshold = 128$), while allowing for 10% of exceptions that need more bits
176 than the core 90% of the *gaps*. Those exceptions are then coded separately using a variant of **Simple9**.
177 In our case, we obtained our best results with $pdfThreshold = 100$. Again, we performed intersections
178 *merge-wise*.

179 • **QMX**. This technique [29] uses SIMD-instructions to boost decoding of large lists,⁵ and was coupled with
180 an intersection algorithm [19] that also benefits from SIMD-instructions.⁶

181 • **Rice-Runs**. We coupled *run-length* encoding with Rice coding to boost both the compression and
182 intersection speed of our **Rice** representation in a repetitive scenario where large sequences of +1 *gaps*
183 could occur. Basically, a sequence of +1 *gaps* of length k is represented as the *ricecode*(+1) followed
184 by *ricecode*(k). Intersections are performed *merge-wise*, yet they take advantage of the fact that a run
185 of +1 values can be skipped/decoded in a unique operation.

186 • **Vbyte-LZMA**. In this representation we independently compress each posting list with a **Vbyte+LZMA**
187 chain. We initially encode the *gaps* in the posting list with **Vbyte** and then compress the output with
188 the **LZMA** variant of LZ77 (www.7-zip.org). We use a threshold parameter (*minbcssize*) so that
189 LZMA is only applied on those lists where their **Vbyte** encoded sequence occupies at least 10 bytes
190 ($minbcssize \leftarrow 10$). Otherwise, the initial **Vbyte** encoded sequence is stored. A bitmap indicates
191 which posting lists were stored compressed with **Vbyte** plus LZMA, and which ones only with **Vbyte**.
192 Since **Vbyte-LZMA** only supports extracting a list from the beginning, intersections can involve a initial

⁵<http://www.cs.otago.ac.nz/homepages/andrew/papers/QMX.zip>.

⁶<https://github.com/lemire/SIMDCompressionAndIntersection>.

193 step that applies LZMA decoder to recover the **Vbyte** encoded sequences, and then we apply *merge*
194 algorithm on those **Vbyte** encoded sequences.

195 • **Vbyte-Lzend**. This representation follows the same idea of **Vbyte-LZMA** but uses LZ-End [16] instead
196 of LZMA to compress the posting lists. Since LZ-End allows random access, the sequence of all the
197 concatenated lists can be compressed as a whole, not list-wise. Therefore, we first concatenate the
198 **Vbyte** representation of all the posting lists (keeping track of the offset where each posting list started
199 in the **Vbyte** encoded sequence), and then use LZ-End to represent it. At construction time, we can
200 obtain different space/time tradeoffs by tuning the *delta-codes-sampling* parameter (ds) of LZ-End (see
201 [16] for details). In our experiments we set it to $ds = \{4, 16, 64, 256\}$. At intersection time, we first
202 recover the **Vbyte** encoded sequences and then proceed *merge*-wise as in **Vbyte-LZMA**.

203 • **RePair**-based representations. As in the LZ-End-based representation we compress all the posting lists
204 as a whole, yet we directly work on the sequences of *gaps* rather than on their **Vbyte** encodings. We
205 use the grammar-based compressor Re-Pair that recursively replaces the most frequent pair of symbols
206 (initially called *terminals*, which are those *gaps* from the concatenation of all the posting lists) by a new
207 *non-terminal* symbol not occurring before in the original sequence. Re-pairing process ends when no
208 pair occurs at least twice. The result is a sequence that could contain both terminal and non-terminal
209 symbols, and a dictionary of substitution rules which is represented in a compact format. This makes
210 up our basic **RePair** representation. It allows us to extract individual posting lists,⁷ and to perform
211 intersections in a *merge*-type fashion over the compressed representation of the lists (non-terminals are
212 expanded/decoded recursively during the traversal of the lists). In addition, in our implementation
213 we added a parameter *repairBreak* that allows us to stop the recursive pairing when the gain in
214 compression ratio (reduction of size of the compressed sequence) of the current step with respect to
215 the previous step is smaller than *repairBreak*.⁸ In practice, for **RePair** we set *repairBreak* = 0.0 so
216 that the Re-pairing process is not broken at all.

217 The compact representation of the dictionary of rules allows us to add additional *skipping* data (i.e.
218 the sum of all the *gaps* included below that non-terminal symbol) with little space cost. This skipping
219 data allows us to improve both decompression time for a list, and intersection time. We adapted the
220 *merge*-based algorithm from **RePair** so that it is boosted by using this skipping data. The resulting
221 algorithm was named *skip*. This makes a new representation of posting lists referred to as **RePair-Skip**.

⁷To avoid creating Re-Pair phrases that span along two different posting lists, we add a non-repeating separator to mark the beginning of each posting list.

⁸Being n and m respectively the length of the original sequence and the length of the compressed sequence, we initially set $prevRatio \leftarrow (100.0 * m/n)$. Then, after each substitution, if $((prevRatio - (100.0 * m/n)) < repairBreak)$ we break the Re-pairing process. Otherwise, we simply update $prevRatio \leftarrow (100.0 * m/n)$.

222 In our experiments, we tuned `RePair-Skip` with $repairBreak = 4 \times 10^{-7}$ and $repairBreak = 5 \times 10^{-7}$
223 respectively for the non-positional and positional scenarios.

224 Initially, neither `RePair` nor `RePair-Skip` had sampling to speed up intersections, and only `RePair-Skip`
225 could benefit of the additional data to boost the intersections. However, since `RePair-Skip` was shown
226 to outperform `RePair`, we also created two `RePair-Skip` variants using both `CM`-type and `ST`-type
227 sampling. They are named `RePair-Skip-CM` (where we set $k = \{1, 64\}$) and `RePair-Skip-ST` (with
228 $B = \{256, 1024\}$).

229 • **Ottaviano&Venturini’s Partitioned Elias-Fano indexes.** In [24], Ottaviano and Venturini pre-
230 sented several posting list representation alternatives based on Elias-Fano codes. We used the best of
231 them (`EF-opt`). The source code is available at authors’ website.⁹ In addition, they included in their
232 framework implementations of some well-known representations such as `OPT-PFD` (optimized `PforDelta`
233 variant [34]), `Interpolative` (Binary Interpolative Coding [22]), and `varintG8IU` (SIMD-optimized
234 Variable Byte code [27]). All of them were compared in our parent paper.

235 *Our implementation of non-positional and positional inverted indexes.* Given a text T that is composed of
236 a set of documents, to create our inverted indexes, we process T to gather the different words/terms in T ,
237 and for each term we obtain the positions where they occur. The text itself can be processed with a suitable
238 compressor in order to reduce space needs. In our implementation, we used `Re-Pair`¹⁰ as text compressor,
239 which reached compression ratios below 2% (we kept the rules uncompressed -as a pair of integers- to speed
240 up decompression), and we added sampling for every $\{1, 2, 8, 32, 64, 256, 4096\}$ symbols of the final `Re-Pair`
241 sequence to allow decompressing arbitrary parts of the original text.¹¹ Yet, the main focus of our parent
242 paper was set on how to efficiently deal with the set of posting lists. We included all the above posting
243 lists representations, with the exception of Ottaviano&Venturini’s implementations, within a package that
244 provides us methods to build a compressed representation for a set of posting lists, extract a list, intersect
245 2- n lists, show the compressed size, etc. We did this to obtain compressed representations for posting lists
246 to be used in a non-positional scenario, where we indexed document-ids and intersections were used to
247 solve AND-conjunctive queries for the elements of our query patterns, but also for a positional scenario,
248 where those patterns were taken as phrase queries. In this positional scenario, we indexed actual word/term
249 positions and solving phrase queries implied that the intersections had to consider the order of the terms in
250 the query patterns to return the documents where they occur at consecutive positions. The building process
251 is handled by a **build** program that processes the source text and saves all the required structures to disk.

⁹https://github.com/ot/partitioned_elias_fano.

¹⁰We have also the option to use: a LZ-End compressor; keeping the text in plain form; or even not storing the text at all.

¹¹This was of interest in the positional scenario where, as we will see in the next section, we are interested in comparing the snippet extraction time with that of self-indexes.

252 A **searcher** program is on charge of loading those structures into main memory and then allowing us to
253 perform queries. Both programs are linked a posting list representation and a text compressor.

254 According to the experiments reported in our parent paper, our compressed inverted index variants allow
255 us to perform two types of queries:

256 • **locate**(*p*) returns the positions of all the occurrences of *p* in *T*. In the non-positional scenario our
257 experiments (see Section 2.2) were focused on comparing both the fetch time and the intersection time
258 of the different posting lists representations, so basically **locate**(*p*) returns the documents where *p*
259 occurs.

260 In the positional scenario, **locate**(*p*) returns a pair of the form (*doc, pos_in_doc*). Note that our posting
261 list representation indexes *absolute* positions in *T* rather than pairs (*doc, pos_in_doc*). Therefore, to sup-
262 port returning relative positions within each document, we enhanced our inverted index with an array of
263 offsets that mark the beginning of each document in *T*, and provided an operation **merge-occs-to-docs**
264 which efficiently maps a sequence of absolute positions into pairs (*doc, pos_in_doc*). Therefore, to solve
265 **locate**(*p*) we initially perform an intersection to gather the positions in *T* where *p* occurs, and then
266 perform **merge-occs-to-docs** operation to obtain the final output.

267 • **extract**(*a, b*) retrieves the text fragment in $T[a, b]$. This is efficiently supported by using the regular
268 sampling added to the (Re-Pair) text representation.

269 2.1.2. Self-Indexes

270 A *self-index* [23] is a compact data structure that enables efficient searches over an string collection (called
271 the text, *T*), and also replaces the text by supporting extraction of arbitrary snippets. Thus, self-indexes
272 provide, at least, two basic queries:

273 • **locate**(*p*) returns the positions of all the occurrences of *p* in *T*.

274 • **extract**(*a, b*) retrieves the text fragment in $T[a, b]$.

275 This functionality allows self-indexes to be considered as an alternative choice to positional inverted
276 indexes.

277 Our original setup comprised three families of self-indexes that were able to capture high repetitiveness.
278 All of them have a representative proposal which regards *T* as a sequence of characters, but two word-oriented
279 approaches (which process *T* as a sequence of words) were also proposed. It is worth noting that, in all cases,
280 **locate**(*p*) returns *absolute* positions in *T* that must be then converted into (document,offset) pairs. The
281 aforementioned **merge-occs-to-docs** position-document mapping is used again for such purpose.

282 *CSA-based Self-Indexes*. The Compressed Suffix Array (CSA) by Sadakane [25] is one of the pioneering
 283 self-indexes and proposes a succinct suffix array (SA) encoding. It retains the original locate functionality
 284 of the suffix array (based on binary search), while providing snippet extraction in compressed space. CSA
 285 encompasses two main structures: a bitmap B , which marks where the first symbol of the suffixes changes
 286 in SA, and Ψ , an array which stores the position in SA pointing to the next character of a suffix. Ψ is
 287 highly compressible, but it is also massively used to decode a portion of the text. Thus, self-indexes dealing
 288 with highly repetitive collections must balance Ψ compression and decoding efficiency to be competitive with
 289 respect to compressed positional inverted indexes.

290 Two different CSA-based approaches were evaluated in our benchmark: RLCSA and WCSA.

- 291 • **RLCSA**. The Run-Length Compressed Suffix Array, by Mäkinen et al. [21], exploits that Ψ contains
 292 long runs of successive values in highly repetitive collections. It performs run-length encoding of
 293 $\Psi(i) - \Psi(i - 1)$ and stores regular samples to allow efficient access to absolute Ψ values. This *sample*
 294 value must be provided to build RLCSA, and yields different space/time tradeoffs: larger *sample* values
 295 are used to reduce space requirements, but it penalizes data access speed; on the contrary, small *sample*
 296 values increase efficiency at the price of less compression. The authors proposed *sample*= 512, but our
 297 experiments consider 7 different values of the form 2^i , from $i = 5$ (*sample* = 32) to $i = 11$ (*sample*
 298 = 2048), to analyze particular tradeoffs. A *blocksize* parameter is also required to configure internal
 299 bit vectors blocks (it sets the number of reserved bytes per block). We use *blocksize* = 32, as suggested
 300 by the authors.

- 301 • **WCSA**. The Word Compressed Suffix Array, by Fariña et al. [11] adapts CSA to cope with particu-
 302 lar features of natural language; i.e. it processes the input text as a sequence of words instead of
 303 characters. WCSA transforms the original text into an integer sequence where each position refers to
 304 a word/separator, but it does not provide any particular optimization to manage highly-repetitive
 305 texts. We included it in our original experiments because it acts as a bridge between inverted in-
 306 dexes and more sophisticated self-indexes, both in space and time complexity. WCSA builds B and
 307 Ψ , but over the integer sequence, and provides word-based location and extraction. This requires
 308 keeping samples of SA and samples of the inverse of SA (indicating which position of SA points
 309 to the j -th word) at regular intervals. The non-sampled values can be recovered with Ψ . There-
 310 fore, three different parameters must be provided at construction time: $\langle sPsi, sA, sAinv \rangle$. Note
 311 that the inverse of SA is only needed for extract, SA is used both for extract and locate, and Ψ
 312 is used in all search operations. Therefore, even though we can tune different setups yielding similar
 313 space requirements, it is worth to keep a rather small value of $sPsi$, a larger value of sA , and an
 314 even larger value of $sAinv$. We evaluated seven different configurations for the sampling parameters
 315 ranging from $\langle sA, sAinv, sPsi \rangle = \langle 8, 8, 8 \rangle$ to $\langle 2048, 2048, 2048 \rangle$. In particular, we used the values:

316 $\{\langle 8, 8, 8 \rangle, \langle 16, 64, 16 \rangle, \langle 32, 64, 32 \rangle, \langle 64, 128, 32 \rangle, \langle 128, 256, 128 \rangle, \langle 512, 512, 512 \rangle, \langle 2048, 2048, 2048 \rangle\}$.

317 *SLP-based Self-Indexes.* We previously motivated that grammar-based compression is a good choice for
318 posting list encoding, but it is also promissory for self-indexes. Our original paper focused on a particular
319 type of grammar compressor built around the notion of straight-line program (SLP).¹² We explored two
320 SLP-based self-indexes, referred to as **SLP** and **WSP**, which were carefully tuned for our experiments.

321 • **SLP.** Claude et al. [6] proposed originally a character oriented SLP self-index to manage (highly
322 repetitive) biological databases, and then optimized it to cope with natural language collections [7].
323 SLP indexes the set of rules as a labeled binary relation, and the reduced sequence (obtained by RePair)
324 using a varied configuration of bit-based structures. SLP requires space proportional to that of a Re-
325 Pair compression of the text, but its performance is less competitive than that reported by other
326 self-indexes. The original SLP implementation was improved in our parent paper, where we tuned
327 some of its algorithms and internal data structures. A single q parameter is required to build an
328 SLP index; it sets the lengths of the q -grams that are indexed, by an internal dictionary, to improve
329 prefix/suffix searches. q is set by default to 4 characters because no relevant improvements have been
330 found for larger values.

331 • **WSP.** The word-oriented SLP was proposed in our parent paper, and adapted SLP to perform on a
332 sequence of integers (word identifiers). Its internal configuration is similar to that of the SLP, but it
333 does not use the q -gram dictionary because it was not competitive for words. Thus, WSP exposes a
334 simple build interface which does not require any parameter.

335 *LZ-based Self-Indexes.* Finally, we evaluated two different approaches of self-indexing based on two variants
336 of LZ77-like parsing: **LZ77-index** and **LZend-index**. In short, an LZ77-like parsing transforms a text into
337 a sequence of phrases, each one encoding the first occurrence of a substring. Each substring concatenates a
338 maximal substring previously used in the text and a trailing character. **LZ77-index** and **LZend-index** uses
339 a similar configuration of succinct data structures to encode their corresponding structure of phrases, and
340 to allow fast search and decode capabilities.

341 Both self-indexes were originally proposed by Kreft and Navarro [16], but we tuned them to improve their
342 original space/time tradeoffs. Besides, it is worth noting that they support five different configurations.¹³ The
343 *default* configuration (**Conf.#1**) is considered if no parameters are provided. It reports the best performance
344 but also the worst space. The following parameter configurations are also allowed: ‘‘**bsst ssst**’’ (for

¹²SLPs are grammars in which a rule X_i generates i) a terminal j , or ii) a pair of non terminals $X_l X_r$.

¹³These configurations are described in our original paper and basically consider different structures and search algorithms to provide **locate** and **search** functionality.

345 `Conf.#2`), ‘`brev`’ (for `Conf.#3`), ‘`bsst brev ssst`’ (for `Conf.#2`), and ‘`bsst brev`’ (for `Conf.#5`).
346 The latter reports larger space savings at the price of a bit slowdown, but ensuring a competitive performance.
347 We chose it in our benchmark.

348 *2.2. Experimental Results*

349 We experimentally studied the space/time tradeoffs obtained when performing `locate` operation with the
350 described posting list representations, in both the non-positional and positional scenarios. In the positional
351 scenario we also added a comparison with the proposed self-indexes, and we finally included results regarding
352 the time needed to `extract` snippets.

353 In this section we provide a summary of the results that can be obtained using our replication framework
354 `uiHRDC` (see Section 4). In particular, for `locate(p)` operation we include results for the case where the
355 patterns are 2-word phrases, and for `extract(range)` we only include results for snippets of around 13,000
356 characters. In Section 3.1 we discuss that there are actually four types of patterns for `locate` and two
357 different snippet lengths for `extract` operation. In that section we also discuss details of the text collections
358 used. In addition, we have also included some fixes to errors that were detected in the parent paper during
359 the reviewing process.

360 The time measures included here are referred to CPU user-times and were obtained using our Docker
361 instance in an Intel(R) i7-8700K@3.70GHz CPU with 64 GB of DDR4@2400MHz memory.

362 *2.2.1. Fixes to the Parent Paper*

363 Before delving into the details of the experiments described in the parent paper, it is worth noting that
364 developing `uiHRDC` helped us to detect some minor errors in the results reported in the parent paper. More
365 precisely:

- 366 • `QMX` performance was wrong in the experiment dealing with `locate` operation for 5-words phrase
367 patterns (Figure 3 of the parent paper) in the non-positional scenario. Its actual performance is
368 slightly faster than the one reported by `VbyteB`, but it means that it is half an order of magnitude
369 slower than the fastest choice: `Vbyte-STB`. The `locate` times of `QMX` were also wrong in the case of
370 word queries in the positional scenario (Figure 6 of the parent paper). Although it remains as the
371 fastest choice in both cases, its difference to `Simple9` or `Vbyte` is drastically reduced.
- 372 • The `OPT-PFD` compression ratios are wrong in the experiments for the positional scenario (Figures
373 6 and 9). Its actual compression ratio is 29.742% (instead of 31.848%), so it is more effective than
374 `Interpolative` and `EF-opt`.

375 *2.2.2. Non-Positional scenario: Inverted Indexes*

376 For the non-positional scenario we include a comparison of all the compressed inverted index represen-
 377 tations discussed in Section 2.1.1 when considering `locate` operation. Yet, in this case we only consider
 378 fetch/intersection time and skip the parsing time of the query patterns. That is, we assume all the patterns
 379 have been mapped into *ids* before timing starts.

380 Figure 1 shows the results. In the left part, we show the results for the state-of-the-art techniques. In
 381 the right part we also include our proposals. Recall that for the variants from Ottaviano&Venturini’s
 382 framework [24] we used exactly their source code and simply adapted the format of our data and patterns
 383 to run their `build` and `search` programs. Those techniques are marked with an ‘*’ in the figures.

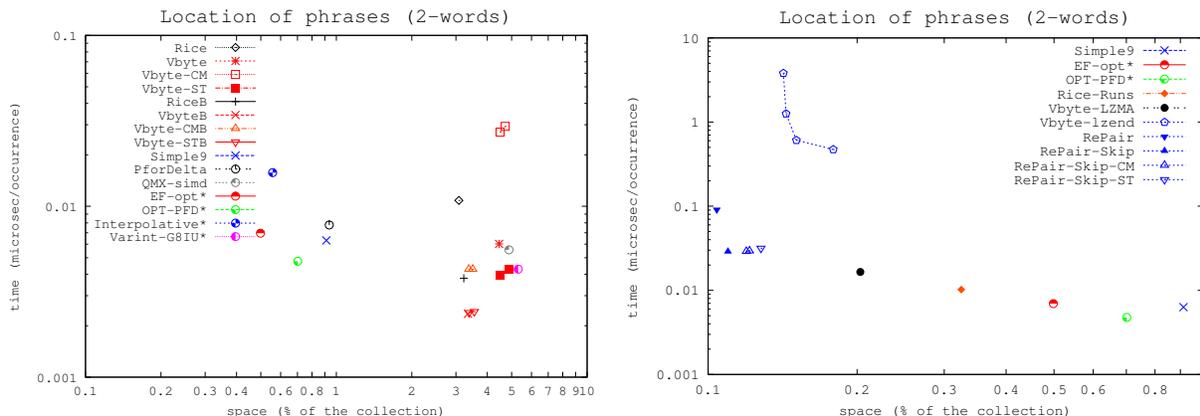


Figure 1: Main results for the non-positional scenario.

384 As shown in the parent paper, our new techniques Rice-Runs Vbyte-LZMA, Vbyte-Lzend, and Re-Pair-
 385 based variants drastically reduced the space needs of the existing techniques. Yet, they were typically also
 386 much slower (particularly in the case of Vbyte-Lzend).

387 *2.2.3. Positional scenario: Inverted Indexes and Self-Indexes*

388 In this scenario we compare both our inverted index representations and self-indexes at both `locate` and
 389 snippet `extract` operations.

390 *Pattern Location.* We did not include all the inverted index variants used in the non-positional scenario, but
 391 only those that could be of interest here. For the self-indexes, timing includes the time needed to perform
 392 (`locate(p)`) operation to obtain the positions where (*p*) occurs, and then converting absolute positions to
 393 (*document, offset*) pairs using `merge-occs-to-docs` operation. For the inverted indexes, timing includes:
 394 (a) the parsing time required to map each pattern into a sequence of *ids*; (b) performing intersection of the
 395 required posting lists (for those *ids*); and (c) running `merge-occs-to-docs` to obtain the final result. For the
 396 techniques from Ottaviano&Venturini’s framework we directly used their sources *with no modifications* to

397 measure (b) times, whereas time measures for stages (a) and (c) where done separately with a program we
 398 also implemented. During step (a), our program not only measures times, but also outputs the sequence of
 399 *ids* obtained after the parsing of each pattern *p*. These *id*-based patterns are used to perform intersections
 400 with the techniques from Ottaviano&Venturini's framework.

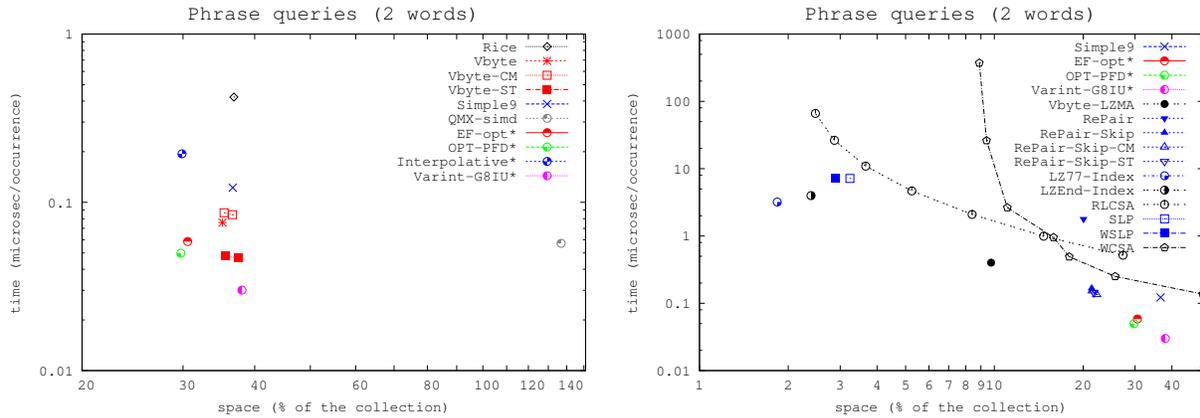


Figure 2: Main results for the positional scenario.

401 Figure 2 shows the results. On the one hand, our Vbyte-LZMA and Re-Pair-based inverted indexes still
 402 improve the space needs of the other representations, while being typically up to one order of magnitude
 403 slower. On the other hand, self-indexes showed to obtain around one order of magnitude reduction on space,
 404 while becoming up to three orders of magnitude slower than the fastest inverted index alternative.

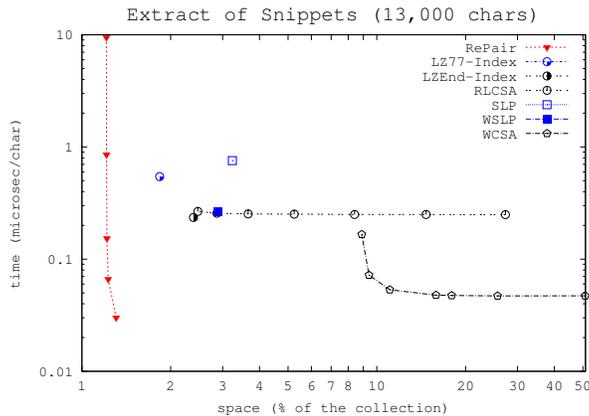


Figure 3: Main results regarding text extraction performance for both self-indexes and text compressed with Re-Pair.

405 *Snippet Extraction.* We report the time needed to extract text snippets $T[s, e]$ from the self-indexes and
 406 compare them with the decoding performance of Re-Pair, that was the text compressor chosen to compress
 407 the document collection for the representations using inverted indexes. Recall that Re-Pair is coupled with

408 the additional sampling that permit partial decompression, as discussed in Section 2.1.1. Results are shown
 409 in Figure 3. We can see that Re-Pair is very fast at decompression (when a dense sampling is used) and
 410 requires less than 1.5% of the size of the document collection. With very similar space needs (around 2–3%)
 411 we find most of the self-indexes (LZend-index, LZ77-index, RLCSA SLP, and WSLP), yet they are typically
 412 much slower. Finally, WCSA competes in speed with Re-Pair, but requires around one order of magnitude
 413 more space.

414 3. The uiHRDC Framework

415 Our experimental framework was named uiHRDC (*universal indexes for Highly Repetitive Document Col-*
 416 *lections*). It is licensed under the GNU Lesser General Public License v2.1 (GNU LGPL), is hosted at a
 417 GitHub repository,¹⁴ and it is also published through Mendeley Data [10]. It includes all the required ele-
 418 ments to reproduce the main experiments in our original paper including datasets, query patterns, source
 419 code, and scripts.

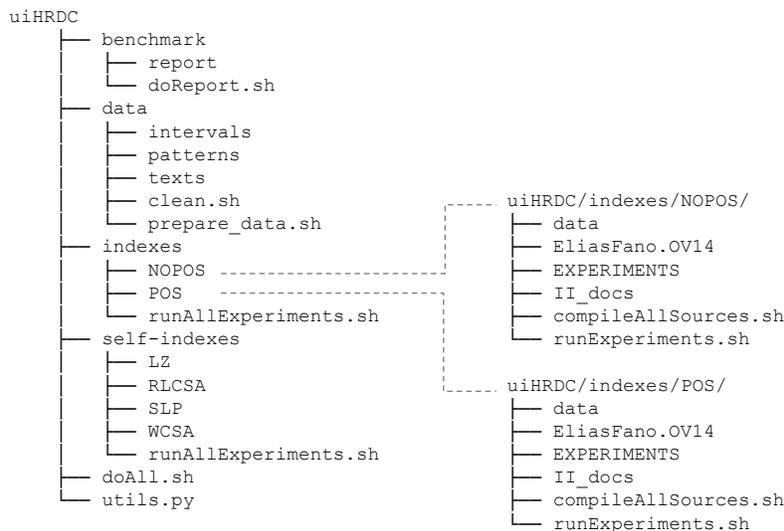


Figure 4: Structure of the uiHRDC repository.

420 The general structure of the uiHRDC repository is shown in Figure 4. It can be seen that, under the
 421 root of the repository, we include: *i*) directory **benchmark** which includes a L^AT_EX formatted report and a
 422 script that will collect all the data files resulting from running all the experiments and will generate a PDF
 423 report with all the relevant figures (including those in Section 2.2). Further python scripts required for such
 424 task are also included in directory **utils.py**. *ii*) Directory **data**, which includes the text collections (7z

¹⁴<https://github.com/migumar2/uiHRDC/>

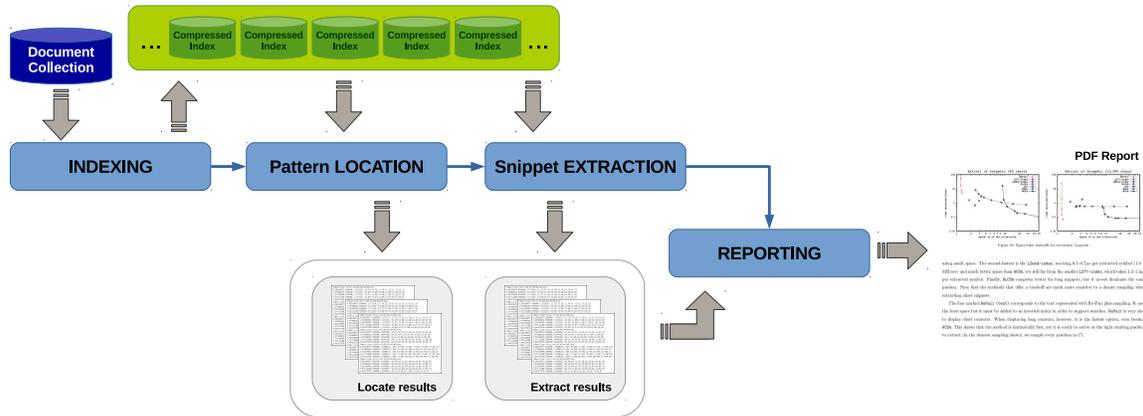


Figure 5: Workflow used in uiHRDC to reproduce our experiments.

425 compressed), and the query patterns. *iii*) Directories `indexes` and `self-indexes` that contain the source
 426 code for each indexing alternative, and scripts that permit to run all the experiments for each technique. This
 427 includes the construction of each compressed index of interest (using a `builder` program) and then performing
 428 both `locate` and `extract` operations over that index (using the corresponding `searcher` program). Each
 429 experiment will output relevant data to a results-data file. And *iv*) a script `doAll.sh` that will drive all
 430 the process of decompressing the source collections; compiling the sources for each index and running the
 431 experiments with it; and finally, generating the final report. The overall workflow followed is depicted in
 432 Figure 5.

433 In the following section we will include further details regarding the contents of our repository.

434 3.1. Test Data

435 Within uiHRDC we provide in `data` directory, both the document collections and the query sets used in
 436 the experimental setup of our parent paper. They are described below.

437 3.1.1. Document Collections

438 Our document collections were created from the 108.5 GB Wikipedia collection described by He et al. [14],
 439 which contained 10% of the complete English Wikipedia from 2001 to mid 2008. It contains 240,179 articles,
 440 and each of them has a number of versions. Its statistics are shown in Table 2. Note that we did not have
 441 the original 108.5 GB collection, but a filtered (tag-free) version of it whose size totaled 85.58 GB. Therefore,
 442 the original collection was 1.27 times larger than ours.

443 From the `Full` collection of articles, we chose two subsets of the articles, and collected all the versions
 444 of the chosen articles. For the non-positional setting our subset (`Non-pos`) contains a prefix of 19.53 GB of
 445 the full collection, whereas for positional indexes we chose 1.94 GB of random articles. However, for a fair

| Subset | Size (GB) | Articles | Number of versions | Versions / article | Filename (within uiHRDC) | Filesize (GB) |
|---------|--------------|----------|-----------------------|-----------------------|-----------------------------|------------------|
| Full | 108.50 | 240,179 | 8,467,927 | 35.26 | -- | 85.58 |
| Non-pos | 24.77 | 2,203 | 881,802 | 400.27 | text20gb.txt | 19.53 |
| Pos | 1.94 | 4,327 | 149,761 | 34.61 | wiki_src2gb.txt | 1.94 |

Table 2: Detailed statistics of the document collections used.

446 comparison with the techniques from [14], in the non-positional scenario we scaled the size of the `Non-pos`
447 subset using the above 1.27 factor when referring to its space requirements. Additional statistics of our two
448 subsets are also included in Table 2.

449 3.1.2. Query sets

450 Since the experiments target at providing space/time for the different indexing alternatives when per-
451 forming `locate(pattern)` and `extract(interval)` queries we provide two types of query sets for each
452 document subcollection.

- 453 • Query sets for `locate`: We provide four query sets, each of them containing 1,000 queries, which
454 include: *i*) two query sets composed of one-word patterns chosen at random from the vocabulary of
455 the indexed subcollection. In the first case (W_a), it includes low-frequency words occurring less than
456 1,000 times. In the second case, the query set (W_b) includes high-frequency words occurring more than
457 1,000 times; *and ii*) two query sets with 1,000 phrases composed of 2 and 5 words that were chosen
458 randomly from the text of the subcollection (with no restrictions on its frequency).
- 459 • Interval sets for `extract`: Aiming at measuring extraction time when recovering snippets of length 80
460 (around one line) and 13,000 (around one document, in our collection) characters, we generated: *i*) a
461 set of 10,000 intervals of width 13,000 characters from the POS text collection, and *ii*) a set containing
462 100,000 intervals of width 80 characters. Since these intervals are not suitable for our word-based
463 self-indexes (WCSA and WSLP), and assuming that the average word length is around 4 in our datasets,
464 we also generated two additional sets composed respectively of 10,000 intervals containing 3,000 words
465 each, and 100,000 intervals containing 20 words each.

466 3.2. Source Code and scripts

467 As indicated above, the source code provided in our uiHRDC repository has two main directories `indexes`
468 and `self-indexes`. Those directories include both the source code and the scripts required to reproduce our
469 experiments. In this section we include more details regarding their structure so that an interested reader
470 can rapidly understand how they are organized.

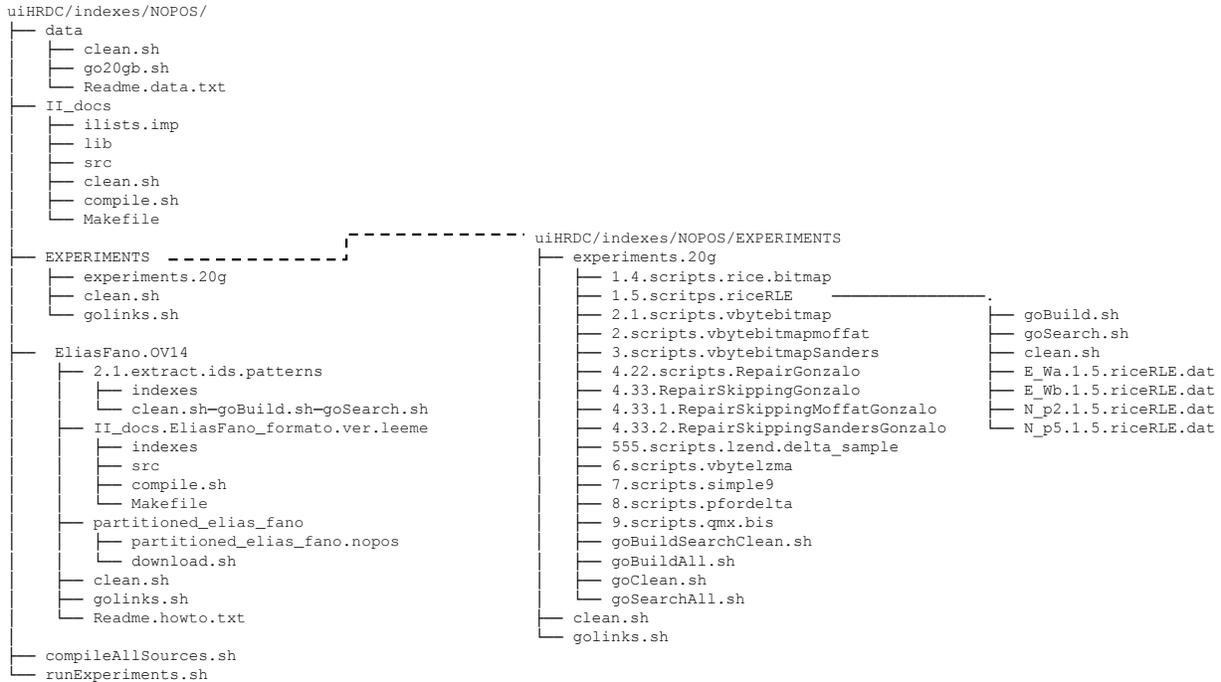


Figure 6: Structure of the `uiHRDC/indexes` directory in the `uiHRDC` repository.

471 3.2.1. Indexes

472 Under `uiHRDC/indexes` directory, as it is shown in Figure 4, we can find a script `runAllExperiments.sh`
473 and two directories `NOPOS` and `POS`. Basically, that script enters both directories, and then runs two scripts:
474 one to compile the source codes (`compileAllSources.sh`) and another one to launch the experiments
475 (`runExperiments.sh`) in that directory. An interested reader should probably start by opening those small
476 scripts. These scripts are included in Figure 6, where we show the structure of directory `NOPOS` (the structure
477 of `POS` is almost identical).

478 *Directory `uiHRDC/indexes/NO-POS: Non-Positional Inverted Indexes`.* Under this directory, we can find the
479 scripts `compileAllSources.sh` and `runExperiments.sh` discussed above, a `data` directory were links to
480 `uiHRDC/data` will be created by the scripts under this directory, and finally two main parts: *i)* `EliasFano.OV14`
481 with the sources and scripts needed to reproduce our experiments for techniques `EF-opt`, `Interpolative`,
482 `varintG8IU`, and `OPT-PFD`; and *ii)* Directories `II_docs` and `EXPERIMENTS` that include respectively the source
483 code of the remaining inverted index variants and the corresponding scripts to run the experiments. We
484 include more details below.

- 485 • Our implementation and scripts within directories `II_docs` and `EXPERIMENTS`. We have organized
486 the source code for our inverted indexes within `II_docs` directory. In Figure 6, we can see a script

487 `compile.sh` and three subdirectories: `ilists.imp`, `lib`, and `src`. The implementation for all our
488 types of compressed posting list representations is contained within directory `ilists.imp`. Script
489 `compile.sh` will create a package for each of them and will move such package into `lib` directory.
490 Finally, the source code for our non-positional compressed inverted index, located within `src` directory,
491 will be linked with each of those posting list representations to obtain the final `BUILD*` and `SEARCH*`
492 binaries for each all our variants of non-positional inverted index.

493 In Figure 6, we can also see the contents of directory `EXPERIMENTS/experiments20gb`. Basically, there
494 is a subdirectory for each technique with scripts `goBuild.sh` (to create the indexes), `goSeach.sh`
495 (to perform query operations), and `clean.sh` (to clean temporal stuff). Those techniques are re-
496 spectively (top-to-bottom in the figure): 1.4) Rice and RiceB, 1.5) Rice-Runs, 2.1) Vbyte and
497 VbyteB, 2)Vbyte-CM and Vbyte-CMB, 3) Vbyte-ST and Vbyte-STB, 4.22) RePair, 4.33) RePair-Skip,
498 4.33.1) RePair-Skip-CM, 4.33.2) RePair-Skip-ST, 555) Vbyte-Lzend, 6) Vbyte-LZMA, 7) Simple9,
499 8) PforDelta, and 9) QMX. In addition, a script `goBuildSearchClean.sh` is in charge of entering each
500 directory running the experiments for the corresponding technique. The output of those experiments
501 are data files containing both space and time statistics (see files `E_Wa.1.5.riceRLE.dat`, ...).

- 502 • Ottaviano&Venturini’s variants within directory `EliasFano.OV14`: This directory contains: *i*) A sub-
503 directory `partitioned_elias_fano` containing the source code from Ottaviano&Venturini’s framework
504 and a script to run the experiments for variants `EF-opt`, `Interpolative`, `varintG8IU`, and `OPT-PFD`;
505 and *ii*) two additional directories `II_docs.EliasFano_formato.ver.leeme` and `2.1.extract.ids.patterns`
506 that contain source code and scripts to transform the text-based patterns into the *id*-based patterns
507 that will be used at query time. In this case, the output of the experiments is written to a log-file and
508 finally a Python script parses such file to gather the values regarding space and time measures for each
509 technique.

510 *Directory `uiHRDC/indexes/POS`: Positional Inverted Indexes.* As indicated above, the structure of this di-
511 rectory is almost identical to that of directory `uiHRDC/indexes/NO-POS`. Therefore, we include no further
512 details here.

513 3.2.2. Self-Indexes

514 Under directory `uiHRDC/self-indexes`, as it is shown in Figure 7, we can find a script `runAllExperiments.sh`
515 and one directory for the different types of self-indexes. In particular, we find directories `RLCSA`, `WCSA`, `LZ`
516 (for `LZ77-index` and `LZend-index`) and `SLP` (for `SLP` and `WSLP`). In those directories we will find, among
517 others, the source code and scripts to run the experiments for each technique.

518 In addition, directory `uiHRDC/self-indexes/collectResults` contains scripts to process the log files
519 obtained when we run the script `runAllExperiments.sh` and create gnuplot-type data files for each technique

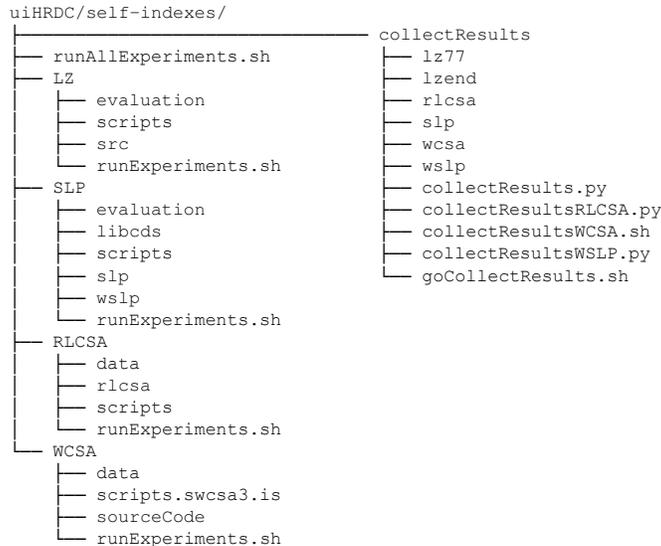


Figure 7: Structure of the `uiHRDC/self-indexes` directory in the `uiHRDC` repository.

520 that are used to make up the final report for the experiments.

521 4. Deploying the Experimental Setup with Docker

522 For those readers interested in reproducing our experiments in their machine, we provide a *Docker*
523 environment that will allow them to:

- 524 1. Reproduce our test framework. We create a docker image with Ubuntu 14 (`ubuntu:trusty`) that in-
525 cludes all the libraries and software requirements to compile/run our indexing alternatives and also
526 build the final report. These includes packages that are installed via `apt` such as: `gcc-multilib`, `g++-`
527 `multilib`, `cmake`, `libboost-all-dev`, `p7zip-full`, `openssh-server`, `screen`, `texlive-latex-base`, and `texlive-`
528 `fonts-recommended`; and finally, `snappy-1.1.1`¹⁵, which we included in a `snappy-1.1.1.tar.gz` file, and
529 `gnuplot-qt`, which in included in `gnuplot-4.4.3.tar` file. In addition, the contents of the `uiHRDC` frame-
530 work (downloaded from our repository) and copied into `/home/user/uiHRDC` directory of our docker
531 image.
- 532 2. Connect to an instance of our docker image using `ssh/sftp`. Basically, we **expose** port 22 and create
533 a user `'user'` with password `'userR1'` who can connect via `ssh`. This will allow the reader to connect
534 to the docker container as to any remote server and to retrieve the final report by `sftp`. Once connected,
535 `user` has `sudo` priviledges and, for example, can become `root` simply entering `sudo su`.

¹⁵<https://github.com/google/snappy>

536 3. Run `doAll.sh` script to automatically run all our experiments and generate our final report. This
537 script must be run by `root` user. Note that we have installed (via `apt-get`) `screen` virtual terminal so
538 that the user can disconnect from the docker container and still keep `doAll.sh` script running.

539 The minimum hardware requirements to run `doAll.sh` script would be to have a machine with at least
540 32GB RAM (and 16GB swap) and around 200GB of free disk space (in the file system where Docker keeps
541 its files¹⁶). In our machine, i7-8700K@3.70GHz CPU (6 cores/12 siblings) with 64 GB of DDR4@2400MHz
542 memory and a 7200rpm SATA disk, it took around 40 hours to run all the experiments from `doAll.sh` script.

543 4.1. Running the experiments: step-by-step

544 Below we show all the commands needed to use Docker¹⁷ to reproduce our framework within a Docker
545 container, then run the experiments within that container, and retrieve the final report (including also the
546 gnuplot-type result data files).

- 547 1. Create a temporal working directory (i.e. `$TMP`) and move to it: `mkdir $TMP` and `cd $TMP`
- 548 2. Clone the `uiHRDC` GitHub repository at `https://github.com/migumar2/uiHRDC/`. It will create a
549 directory `$uiHRDC`. Please rename it as `$DIR`. Then move to directory `$DIR`. Within `$DIR` you will find
550 a file `Dockerfile` and two directories: `docker` and `uiHRDC`.
551 `$TMP> git clone https://github.com/migumar2/uiHRDC.git`
552 `$TMP> mv uiHRDC $DIR`
553 `$TMP> cd $DIR`
- 554 3. Creating a docker image named `repet`: (you must be at `$DIR` directory)
555 `$DIR> sudo docker build -t repet --rm=true . # Note: there is a final 'dot' (.)`
- 556 4. Running a docker container named `repet-exp` that exports `ssh/sftp` port (22) via port 22222:
557 `$DIR> sudo docker run --name repet-exp -p 22222:22 -it repet`
558 Now, you can press `CTRL+P CTRL+Q` to detach from docker-container.
- 559 5. Connecting via `ssh` to the container and becoming `root`.
560 `$DIR> ssh user@localhost -p 22222 # enter password userR1 when prompted.`
561 `$user@docker> sudo su # enter password userR1 when prompted.`
562 `$root@docker> screen -t EXPERIMENTS # optional, to launch screen virtual terminal.`
- 563 6. Move to the working directory `/home/user/uiHRDC` where `doAll.sh` script is located.
564 `$root@docker> cd /home/user/uiHRDC`

¹⁶If you experiment space issues, a simple workaround could be to check which is the directory where Docker creates its images (e.g. `/home/shared/docker/tmp/` or `/var/lib/docker`), and to replace it by a symbolic link pointing to a location in a larger partition/disk.

¹⁷To install Docker, please refer to the installation guide for your host operating system: <https://docs.docker.com/install/>. In our Ubuntu system it simply consisted in running `sudo apt-get install docker.io`

565 7. Now we run `doAll.sh` script. `$root@docker> sh doAll.sh`
566 and wait until `doAll.sh` had completed.

567 8. If you want to retrieve the final report and the results (gnuplot-type data files) via sftp you must grant
568 access to those files to user `user`.

569 `$root@docker>tar czvf /home/user/uiHRDC/report.tar.gz /home/user/uiHRDC/benchmark`
570 `$root@docker>chown user:user /home/user -R`

571 9. Now we can close the ssh session (or detach from the docker container).

572 `$root@docker> exit`
573 `$user@docker> exit`

574 10. Now we can connect by sftp or scp using again port 22222 to download `/home/user/uiHRDC/report.tar.gz`
575 from the docker container.

576 `$DIR> sftp -P 22222 user@localhost:/home/user/uiHRDC # enter password userR1 when prompted.`
577 `$sftp> get report.tar.gz`

578 or alternatively:

579 `$DIR> scp -P 22222 user@localhost:/home/user/uiHRDC/report.tar.gz . # use password userR1`
580 When decompressed, you will find the report here `$DIR/benchmark/report/report.pdf` and all the
581 gnuplot data files within `$DIR/benchmark/report/figures` directory.

582 11. Finally you can stop `repet-exp` container, and if needed remove it and also `repet` image:¹⁸

583 `$DIR> sudo docker stop repet-exp`
584 `$DIR> sudo docker rm repet-exp`
585 `$DIR> sudo docker rmi repet`

586 5. Conclusions and Future Work

587 We have briefly described all the techniques used in our original paper. In total, we had twenty two
588 variants of posting list representations available that were used to create both non-positional and positional
589 inverted indexes. In addition, we had six self-indexing techniques. Since those techniques have their own
590 configuration parameters, and in some cases dependencies of libraries/software, it would be hard to replicate
591 the results and to reuse our techniques by simply reading our original paper and cloning the (GitHub)
592 repository were we had made our sources available. To overcome those limitations, this paper includes a
593 detailed description of our replication framework `uiHRDC`. An interested reader could find not only our source
594 code, our document collections and the query patterns used in our experiments, but also a set of related
595 scripts. Those scripts permit us to replicate all our experiments with little effort and, additionally, generate

¹⁸The reader can use `sudo docker ps` and `sudo docker images` to see respectively the active containers and existing images.

596 a PDF report (using python, gnuplot, and latex) that contains the figures with the experimental results
597 from our parent paper. In addition, we provide some configuration files to create a Docker container that
598 reproduces our test environment (including all the required dependencies), and instructions regarding how
599 to start the container, run the experiments, and finally download a copy of the final report from the Docker
600 container. We hope that the descriptions and instructions provided along this paper will simplify the work
601 of any reader interested in reusing the experiments from our original paper.

602 An interesting line of future work is to redesign uiHRDC to facilitate that new indexes to be added to
603 the framework. The resulting benchmark framework would ensure reproducible experimentation for ongoing
604 research about compression and indexing of highly repetitive collections.

605 **6. Revision Comments**

606 We would like to thank the authors for providing this valuable reproducibility platform, which allows
607 both reproducing the results of its parent paper and encouraging the evaluation of new indexes for highly
608 repetitive document collections in an easy and systematic way. Undoubtedly, this will be a topic of active
609 research in the coming years given, for example, the cheapening of gene sequencing technology. Hence,
610 the public availability of this type of benchmarking platform is becoming imperative. For this reason, we
611 sincerely expect that uiHRDC sets an experimental standard for this line of research.

612 On the other hand, this work confirms again that the production of reproducible science is a difficult
613 task, thus reproducibility of any research work must be considered and planned since the very beginning
614 stages of development. Despite the reviewers reached a consensus about the reproducibility of the paper,
615 the rigorous review process raised some reproducibility issues that left us some significant lessons on the
616 difficulties of producing fully reproducible science and developing such ambitious reproducible platform as
617 introduced herein. Next, we discuss our experience with this work as well as the lessons learned from our
618 review.

619 We fully support the choice of using Docker for this kind of reproducibility experiments. Using a container
620 makes it extremely simple to reproduce the experimentation framework, which otherwise would require to
621 install the dependencies and compile the tested codes manually. However, review process unveiled other issues
622 on the set-up and running of the uiHRDC experiments, which were fixed by providing specific instructions
623 solving them. For instance, authors added information in the output report to warn the experimenter in
624 those cases in which the experiments failed by lack of resources, as well as to ensure that the folder used by
625 Docker resides in a drive with enough memory, which might not be the case for the default drive used by
626 Docker.

627 In addition to the aforementioned set-up and running issues, some others arose with the experimentation
628 itself in a first review: (1) a significant mismatch in the ordering of time-axis results derived from the

629 uncertainty in the evaluation of running time values; (2) a difference in the compression ratios obtained
630 for the OPT-PFD indexing method; (3) missing values for Vbyte-Lzend and QMX-SIMD methods derived
631 from software problems which did not provide any warning to the experimenter; and (4) differences of scaling
632 between the original figures and those generated by uiHRDC, which made a detailed checking process difficult.
633 Subsequently, all experimentation issues detailed above were either fixed or honestly and rigorously justified
634 in the paper in a sample of best research practices. First, time measurement was significantly improved
635 by increasing the number of evaluations with the aim of reducing its fluctuation. It worths to highlight
636 that special attention should be put in any research for the reproducibility of running time values and their
637 conclusions. Second, a few code bugs were fixed and the output experimentation report was extended to
638 provide detailed information on the execution of the experiments. This further information shown to be very
639 useful to guide the review process, and also to detect possible bugs for particular execution environments.
640 Finally, axis scaling was revised to match the original paper exactly, and a bug with gnuplot was fixed so the
641 generated report is produced with the same symbology of the parent paper. These later two issues highlight
642 the importance of the presentation of the results to properly communicate research findings in a clear way.

643 The uiHRDC framework is a first try of standardizing experiments that, in absence of such a framework,
644 would need to be repeated each time a new index is introduced in the literature. Thus, uiHRDC will be
645 a very valuable resource to the research community by avoiding this tedious and costly task, which is also
646 prone to errors.

647 As forthcoming activities, we invite the authors to extend uiHRDC with the aim of removing some
648 drawbacks that hinder the integration of new indexes in their platform. The main drawback is that adding
649 a new index is not trivial and requires editing several scripts; thus, any improvement in this sense will
650 contribute to turn uiHRDC into a standard for comparison by future researchers.

651 **Acknowledgments**

652 This paper is funded in part by European Union’s Horizon 2020 research and innovation programme under
653 the Marie Skłodowska-Curie grant agreement No 690941 (project BIRDS). Antonio Fariña is funded by Xunta
654 de Galicia/FEDER-UE [CSI: ED431G/01 and GRC: ED431C 2017/58] and by MINECO-AEI/FEDER-UE
655 [ETOME-RDFD3: TIN2015-69951-R]. Miguel A. Martínez-Prieto is funded by MINECO-AEI/FEDER-UE
656 [Datos 4.0: TIN2016-78011-C4-1-R]. Gonzalo Navarro is funded by the Millennium Institute for Foundational
657 Research on Data (IMFD), and by Fondecyt Grant 1-170048, Conicyt, Chile.

658 **References**

- 659 [1] ACM. Artifact Review and Badging. [https://www.acm.org/publications/policies/artifact-review-](https://www.acm.org/publications/policies/artifact-review-badging)
660 [badging](https://www.acm.org/publications/policies/artifact-review-badging), 2018.

- 661 [2] V. Anh and A. Moffat. Inverted index compression using word-aligned binary codes. *Information*
662 *Retrieval*, 8:151–166, 2005.
- 663 [3] R. A. Baeza-Yates and B. Ribeiro-Neto. *Modern Information Retrieval*. Addison-Wesley Publishing
664 Company, 2nd edition, 2011.
- 665 [4] F. Chirigati, R. Capone, R. Rampin, J. Freire, and D. E. Shasha. A Collaborative Approach to Com-
666 putational Reproducibility. *Information Systems*, 59:95–97, 2016.
- 667 [5] F. Claude, A. Fariña, M. A. Martínez-Prieto, and G. Navarro. Universal Indexes for Highly Repetitive
668 Document Collections. *Information Systems*, 61:1–23, 2016.
- 669 [6] F. Claude, A. Fariña, M. Martínez-Prieto, and G. Navarro. Compressed q -gram indexing for highly
670 repetitive biological sequences. In *Proc. 10th International Conference on Bioinformatics and Bioengi-*
671 *neering (BIBE)*, pages 86–91, 2010.
- 672 [7] F. Claude, A. Fariña, M. Martínez-Prieto, and G. Navarro. Indexes for highly repetitive document
673 collections. In *Proc. 20th ACM International Conference on Information and Knowledge Management*
674 *(CIKM)*, pages 463–468, 2011.
- 675 [8] C. Collberg, T. Proebsting, and A. M. Warren. Repeatability and Benefaction in Computer Systems
676 Research: A Study and a Modest Proposal. Technical Report TR 14-04, University of Arizona, 2015.
- 677 [9] J. S. Culpepper and A. Moffat. Efficient set intersection for inverted indexing. *ACM Transactions on*
678 *Information Systems*, 29(1):article 1, 2010.
- 679 [10] A. Fariña, M. A. Martínez-Prieto, F. Claude, and G. Navarro. *uiHRDC (Mendeley Data v1)*. 2018.
680 <http://dx.doi.org/10.17632/xxntkjvtxw.1>.
- 681 [11] A. Fariña, N. Brisaboa, G. Navarro, F. Claude, A. Places, and E. Rodríguez. Word-based self-indexes
682 for natural language text. *ACM Transactions on Information Systems*, 30(1):article 1, 2012.
- 683 [12] N. Ferro, N. Kando, N. Fuhr, M. Lippold, K. Jrvelin, and J. Zobel. Increasing Reproducibility in IR:
684 Findings from the Dagstuhl Seminar on ”Reproducibility of Data-Oriented Experiments in e-Science”.
685 *ACM SIGIR Forum*, 50:68–82, 2016.
- 686 [13] N. Ferro and G. Silvello. Rank-Biased Precision Reloaded: Reproducibility and Generalization. In
687 *Proceedings of the 37th European Conference on IR Research, ECIR*, page 768780, 2015.
- 688 [14] J. He, J. Zeng, and T. Suel. Improved index compression techniques for versioned document collections.
689 In *Proc. 19th ACM International Conference on Information and Knowledge Management (CIKM)*,
690 pages 1239–1248, 2010.

- 691 [15] S. Heman. *Super-scalar database compression between RAM and CPU-cache*. PhD thesis, Centrum voor
692 Wiskunde en Informatica (CWI), Amsterdam, 2005.
- 693 [16] S. Kreft and G. Navarro. On Compressing and Indexing Repetitive Sequences. *Theoretical Computer
694 Science*, 483:115–133, 2013.
- 695 [17] N. Larsson and A. Moffat. Offline Dictionary-Based Compression. pages 296–305. IEEE Computer
696 Society, 1999.
- 697 [18] J. J. Lastra-Díaz, A. García-Serrano, M. Batet, M. Fernández, and F. Chirigati. HESML: A scal-
698 able ontology-based semantic similarity measures library with a set of reproducible experiments and a
699 replication dataset. *Information Systems*, 66:97 – 118, 2017.
- 700 [19] D. Lemire, L. Boytsov, and N. Kurz. Simd compression and the intersection of sorted integers. *Software:
701 Practice and Experience (published online)*, 2015.
- 702 [20] J. Lin, M. Crane, A. Trotman, J. Callan, I. Chattopadhyaya, J. Foley, G. Ingersoll, C. Macdonald,
703 and S. Vigna. Toward Reproducible Baselines: The Open-Source IR Reproducibility Challenge. In
704 *Proceedings of the 38th European Conference on IR Research, ECIR*, pages 408–420, 2016.
- 705 [21] V. Mäkinen, G. Navarro, J. Sirén, and N. Välimäki. Storage and Retrieval of Highly Repetitive Sequence
706 Collections. *Journal of Computational Biology*, 17(3):281–308, 2010.
- 707 [22] A. Moffat and L. Stuiiver. Binary interpolative coding for effective index compression. *Information
708 Retrieval*, 3(1):25–47, 2000.
- 709 [23] G. Navarro. *Compact Data Structures – A practical approach*. Cambridge University Press, 2016.
- 710 [24] G. Ottaviano and R. Venturini. Partitioned elias-fano indexes. In *Proc. 37th International ACM SIGIR
711 Conference on Research and Development in Information Retrieval (SIGIR)*, pages 273–282, 2014.
- 712 [25] K. Sadakane. New Text Indexing Functionalities of the Compressed Suffix Arrays. *Journal of Algorithms*,
713 48(2):294–313, 2003.
- 714 [26] G. K. Sandve, A. Nekrutenko, J. Taylor, and E. Hovig. Ten Simple Rules for Reproducible Computa-
715 tional Research. *PLOS Computational Biology*, 9:1–4, 2013.
- 716 [27] A. A. Stepanov, A. R. Gangolli, D. E. Rose, R. J. Ernst, and P. S. Oberoi. SIMD-based decoding of
717 posting lists. In *Proc. 20th ACM International Conference on Information and Knowledge Management
718 (CIKM)*, pages 317–326, 2011.

- 719 [28] F. Transier and P. Sanders. Engineering basic algorithms of an in-memory text search engine. *ACM*
720 *Transactions on Information Systems*, 29:article 2, 2010.
- 721 [29] A. Trotman. Compression, simd, and postings lists. In *Proc. 19th Australasian Document Computing*
722 *Symposium (ADCS)*, pages 50–57. ACM, 2014.
- 723 [30] P. Vandewalle, J. Kovacevic, and M. Vetterli. Reproducible Research in Signal Processing - What, Why,
724 and How. *IEEE Signal Processing Magazine*, 26:37–47, 2009.
- 725 [31] H. Williams and J. Zobel. Compressing integers for fast file access. *The Computer Journal*, 42:193–201,
726 1999.
- 727 [32] I. Witten, A. Moffat, and T. Bell. *Managing Gigabytes*. Morgan Kaufmann, 2nd edition, 1999.
- 728 [33] A. Wolke, M. Bichler, F. Chirigati, and V. Steeves. Reproducible experiments on dynamic resource
729 allocation in cloud data centers. *Information Systems*, 59:98–101, 2016.
- 730 [34] H. Yan, S. Ding, and T. Suel. Inverted index compression and query processing with optimized document
731 ordering. In *Proc. 18th International Conference on World Wide Web (WWW)*, pages 401–410, 2009.
- 732 [35] M. Zukowski, S. Heman, N. Nes, , and P. Boncz. Super-scalar RAM-CPU cache compression. In *Proc.*
733 *22nd International Conference on Data Engineering (ICDE)*, page 59, 2006.

734 A. Compression Ratios

735 This appendix shows the exact compression ratios for all the techniques included in our uiHRDC framework.
736 The corresponding tuning parameters are provided for each case (\times is used for those techniques that do not
737 require any parameter). Those values are included in Table 3.

738 Note that in the positional scenario, when the source text collection is compressed with Re-Pair, the parent
739 paper allowed the extraction of snippets consisting of either 80 or 13,000 chars. We considered sampling
740 parameter values $sample_ct = \{1, 2, 8, 32, 64, 256, 4096\}$. Table 4 shows the compression ratios of the RePair-
741 compressed text that corresponds to each sampling configuration. Yet, the plot included in the parent paper
742 (Fig.10-left) when 80 chars included only points for $sample_ct = \{1, 2, 8, 32, 64\}$, and when extracting 13,000
743 chars (see Figure 3, or Fig.10-right in the parent paper) we tuned $sample_ct = \{1, 8, 32, 256, 4096\}$.

744 B. Using Self-Indexes for other types (non-document oriented) of Highly-Repetitive Collec- 745 tions

746 While self-indexes like WCSA or WSLP are designed to operate on words, the remaining approaches in
747 our benchmark can effectively operate on an universal scenario; i.e. they are able to exploit repetitiveness

| Non-positional indexes | | | Positional indexes | | |
|------------------------|-------------------|--|--------------------|-------------------|--|
| Method | Compression ratio | | Method | Compression ratio | |
| | Value | Parameterization | | Value | Parameterization |
| Vbyte | 4.4592% | × | Vbyte | 35.1530% | × |
| VbyteB | 3.3522% | $lenBitmapDiv = 8$ | Vbyte-CM | 35.3805% | $k = 32$ |
| Vbyte-CM | 4.4956% | $k = 32$ | | 36.5899% | $k = 4$ |
| | 4.7095% | $k = 4$ | Vbyte-ST | 35.4935% | $B = 128$ |
| Vbyte-CMB | 3.3754% | $k = 32, lenBitmapDiv = 8$ | | 37.4281% | $B = 16$ |
| | 3.4996% | $k = 4, lenBitmapDiv = 8$ | Rice | 36.7974% | × |
| Vbyte-ST | 4.5147% | $B = 128$ | Simple9 | 36.6233% | × |
| | 4.8650% | $B = 16$ | QMX | 136.8762% | × |
| Vbyte-STB | 3.3820% | $B = 128, lenBitmapDiv = 8$ | Vbyte-LZMA | 9.7539% | $minbcssize = 10$ |
| | 3.5488% | $B = 16, lenBitmapDiv = 8$ | RePair | 20.0641% | × |
| Rice | 3.0807% | × | RePair-Skip | 21.3769% | $repairBreak = 5 \times 10^{-7}$ |
| RiceB | 3.2235% | $lenBitmapDiv = 8$ | RePair-Skip-CM | 20.0786% | $k = 64, repairBreak = 5 \times 10^{-7}$ |
| Simple9 | 0.9130% | × | | 20.8584% | $k = 1, repairBreak = 5 \times 10^{-7}$ |
| PforDelta | 0.9372% | $pdfThreshold = 100$ | RePair-Skip-ST | 21.7456% | $B = 256, repairBreak = 5 \times 10^{-7}$ |
| QMX | 4.8825% | × | EF-opt | 30.5630% | × |
| Rice-Runs | 0.3247% | × | OPT-PFD | 29.7424% | × |
| Vbyte-LZMA | 0.2030% | $minbcssize = 10$ | Interpolative | 29.8820% | × |
| | 0.1420% | $ds = 256$ | varintG8IU | 37.9975% | × |
| | 0.1437% | $ds = 64$ | RLCSA | 2.4735% | $sample = 2048$ |
| | 0.1508% | $ds = 16$ | | 2.8686% | $sample = 1024$ |
| | 0.1790% | $ds = 4$ | | 3.6630% | $sample = 512$ |
| Re-Pair | 0.1040% | × | | 5.2489% | $sample = 256$ |
| RePair-Skip | 0.1097% | $repairBreak = 4 \times 10^{-7}$ | | 8.4048% | $sample = 128$ |
| RePair-Skip-CM | 0.1195% | $k = 64, repairBreak = 4 \times 10^{-7}$ | | 14.6924% | $sample = 64$ |
| | 0.1212% | $k = 1, repairBreak = 4 \times 10^{-7}$ | | 27.2996% | $sample = 32$ |
| RePair-Skip-ST | 0.1279% | $B = 1024, repairBreak = 4 \times 10^{-7}$ | WCSA | 8.8921% | $\langle sA, sAinv, sPsi \rangle = \langle 2048, 2048, 2048 \rangle$ |
| EF-opt | 0.4984% | × | | 9.4071% | $\langle sA, sAinv, sPsi \rangle = \langle 512, 512, 512 \rangle$ |
| OPT-PFD | 0.7015% | × | | 11.0718% | $\langle sA, sAinv, sPsi \rangle = \langle 128, 256, 128 \rangle$ |
| Interpolative | 0.5573% | × | | 15.8673% | $\langle sA, sAinv, sPsi \rangle = \langle 64, 128, 32 \rangle$ |
| varintG8IU | 5.3144% | × | | 17.9458% | $\langle sA, sAinv, sPsi \rangle = \langle 32, 64, 32 \rangle$ |
| | | | | 25.6777% | $\langle sA, sAinv, sPsi \rangle = \langle 16, 64, 16 \rangle$ |
| | | | | 50.8565% | $\langle sA, sAinv, sPsi \rangle = \langle 8, 8, 8 \rangle$ |
| | | | SLP | 3.2374% | × |
| | | | WSLP | 2.8962% | × |
| | | | LZ77-index | 1.8357% | × |
| | | | LZend-index | 2.3894% | × |

Table 3: Summary of compression ratios for all indexes in the framework.

| Compression ratio | Parameterization |
|-------------------|------------------|
| 1.306% | sample_ct = 1 |
| 1.265% | sample_ct = 2 |
| 1.227% | sample_ct = 8 |
| 1.215% | sample_ct = 32 |
| 1.213% | sample_ct = 64 |
| 1.211% | sample_ct = 256 |
| 1.210% | sample_ct = 4096 |

Table 4: Summary of compression ratios for the RePair-compressed version of the text with different sampling values.

748 underlying to any arbitrary sequence of chars.

749 Focusing on SLP, LZ77-index, and LZ-End, the implementations provided in uiHRDC can be easily used
750 to self-index any data collection. Scripts for indexing and searching can be reused, although the number
751 of parameters passed to `locate` changes. Currently, `locate` requires a parameter (called *doc_boundaries*)
752 which sets the path to a file that contains the array of offsets that mark the beginning of each document in
753 the text. It is required to map absolute positions to *(doc, offset)* pairs which indicate the position of each
754 pattern occurrence within a document. However, this operation is not usually needed in a general scenario,
755 so the corresponding parameter is neither needed.

756 Therefore, to invoke `locate` in a general scenario we have simply to remove *doc_boundaries*, and the
757 script will deliver absolute positions of the pattern in the text.